

RとSTANで学ぶフリーで楽しい心理統計の世界

心理学データ解析 応用

.....
データの背後のメカニズムを解析する方法



小杉考司

この本は Creative Commons BY-SA(CC BY-SA) ライセンス Version 4.0 に基づいて提供されています。著者に適切なクレジットを与える限り、この本を再利用、再編集、保持、改訂、再頒布 (商用利用を含む) をすることができます。もし再編集したり、このオープンなテキストを変更したい場合、すべてのバージョンにわたってこれと同じライセンス、CC BY-SA を適用しなければなりません。

<https://creativecommons.org/licenses/by-sa/4.0/deed.ja>

This book is published under a Creative Commons BY-SA license (CC BY-SA) version 4.0.

This means that this book can be reused, remixed, retained, revised and redistributed (including commercially) as long as appropriate credit is given to the authors. If you remix, or modify the original version of this open textbook, you must redistribute all versions of this open textbook under the same license - CC BY-SA.

<https://creativecommons.org/licenses/by-sa/4.0/>

心理学データ解析応用

小杉 考司

Last Compiled on 2024.1.11

はじめに

昨今はデータサイエンス、情報科学の領域が非常に隆盛で、コンピュータを使ってデータを分析し、経済の動向や購買行動などの予測に用いられることが広く行われている。

人の行動や考え方をどのようにデータにするかについては、当然ながら心理学には一日の長がある。また、人が頭の中でどのような考え方のプロセスをたどるのか、それをどのように検証するのかについても、心理学はその短い歴史の中で徹底的にその技法を洗練させてきた。このような根源的なレベルでの理論や方法論は時代が変わっても色褪せることなく、また今後ますます必要とされてくる時代になっている。

本講ではデータ解析の応用段階として、より実践的なテーマを扱う。すなわち、**心理尺度が作られる理論的背景と、データの背後のメカニズムを解析する方法**を知ることである。

心理学研究法の1つとして、調査研究がある。紙とペンで回答を集めた時、回答者がある反応カテゴリにまるをつけたことが、どうして数値処理の対象になるのか。そこには数字を割り振るルールとしての「尺度化」の手続きがある。残念ながら応用的側面が発展しすぎたため、回答に数字を割り振る原理について語られることが少なくなってきてしまい、それに対する反動からか、近年改めてこの根本原理についての理解と解説が求められている。この講義では、尺度化の原理や目に見えない潜在変数を想定して分析するとはどういうことかについて、理論と演習を交えながら習得することを目指す。この理論的側面を考えるためには、どうしても線形代数・行列計算の知識が必要になってくる。線形代数については特別な事前知識は不要で、定義から改めて解説するので安心してもらいたい。

心理学研究法のまた1つの大きな柱として、実験的研究がある。実験的研究はその手続きが厳格に準備されていることで、結果は群間の平均値差を推定すれば十分である、とされてきた。心理学の基礎領域では、そのため、標本の平均値から母集団の平均値を推測する方法が主なテーマになる。しかしこの方法は逆に、結果を群間の平均値差に帰着させるための実験計画を必要とする。結果に方法が規定されているのである。本来考えたかったことは群間の差だけではなく、どのようにして心理的なメカニズムからデータが生成されているか、という問いであったことを思い出そう。そして群間の平均値に拘ることなく、心理的なメカニズムを数式的に表現することで分析する**数理モデリング**というアプローチがある。このモデリングの基礎的な知識、方法論の習得を目指す。

授業のテーマ

データから意味のある情報を取り出すための、さまざまな分析法を習熟するにあたって、その背後にあって語られることのない「発想」の観点から理解する。数値だけに振り回される状態から脱却し、数値を算出する数式に込められた意味について考える視点を持つ。またこれらに習熟することで、どのような研究対象に対してどのような心理統計的アプローチができるかを、俯瞰的に見れるようになる。

一年を通じて伝えたいポイント

尺度化とは何か 心理学で行われるアンケート調査やその後の分析はどのような原理があって「心を測定した」といえるのか。その原理やモデルを理解して利用できるようになる。

多変量解析から何がわかるのか 調査研究などで得られた多変量を分析することで何がわかるのか。あるいは何をしてわかったというのか。

多変量解析の基礎となる数式的原理 多変量解析の背景にあるのは線形代数という数学であり、線形代数の基礎を学ぶことで多変量解析のメカニズムを統合的に理解できる。

データ生成メカニズム データから情報を取り出す受け身の分析ではなく、データに数字を与えたり、データが生まれてくるメカニズムをリバースエンジニアリングすることで、さらに積極的にデータ解析に立ち向かおう。

統計環境 R と確率的プログラミング言語 Stan による実践 統計環境 R と確率的プログラミング言語 Stan に習熟することで実際に計算し、確認しながら分析を進めることができる。

その他

授業シラバスとこの講義資料を掲載したサイト (https://kosugitti.github.io/psychometrics_syllabus/) で、最新版のシラバスと授業資料、授業で用いるサンプルデータやコードの配布を行なっています。

目次

第 I 部	心理学データ解析応用 1	13
第 1 章	導入；多変量データと心理学	15
1.1	正規線形モデルの世界	16
1.2	尺度の四水準	18
1.3	平均と分散	20
1.4	共分散と相関係数	22
1.5	課題	23
第 2 章	心理尺度を作る	25
2.1	はじめに	25
2.2	サーストンの等現間隔法	27
2.3	リッカートのシグマ法	28
2.4	尺度を評価する	30
2.5	課題	32
第 3 章	テスト理論と因子分析	35
3.1	古典的テスト理論	35
3.2	因子分析モデル	36
3.3	因子分析の定理	38
3.4	課題	42
第 4 章	現代テスト理論	43
4.1	因子分析とテスト理論	43
4.2	通過率と累積正規分布	43
4.3	項目母数の特徴	46
4.4	被験者母数の推定	48
4.5	課題	50
第 5 章	現代テスト理論その 2	51
5.1	現代テスト理論の特徴	51
5.2	段階反応モデル	55
5.3	因子分析の歴史と展開	57
5.4	課題	59

第 6 章	行列計算の基礎	61
6.1	行列とベクトル	61
6.2	行列の四則演算と操作	63
6.3	行列を使うと便利なこと	67
6.4	課題	69
第 7 章	行列による関係の表現	71
7.1	データの行列表現	71
7.2	線形モデルの行列表現	73
7.3	デザイン行列	74
7.4	因子分析モデルの行列表現	77
7.5	課題	78
第 8 章	固有値と固有ベクトルと因子分析モデルの関係	79
8.1	固有値と固有ベクトル	79
8.2	固有値と固有ベクトルを求める	81
8.3	固有値と固有ベクトルの幾何学的意味	83
8.4	因子分析の数学的理解	84
8.5	課題	85
第 9 章	R をつかっての行列計算	87
9.1	R による行列計算	87
9.2	データの行列表現	92
9.3	R による固有値計算	94
9.4	課題	95
第 10 章	R をつけた因子分析と尺度作成法	97
10.1	調査研究の手順	97
10.2	共通性の推定	98
10.3	因子数の決定	99
10.4	探索的因子分析の実際	100
10.5	因子分析の後で	107
10.6	さいごに	108
10.7	課題	109
第 11 章	R をつけた項目反応理論	111
11.1	項目反応理論の実際	111
11.2	段階反応モデルの実際	119
11.3	カテゴリカル因子分析との対応	121
11.4	課題	122
第 12 章	構造方程式モデリング	123
12.1	パスダイアグラムの書き方	123
12.2	パスダイアグラムによるさまざまなモデル	125

12.3	構造方程式モデルによる未知数の推定	127
12.4	適合度によるモデルの評価	130
12.5	課題	131
第 13 章	R による構造方程式モデリング	133
13.1	モデル式の入力	133
13.2	実践上の注意点	140
13.3	そのほかの統計パッケージ	141
13.4	課題	141
第 14 章	双対尺度法	143
14.1	直線的ではない関係	143
14.2	林の数量化理論	145
14.3	双対尺度法による分析	146
14.4	テキストマイニングへの応用	148
14.5	課題	149
第 15 章	多次元尺度構成法	151
15.1	多次元尺度構成法	151
15.2	距離と心理学のデータ	153
15.3	非計量多次元尺度法	154
15.4	多次元尺度法の展開	158
15.5	課題	162
第 II 部	心理学データ解析応用 2	163
第 16 章	プログラミングの基礎	165
16.1	プログラミングの基礎	165
16.2	プログラミング言語	168
16.3	プログラミング言語の基本的な働き	169
16.4	まとめ	174
16.5	課題	174
第 17 章	データ生成モデリング	175
17.1	データ生成モデリング	175
17.2	ベイズ推定の基礎	176
17.3	マルコフ連鎖モンテカルロ法	178
17.4	乱数によるアプローチの例	180
17.5	課題	184
第 18 章	いんたーみっしょん；Stan の概略と環境の準備について	185
18.1	はじめに	185
18.2	Stan の位置付け	186

18.3	導入の概略	191
18.4	導入方法 3;外部サーバの利用	192
18.5	Stan を使ってみよう	193
第 19 章	ベイジアンアプローチと確率的プログラミング 1	203
19.1	7 人の科学者	203
19.2	Stan コードの書き方	205
19.3	Stan を使った MCMC の実践	208
19.4	MCMC 結果の診断	210
19.5	MCMC の結果の解釈	213
19.6	課題	214
第 20 章	モデリングの目から見た検定 1 ; 二群の平均値の差	215
20.1	t 検定の過程と実際	215
20.2	差の分布	220
20.3	帰無仮説検定を省みる	223
20.4	今回のまとめ	225
20.5	課題	225
第 21 章	モデリングの目から見た検定 2 ; パラメータの世界とデータの世界	227
21.1	事後予測分布	227
21.2	データレベルの仮説	230
21.3	パラメータ・リカバリ	233
21.4	今回のまとめ	237
21.5	課題	237
第 22 章	モデリングの目から見た検定 3 ; 多群の平均値差モデル	239
22.1	要因計画モデル	239
22.2	パラメータの変形と制約	241
22.3	モデルの洗練	245
22.4	パラメータリカバリ	249
22.5	課題	250
第 23 章	モデリングの目から見た検定 4 ; 対応のある群の比較	251
23.1	対応のある群	251
23.2	ID をもったデータ構造	257
23.3	個人差と変化量のモデルへ	260
23.4	課題	262
第 24 章	モデリングの目から見た検定 5 ; カテゴリカル分布をつかって	263
24.1	離散的な分布	263
24.2	χ^2 検定	264
24.3	カテゴリカル分布のモデリング	266
24.4	κ 係数の算出	268

24.5	課題	271
第 25 章	一般化線形モデル	273
25.1	一般線形モデル	273
25.2	データに合わせた確率分布	279
25.3	リンク関数とパラメータの解釈	284
25.4	まとめ	286
25.5	課題	287
第 26 章	階層線形モデル	289
26.1	一般化線形混合モデル	289
26.2	ネストされたデータ	295
26.3	階層線形モデル	297
26.4	課題	301
第 27 章	混合分布モデル	303
27.1	混合分布モデル	303
27.2	ターゲット記法と周辺化消去	307
27.3	ゼロ過剰ポアソン分布モデル	313
27.4	課題	318
第 28 章	確率的プログラミング；項目反応理論	319
28.1	ロジスティックモデルの復習	319
28.2	ロジスティックモデルでの実装	321
28.3	整然データでの分析	323
28.4	課題	328
第 29 章	確率的プログラミング；変化点と折線回帰	329
29.1	混合分布モデルの応用	330
29.2	変化点検出	332
29.3	折線回帰	334
29.4	課題	339
第 30 章	確率的プログラミング；状態空間モデル	341
30.1	時系列データの特徴	341
30.2	状態空間モデル	342
30.3	欠損値の補間	346
30.4	状態空間モデルの展開	356
30.5	課題	356
第 31 章	モデル比較	357
31.1	ベイジアンモデリング	357
31.2	帰無仮説検定の代案	359
31.3	モデル比較	362

31.4	おわりに	363
付録 A	よくある質問とミス	365
A.1	Frequently Miss and Comments	365
A.2	Frequently Asked Questions;よくある質問と答え	369
付録 B	標準正規分布から尺度値を求める計算方法	377
付録 C	電子計算機のイロハ	381
C.1	前置き	381
C.2	コンピュータの基礎	381
C.3	コンピュータの歴史	382
C.4	情報の単位	385
C.5	ファイルの種類と拡張子	386
C.6	クラウドとは	388
C.7	ファイルの位置の指定	389
C.8	ファイルのバージョン管理	391
C.9	おわりに	392
付録 D	ギリシア文字一覧	393
付録 E	記号の入力とキーボードの場所	395
付録 F	本講義に対応する詳細シラバス	399
F.1	イントロダクション	399
F.2	心理尺度を作る	400
F.3	テスト理論と因子分析	402
F.4	現代テスト理論	403
F.5	現代テスト理論その 2	404
F.6	行列計算の基礎	406
F.7	行列による関係の表現	407
F.8	固有値と固有ベクトルと因子分析モデルの関係	408
F.9	R をつかっての行列計算	409
F.10	R をつかった因子分析と尺度作成法	411
F.11	R をつかった項目反応理論	412
F.12	構造方程式モデリング	414
F.13	R による構造方程式モデリング	415
F.14	双対尺度法	417
F.15	多次元尺度構成法	419
F.16	プログラミングの基礎	420
F.17	データ生成メカニズムとモデリング	422
F.18	ベイジアンアプローチと確率的プログラミング 1	423
F.19	モデリングの目から見た検定 1;二群の平均値の差	425
F.20	モデリングの目から見た検定 2;パタメータの世界とデータの世界	426

F.21	モデリングの目から見た検定 3;多群の平均値差を求めるモデル	427
F.22	モデリングの目から見た検定 4;対応のある群の比較	429
F.23	モデリングの目から見た検定 5;カテゴリカル分布をつかって	430
F.24	一般化線形モデル	431
F.25	階層線形モデル	432
F.26	混合分布モデル	433
F.27	確率的プログラミングの応用 1; 項目反応理論	434
F.28	確率的プログラミングの応用 2; 変化点と折線回帰	435
F.29	確率的プログラミングの応用 3; 状態空間モデル	437
F.30	モデル比較	438
引用文献		439
索引		444

第 I 部

心理学データ解析応用 1

第 1 章

導入；多変量データと心理学

この授業は基礎的な心理統計を修めた人向けの、応用コースになっています。基礎的な心理統計、という言葉に私が込めた意味は、

- 記述統計；得られたデータの統計量を算出したり、可視化することによってデータの特徴を把握する。
- 推測統計；得られたデータが母集団からの標本であると考え、標本の特徴を使って母数を推測する。さらに標本の特徴から母集団について何らかの判断（意思決定）を行う。

という 2 点です。とくに推測統計学の領域では、確率の話やさまざまな推定法、それに伴う技術などが必要ですから、これだけでも膨大な量だったのではと思います。こうした基礎的な知識や技術を身につけると、とりあえず手元のデータを使って差があるかないとか、どの程度効果があったのか、と言った基本的な判断はできると思います。複雑な計算式のところは機械（統計環境 R など）がやってくれますので、出た結果だけをみて判断すれば良いのです*1。

しかし基礎的なところだけで満足していると、「はて、何がしたかったのかな」とそもそもの問題意識を忘れてしまうことにもなりかねません。とりあえず教わった（膨大な！）プロセスを経て実験結果を見てみれば OK なのでしょう、というだけでは表面的な理解に止まっていると言わざるを得ません。さまざまなシーンに適用される方法論、その計算は何を意味していて、元々の数式にはどのような意味があり、式から得られるものは何をどこまで指し示しているのか、というところまで考えるのが、この講義の狙いです。数式は数式に過ぎない、というのはその通りなのですが、その数式にどのような意味があるのか、数式から得られる結果にどのような意味があるのか、を理解した上で数値（データ）を扱うようになることが目的です。数値だけに振り回される状態からの脱却、が狙いなわけです。

この授業の前半（前期）は、講義の形をとります。テーマとしては、**因子分析 (Factor Analysis)** を扱います。因子分析は、いわゆる質問紙調査を行った後で適用される多変量解析法の一つで、たくさんの質問項目の背後にある「因子」を見つけ出します。その因子には XXX 特性、XXX 傾向といった名前がつけられ、その上で心理学的に考察することが一般的です。因子分析の結果として性格特性が得られる、などといったりするわけですから、これはもう心理学の王道中の王道、と言えるかもしれません。しかし実際は統計ソフトウェアが計算してくれて、何だかわからないまま使っている、という人も少なくありません。質問紙調査で分析して何がわかる、というのはどういうことなのか、原理的なところから解説を始めていきます。またこの講義ではその基礎的なところ、数式レベルでしっかりと把握した上で、数値例に進みます。講義が終わる頃には、意味内容をしっかりと理解したうえで因子分析を使えるようになっていることが期待されます。

この授業の後半（後期）は、実践・演習を主にした学習になります。我々が手にしたデータは、何らかのモデ

*1 もちろん結果の見方が間違っていたり、拙速だったりしてはいけないなど、注意すべきところは多々あります。

33 ル・仮定のもとで生成されたものである, という考え方に立ち, データから生成メカニズムを推測することに
 34 チャレンジします。これは非常に夢のある話です。だってデータ生成メカニズムというのは, 私たちの心の中の
 35 機序そのものだからです。推測に当たってはさまざまな前提・仮定が必要ですが, 得られる結果は非常に含
 36 蓄に富み, さまざまな角度から心を考えさせてくれるヒントになります。線形モデルは直線的な関係でしたが,
 37 より柔軟で豊富な表現力を持つ技術へと理解を進めていきましょう。

38 今回は初回ですので, 基礎的な心理統計で学んだことを改めて確認・復習することを中心に, 今後の講義
 39 でも用いられる数学的記号の準備をします。

40 1.1 正規線形モデルの世界

41 因子分析法や (重) 回帰分析では, 基本的に扱う変数が複数あります。これまでの相関・回帰分析や, 実験
 42 計画の中では, 説明変数と被説明変数がひとつずつ, という二変数の世界でした。これが多くなったシーンは,
 43 一般に**多変量解析 (Multivariate Analysis)**と呼ばれます。変数あるいは**変数 (Variables)**は,
 44 ケースごとに変わる数のことを指します。変わる「数」と言っていますが, この後述するように数字以外のもの
 45 も数字として扱いますので, 「ケースごとに変わるもの」の総称だと思ってください。多変量はそれが多くある
 46 もの, データセット全体を指します。**たくさんのデータセットの中から意味ある情報を引き出すこと**, これが
 47 多変量解析の目的です。

48 イメージとしては, 表計算ソフトのスプレッドシートに数字がずらっと並んでいる世界です。たとえば性格
 49 検査の尺度の一種である, YG 性格検査は被験者に 120 の項目について回答させます。ひとりにつき 120
 50 の変数があり, これを何百, 何千人に対して実施するのですから, 非常に大量のデータになっているわけ
 51 です。スプレッドシートにデータを入れていくときは一般に, 一行 (横方向) に 1 ケース (1 オブザベーション, 1
 52 個人) であるようにし, 列方向 (縦方向) に変数を並べるのが一般的です。数学記号では次のように表現し
 53 ます。

$$x_{ij}$$

54 ここで i は個人を, j は変数を指します。たとえば x_{13} で第 1 番目の個人 (ケース) の第 3 番目の変数 (に
 55 ついての値) を意味するわけです。このように一般化することで, 120 ある変数でも何千分ものデータでも 1
 56 つの記号で表現できますね。

57 さて, このような大量のデータを今から分析していくわけですが, どのような方法があるでしょうか。データ
 58 分析の領域には, さまざまなモデル, 手法があります。数え方にもよりますが, 何百という種類の分析があ
 59 るかもしれません。もちろん似通ったものもありますし, 同じ目的に使う異なる手法もあります。これを大きく
 60 分けるなら, まず「線形モデル」と「非線形モデル」に分割できます。

61 **■線形モデル** 線形モデルは, 回帰分析や要因計画などが含まれます。関係が直線的であること, つまり変
 62 数に 2 乗, 3 乗の項が入っていないので, 同じパターンで先々まで考えることができる, 単純な関係です。線
 63 形というのは $y = ax + b$ のグラフを書くと同様に, 直線的な関係になることからきています。変数が
 64 x, y だけでなく, $y = ax + bz$ のように 2 つあったとしても, グラフでは線が面になるだけで, ある断面で見
 65 と直線関係であることに変わりはありません。実際が多変量データではたくさんの項が含まれ, 関数全体を可
 66 視化することは不可能ですが, 次元が多くなっても直線関係であることに変わりはありません。一般に線形モ
 67 デルは次の形で表現されます。

$$y = a_1x_1 + a_2x_2 + \cdots + a_nx_n + b$$

ここで a_m は第 m 番目の変数 x_m につく**係数 (coefficients)** であり、変数の重要性を示す数字です。足し合わせる時に、変数の重要性を変えるものなので**重み (weight)** と呼ばれることもあります。この**重みつき線型結合**が線形モデルの基本です。

線形モデルの代表的な分析方法としては、**回帰分析**、**重回帰分析**、**パス解析 (Path Analysis)**、**階層回帰分析**、**因子分析**、**主成分分析 (Principle Component Analysis)**、**共分散分析**、**判別分析**などが含まれます。また今あげたモデルをすべて含んだ表現形式である**構造方程式モデリング (Structural Equation Modeling)** があります*2。構造方程式モデリングは総合的な表現方法で、先にあげたモデルを下位モデルとして含むものですから、これをしっかり学べば各手法をいちいち学ぶ必要がない、とも言える究極的なモデルです。本講義では第 12 講で触れることになります。

線形モデルの中には他にも**グラフィカルモデリング (宮川, 1997)** などが含まれますが、それは本講の範囲を超えるので専門書に譲ります*3。さらに言えばこれら線形モデルのほとんどは**正規分布 (Gaussian Curve)** を仮定した確率モデル群だと言えます。合わせて正規線形モデルといえます。

■**非線形モデル** (正規) 線形モデルは、変数間関係を直線的なものだと考えるのでした。しかし、世の中のことは必ずしも線形関係ではありません。むしろ線形でない関係の方が一般的でしょう。線形関係というのは $A \rightarrow B$ のように、「こうすれば、こうなる」というわかりやすい関係ですが、人間の場合とはくに「叩いたら、泣く」「優しい言葉をかければ、喜ばれる」といったことでも成立しないことがいっぱいあるわけです。叩かれた人が、強がって見せるためにグッと我慢するとか、優しい言葉に絆されて泣いてしまうといった人間の感情の機微は、本当に興味深く複雑な仕組みですよね。

線形モデルは、現状に当てはまらないこともありますが、わかりやすさを優先して作られたモデルです。それに対して、現状に当てはまることを目的にすると、とても線形の関係では無理です。そこで非線形な関係でもいいから、データに適したモデルはないか、と考えられているのがこの非線形モデルです。非線形だからといって、曲線である、というだけではありません。たとえば条件分岐のように、枝分かれしていくような関係なども含まれます。

非線形モデルの代表的な分析方法としては、**決定木**、**ランダム森**、**サポートベクターマシン**、**ニューラルネットワーク**、**ベイジアンネットワーク**、**アソシエーションルール**、**自己組織化マップ**などがあります。入門書としては**豊田 (2008)** などが網羅的で良いですが、より専門的には**機械学習 (Machine Learning)** などのキーワードで選書すると良いでしょう。ここでいう学習とは、データに合わせて重みを調節する方法を指し、機械が自動的に重みを調節していく様を「学習している」と表現しているわけです。巷で A.I.(人工知能) と呼ばれているものはこうした手法の総称で、データに適していることを目的にしているのです。パターンがわかれば行動の予測に使えます*4。行動の予測ができれば、たとえば小売業では売り上げに直結する戦略が取れるでしょうし、犯罪者のプロファイリングなど、応用的側面はいろいろ考えられます。

じゃあもう非線形モデルが最強じゃないか、と思うかもしれませんが、この系列の総合的な問題点は「機械がなぜそのような予測をしたのか説明できない」という点にあります。機械には学び方を教えていますが、どう学んだかは機械次第なのです。理屈はわからなくても正解が出せる、というのは実用的にはいいのですが、科学的な研究の場合は少し困ります。機械が勝手に人の心を「うんうん、わかりますよ」と言ったとしても、

*2 これは**共分散構造分析 (Covariance Structural Analysis)** と呼ばれることもあります。

*3 この方法は類似の名称があること、あまりメジャーな分析方法でないこともあって、専門書もすくないのですが、技術としては非常に興味深い手法です。SEM が線形モデルの王道であり線形関係を正面から捉えているのに対し、グラフィカルモデリングは関係を裏から捉える裏線形ともいうべき手法です。もう少し言葉を足すと、SEM がどこに線形関係があるのかを探っていくのに対し、グラフィカルモデリングはどこに関係がないか、どこどこの関係が弱いか、というところを探して無意味な関係を切断していき、残った関係が考察すべきものだ、と考えるのです。

*4 線形回帰モデルで予測しただけでも、知らない人にとっては人工知能＝機械が計算した予測式だ、と思ってしまっているかもしれません。

103 「じゃあどうわかったのか、理屈を教えてください」といっても答えられなかったり、同じアルゴリズムでも違う
 104 機械が学習すると違う重み係数になったりして一般的な理屈が出てこないということがあります。心理学は実
 105 践的な側面もありますが、究極的には人間行動の理論を探しているのです、そういう意味では非線形モデルは
 106 向いていないのです。

107 ■モデルの展開と全体像 多変量解析の分析方法を、線形モデルと非線形モデルに分割して説明してきま
 108 した。一言でいうなら、線形モデルは「当てはまらないこともあるけど、理屈で説明がつくモデル」であり、いわ
 109 ば理屈が先、現象が後です。非線形モデルは「理屈はわからないけど、データには当てはまるモデル」であり、
 110 いわば現象が先、理屈が後なのです。

111 どちらもデータとの当てはまりを基準にはしていますが、その背後の考え方が違っているわけですね。(正
 112 規) 線形モデルも非線形モデルも、制約を増やしたり減らしたりしながら限界とされている点を克服するべく、
 113 日々モデルの改良がなされています。

114 たとえば線形モデルの世界でも、正規分布以外の**確率分布**を仮定できます。それらは**一般化線形モデル**
 115 (**Generalized Liner Model**) と呼ばれ、さまざまな確率分布を仮定した上で、線形関係を見出す手法
 116 です。心理学の研究対象としているデータは、目に見えない心的状態を対象とすることが多いので、正規分
 117 布を仮定することが一般的でした。もちろん数学的に、正規分布を仮定するモデルの方が単純な形になるの
 118 で、そちらの発展が先に進んだという実情もあります。正規分布以外の形をするデータがなかったわけでは
 119 ありません。所得のデータは対数正規分布のようになりますし、比率のデータは 0 から 1 までの範囲にしか
 120 値を取らないので平均が 0.5 からズレれば左右対称の形にはなりません。友人の数を数える時のように、0,
 121 1, 2... と正の整数しか取らない離散的なデータというのがあります。しかし分析モデルが正規分布を仮定し
 122 たものしかなかったので、これまではデータを正規分布の形になるように変換して分析する、ということが行
 123 われてきました。今は一般化線形モデルを使って、データの形式にあった分析をすることが基本です*5。

124 このように、正規分布の線形モデルが非常に多くあるのですが、その制約を外す方向で統計モデルが展開
 125 してきているわけです。制約の外し方としては、ここで挙げた「正規分布以外の確率モデルを使う」ということ
 126 もありますし、分布の歪度・尖度といったより高次の**積率 (moment)** を使う方法があります*6。

127 また、**数理モデリング (Mathematical Modeling)** というアプローチもあります。これは線形の仮定
 128 を外し、理屈の通る変数関係を数式で表現してデータにフィットさせるというアプローチです。当然非線形な
 129 関係になりますし、正規分布以外の確率分布も使います。この時、さまざまな確率分布が入れ子になったモデ
 130 ルになるので、推定方法としてはベイズ法を使うことになります。**ベイズ法**によるモデリングアプローチは最近
 131 の計算機技術革新によってとても身近なものになりました。この授業でも第 16 講から扱うことになります。

132 さて、多変量データを扱う世界の全体像がわかったところで、まずは線形モデルの世界から少しずつ進ん
 133 でいくことにしましょう。そこで、統計モデルの種類にも大きく関わる尺度水準や、記述統計量について、改め
 134 て説明を加えておきたいと思います。

135 1.2 尺度の四水準

136 心理統計のテキストは何を開いても、まず **Stevens (1946)** による**尺度水準 (level of measurement)**
 137 についての言及があります。何を測定した数値であるかとは別に、その数値にどのような算術処理を施すこと
 138 ができるかによって、数字を 4 つのレベルに分けるのでした。

*5 古い論文を読むと、正解率のような比率のデータに対して**分散分析**を行ったりしていましたが、このような理由から今では推奨されません。

*6 積率について、詳しくはこの後の 1.3 節で。

139 ■**名義尺度水準** **名義尺度水準 (nominal scale)** は数字と対象が 1 対 1 で対応していることだけが重
 140 要です。男性を 1, 女性を 2 とコード化するようなもので、この時「女性は男性の 2 倍である」といった数とし
 141 ての意味はありません。男性を 0, 女性を 42 としても本質的に変わりがないからです。このような名前だけの
 142 数字は、計算ができませんので、せいぜい 1 が何件あったかという度数を数えて集計するにとどまります。

143 とはいえ、数字が直接対象を指し示しているわけですから、もっとも意味のある数字かもしれません。

144 ■**順序尺度水準** **順序尺度水準 (ordinal scale)** は、数字が大小関係の意味を持っているものです。レー
 145 スで 1 位 2 位と順番がつくと、1 位のほうが 2 位より優れていることがわかります。人間の心理的な反応、と
 146 くに 5 段階や 7 段階で評定させる心理尺度は、この水準に相当します。選好の順序は明確でも、量的な違い
 147 がわからないからです。

148 ■**間隔尺度水準** **間隔尺度水準 (interval scale)** は、数字と数字の間隔が等しいことが制約として加わ
 149 ります。たとえば気温で 10 °C と 20 °C の差は、25 °C と 35 °C の差に等しいと言えます。これは摂氏が氷点を
 150 0 °C、沸点を 100 °C としたうえで百等分したという定義から明らかかなことです。間隔が整っているので加法・
 151 減法の計算は可能ですが、原点が定かでないので比を考えることはできません。たとえば 10 °C は 20 °C の倍
 152 の熱量を持っている、とは言えないのです。なぜでしょうか。たとえば、同じエネルギー状態を別の温度体系に
 153 置き換えてみるとしましょう。新しい温度体系は、氷点が 100 で沸点が 200 だったとします。そうすると 10 °C
 154 は 110, 20 °C は 120 に該当しますが、120 は 110 の 2 倍にはなっていないからです。

155 ■**比率尺度水準** **比率尺度水準 (ratio scale)** は、さらに絶対 0 点の制約を付け加えたものです。これで
 156 原点からどれ位離れているか、を基準にして計算ができますので、乗法・除法もできることになりました。物理
 157 的な単位系はこの尺度水準にあるものがほとんどですから、緻密な計算モデルを作ることができるのですね。

158 さて早足で 4 つの水準について説明をしてきました。これが重要なのは、尺度水準によってできる計算が
 159 変わってくる点にあります。名義尺度水準は数え上げぐらいしかできません。順序尺度水準も同様で、名義や
 160 順序といった**質的変数 (categorical variables)** の場合、たとえば代表値を求める時も度数を数えて最
 161 頻値を報告する、というぐらいがせいぜいなのです。

162 これに対して、間隔尺度水準や比率尺度水準の**量的変数 (numeric variables)** では、加減乗除の計
 163 算ができますので、平均値を求めたり標準偏差を求めたり、ということができるようになります。

164 先ほど心理尺度は順序尺度水準でしかない、という話をしましたが、質問紙調査の研究例では尺度平均点
 165 を出したり、さらに進んだ統計手法で分析したりします。実はそれができるようになるためには、尺度につけら
 166 れたカテゴリー（「非常に当てはまる」「どちらとも言えない」など）を、尺度値（「非常に当てはまる」を 5、「や
 167 や当てはまる」を 4 とする、など）にする作業を経ており、その時に「非常に当てはまる」を 5 とするのはなぜ
 168 か、という理屈が必要です。それについては心理尺度の作成の折に詳しく説明しますが*7, 少なくとも盲目的
 169 に行っているわけではないことに注意が必要です*8。

170 尺度値の付与の仕方は色々考えられます。順序尺度水準でとられたデータであっても、その背後に連続的
 171 な心理的実態があると考えてその時の相対的な大きさをつけることもできます*9。あるいは 0 か 1 かという
 172 二値データであっても、それを重みづけて合算することで連続的な値にすることもできます*10。さらに名義
 173 的な尺度水準であっても、データ全体なかで直線的な関係が最も大きくなるように数字を与えることもでき

*7 第 3 講を参照してください。

*8 残念ながら、この辺りの理屈を気にせずに分析する人が少なくありません。そんな人に、「数字では心がわからない」なんて言って欲しくないですね。

*9 因子分析法の一種、**段階反応モデル (Graded Response Model)** と呼ばれる手法がこれにあたります。

*10 これはテスト理論の考え方です。0 が誤答、1 が正答としてテストの点数から学力を考えるのですね。

174 ます*11。

175 このように、尺度水準が変わると、適用できる分析モデルが変わります。たとえば因子分析は間隔尺度水準
176 以上のデータにしか適用できませんが、順序尺度水準のデータであれば因子分析ではなく、段階反応モデ
177 ルのようなカテゴリカル因子分析を適用しなければなりません。名義尺度水準のデータであれば双対尺度法
178 を適用しなければなりません。間隔尺度水準以上のデータに対して、双対尺度法を適用することは可能で
179 すが、その逆は不可能です。たとえば身長データとして $X = \{170, 175, 165\}$ というのがあったときに、連続
180 変数として平均値を求めることもできますし、名義尺度水準として各一件とカウントすることもできるので
181 すが、男性 2 名と女性 1 名がいたので平均 1.3 の性別があるというのは意味をなさないからです。

182 データがどの水準にあるのかを見極められないと、間違えた分析をすることにもなりますので、注意してく
183 ださい。

184 1.3 平均と分散

185 間隔尺度水準以上の数字であれば、平均値や標準偏差の計算が可能です。(算術)平均 (mean) は次の
186 式によって計算されるのでした。

$$\bar{x}_j = \frac{1}{N} \sum_{i=1}^N x_{ij}$$

187 総和の記号である \sum の使い方などを再確認しておいてください。また変数 j の平均は \bar{x}_j と表しますが、
188 \sum の記号の中では $i = 1$ から N までと、 i だけが変化しています。 j は変化していません。
189 \sum の記号はこの後も所々出てきます。計算の際に次のような変形をすることがありますので確認してお
190 いてください。

$$\sum_{i=1}^N cx_i = (cx_1 + cx_2 + \cdots + cx_n) = c(x_1 + x_2 + \cdots + x_n) = c \sum x_i$$

191

$$\sum_{i=1}^N (x_i + y_i) = (x_1 + y_1 + x_2 + y_2 + \cdots) = \sum x_i + \sum y_i$$

192 続いて**分散 (variance)** の式を確認しましょう。

$$s_x^2 = \frac{1}{N} \sum (x_i - \bar{x})^2$$

193 言葉でいうなら、平均偏差 $(x_i - \bar{x})$ の二乗の平均です。平均偏差は各点が平均点からどれくらい離れて
194 いるかを表します。その平均をとる操作 $(\frac{1}{N} \sum)$ なので、平均的にどれくらい離れているかがわかるので
195 すが、そのまま平均偏差の平均をとるとゼロになりますので*12、二乗するものをその値にするのでした。

196 この式は次のように展開できます。

*11 双対尺度法という考え方がこれにあたります。詳しくは西里 (2010) を参照。

*12 平均値が全部のデータの真ん中に位置するように撮られた指標だから当然です。

$$\begin{aligned}
s_x^2 &= \frac{1}{N} \sum (x_i - \bar{x})^2 && \text{定義より} \\
&= \frac{1}{N} \sum (x_i - \bar{x})(x_i - \bar{x}) \\
&= \frac{1}{N} \sum (x_i^2 - 2x_i\bar{x} + \bar{x}^2) && \text{(後ろのカッコを展開)} \\
&= \frac{1}{N} \sum x_i^2 - \frac{1}{N} \sum (2x_i\bar{x}) + \frac{1}{N} \sum \bar{x}^2 && \sum \text{記号の分配} \\
&= \frac{1}{N} \sum x_i^2 - 2\bar{x} \frac{1}{N} \sum x_i + \frac{1}{N} N\bar{x}^2 && i \text{ が変化するところだけにつく} \\
&= \frac{1}{N} \sum x_i^2 - 2\bar{x}^2 + \bar{x}^2 \\
&= \frac{1}{N} \sum x_i^2 - \bar{x}^2
\end{aligned}$$

197 計算するときには最後の形を使うことがあります。

198 この式で表される分散は、変数がどの程度変化しうるかということを表す値であり、変数から引き出せる情
199 報の上限でもあります。分散が 0、つまり変数がまったく変化しなければ、どういときに値が大きくなってど
200 ういときに値が小さくなるのかという違いがわからないということです。違いがないものについては、それ以
201 上考察のしようがないか、当たり前のことを言っているに過ぎないということになります。たとえば質問紙調
202 査で「他人に激しく殴られるのが好きですか」という聞き方をすると、誰も「まったく当てはまらない」と答え
203 ると思います*13。この項目から何かがわかるか、といわれても当たり前のことすぎて何もわからない、としか
204 言えないでしょう。これに対して「対面している相手の手足の動きが気になりますか」というような項目であ
205 れば、気にする人もいるでしょうし、気にしない人もいるでしょうから、回答に分散が生まれます。そうすると、気
206 になった人はどういう人なのか、気にならなかった人はどういう人なのか、という考察に進むことになるわけ
207 です。このように、調査データから意味のある考察をするためには、分散の大きな項目を作らなければならない
208 のです。

209 ところで、分散の式の中には二乗の項が入っていますから、このままでは元のデータと単位が異なっ
210 てしまいます。そこでこの単位を整えるために、分散の正の平方根をとったものを**標準偏差 (Standard**
211 **Deviation)** といって散らばりの指標に使います。本講では標準偏差の記号を s_x と表すことにします。

212 平均は中心化傾向の指標の一種で、他にも中央値、最頻値などがあります。分散や標準偏差は散らばりの
213 指標の一種で、他に最大値、最小値、範囲、IQR、パーセンタイルなどがあります。これら記述統計量は、デー
214 タの特徴を記述するため、データ分析の最初のステップで確認すべき数字です。また、分散は平均偏差を二
215 乗したものの平均でしたが、これを三乗したものは**歪度 (skewness)** といいます。歪度 s_x^3 の式は次のとお
216 りです。

$$s_x^3 = \frac{1}{N} \sum (x_i - \bar{x})^3$$

217 歪度はマイナスになると左方向に、プラスになると右方向に、分布が歪んでいることを示します。ゼロに近
218 れば左右対称に近いことがわかります。これも記述統計量としてデータの特徴記述に利用できるでしょう。さ
219 らに四乗したものは、**尖度 (kurtiosis)** と言われます。

$$s_x^4 = \frac{1}{N} \sum (x_i - \bar{x})^4$$

*13 中にはマゾヒズムの人がいるかもしれない、という屁理屈はここでは脇に置いておいてください。

220 尖度がゼロであれば、分布の山の尖りぐあいが平均的で、マイナスになれば潰れた山、プラスになれば尖っ
221 た山の形をした分布になることがわかります。

222 分散が二乗、歪度が三乗、尖度が四乗でした。これらは分布の中心からどれくらい離れているかにつ
223 いての累乗で、平均値は一乗したものと理解することができます。これを力学の用語を借りて**モーメント**
224 (**moment**) といいます。日本語では**積率**と訳されています。重心からの距離に関する指標という意味で一
225 般化された表現です。多変量解析のほとんどは、二次のモーメント (分散) までしか活用しませんが、3 次、4
226 次のモーメントを使うとなるとデータから得られる情報が増えることになり、より表現力を増した分析ができる
227 ようになります*14。

228 1.4 共分散と相関係数

229 さて、分散が 1 変数の変動を表現する数字であり、そこから得られる情報の上限であるという説明をしま
230 した。実際に調査的な研究をするときは、いくつもの変数を同時に扱うことになるわけですが、そうすると変数
231 と変数の関係を考える必要があります。

232 改めて分散の式を見ると、次のようになっているのでした。

$$s_x^2 = \frac{1}{N} \sum (x_i - \bar{x})^2 = \frac{1}{N} \sum (x_i - \bar{x})(x_i - \bar{x})$$

233 この式を少し変えて次のようにすることで、変数 x と y の関係を考える式になります。

$$s_{xy} = \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})$$

234 この式で表される数字を**共分散 (covariance)** といいます。異なる変数ですので異なる記号で現しましたが
235 が、変数を x_{ij} のように個人 i と変数 j という形で表すならば、変数 j と k の共分散を表す式は次のように
236 書けます。添字のつき方に注意して理解してください。

$$s_{jk} = \frac{1}{N} \sum (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)$$

237 さてこの共分散は、カッコの中がそれぞれの変数の平均偏差を現しており、それを個人ごとに集めて平均し
238 ています。あるケースにおいて、平均偏差が同じ方向、すなわちどちらも平均より上であるか、どちらも平均よ
239 り下であれば、積の符号は正になるのでこの数字は増えます。逆にあるケースにおいて平均偏差が異なる方
240 向、つまり一方は平均より上なのに他方は平均より下であるようなことがあれば、積の符号は負になりますの
241 で、この数字は減ります。これを平均するわけですから、データ全体で同じ方向にずれる傾向があるのか、そ
242 うでないのかを表す指標ということになります。

243 同じ方向にずれるということは、一方の動きがわかれば他方の動きが推測できます。このように、共分散の
244 数字は変数間関係から予測、考察をするための情報量を現しているとも言えます。共分散がまったくない、0
245 であれば、ケースごとにさまざまな可能性があるので一般的なことが言えないわけです。

246 この共分散は、分散と同じで単位に依存します。たとえば身長と体重の共分散を計算すると、メートルとキ
247 ログラムをかけ合わせた単位になります。このように、変数ごとに単位が変わると共分散同士の比較が難し
248 くなりますので、データの標準化を考えましょう。すなわち、どのようなデータであっても単位に捉われずに比
249 較できるようにします。スコアの標準化は次の式で行います。

$$z_i = \frac{x_i - \bar{x}}{s_x}$$

*14 詳しくは豊田 (2007) を参照してください。

250 このように平均偏差を標準偏差で割ることで、あるデータがその変数の平均からどれくらい離れているか、
 251 標準偏差を単位とした相対的な大きさを表すことができます。標準化されたデータは、平均が 0、分散が 1 に
 252 なります。このように、標準化されたスコアは相対的に同じサイズに整えられ、単位に依存しませんので、標準
 253 化したスコア同士は相対的に比較可能です。

254 この標準化スコアしたスコアで共分散を計算したものが、**相関係数 (correlation coefficient)** です。

$$r_{jk} = \frac{1}{N} \sum z_{ij} z_{ik}$$

255 この相関係数はどんな変数であっても、 -1 から $+1$ の範囲に入りますので、変数間関係の相対的な比較
 256 が可能です。つまり身長と体重の関係の強さは、身長と足のサイズの関係の強さよりも大きいとか小さい、
 257 といった表現が可能になるわけです。

258 多変量解析の世界は、変数がたくさんありますから、この共分散の組み合わせもたくさんあります。分散が
 259 変数から得られる情報、共分散が変数間関係から得られる情報ですから、あるデータセットから得られるこれ
 260 らの情報の数はどれくらいになるか計算してみましょう。

261 変数が M 個あったとします。変数番号が $1, 2, 3, \dots, M$ とすると、組み合わせは 1 と 2 , 1 と $3, \dots$ となり
 262 ます。このとき i と j は分散、 i と j が共分散になりますが、 s_{ij} は s_{ji} と同じですから、 $M \times (M - 1)/2$
 263 個の共分散と M 個の分散が、このデータセットから得られる情報のすべてということになります。合計で
 264 $M \times (M + 1)/2$ 個になりますね。このように組み合わせまで考えると、変数の数がひとつ増えるだけでも得
 265 られる情報、分析のヒントがぐっと増えることになります。多変量解析は、こうした変数間関係の情報を使って
 266 分析を進めていくことになります。

267 ちなみに、変数間関係を表す方法は共分散や相関係数だけではありません。これらはあくまでも変数間の
 268 直線的な関係を表す指標であることにも注意が必要です。多変量解析全体としては、変数同士の数字のズレ
 269 をたしあわせた**距離 (distance)** や、ある変数が他の変数と同時に出現した回数をカウントした**共頻度**など
 270 を扱うモデルもあります。

271 1.5 課題

272 ■**尺度水準** 4つの尺度水準の具体例をそれぞれ挙げなさい。

273 ■**平均と分散** 平均、分散、標準偏差、標準得点の定義を式で表現しなさい。

274 ■**共分散と相関係数** 共分散の定義式を書きなさい。そのとき、スコアが標準化されていると相関係数にな
 275 ることを示しなさい。

第 2 章

心理尺度を作る

2.1 はじめに

心理学の研究法は調査、実験、観察の 3 つが代表的なものです。とくに調査法は、質問紙を用いて多くの人に回答を依頼し、それを統計的に分析することで研究仮説を確認しようとするものです。調査対象者をどのように集めるか、調査票をどのようにデザインするか、得られたデータをどのように分析するか、といったことでそれぞれ 90 分以上話せるぐらいに色々考えるべきことはありますが、ここではとくにその本質について考えてみましょう。すなわち、なぜ紙で用紙に丸をつけたものを集めただけで、心の何かがわかったと思えるのか、ということです。

アンケート調査は一見すると、なんのひねりもない研究法に思えます。すなわち、聞きたいことを当の本人に聞いてみる、これだけです。しかも答えやすいように（集計しやすいように？）紙に問題が書いてあって、該当するところに丸をつけるだけでよかったです。この手の研究は少なくとも 100 人、できれば数百人、大きい規模だと千以上の桁数の協力者に回答を求めます。それを PC に入力して集計するわけです*1。しかし「そう思う」を 5 点、「ややそう思う」を 4 点、以下 3, 2, 1 点・・・と入力するのはなぜでしょう。これらの反応は**順序尺度水準**でしかない、と言われますが、それを**間隔尺度水準**であると「みなし」て、平均値を求めたりします。なんでそんなことが許されるのでしょうか。

本稿ではその謎について迫っていきたいと思います。

2.1.1 測ろうとしているもの；態度

調査票で質問する内容は、（お客様アンケートとかではなく）心理学の研究であれば**尺度 (scale)**を使います。尺度とは、測定したい対象に数字を割り振るルールのことであり、心理尺度は項目とその採点方法が規格化されたものを指します。心理尺度の測ろうとしているものの多くは**態度 (Attitude)**と呼ばれるものです。これは社会心理学の用語で、簡単に定義すると態度とは「行動の準備状態」のことです*2。

社会心理学も心理学ですから、研究は基本的に観察可能で客観的なものを対象にします。社会的な行動をテーマにするというのがもちろんですが、社会的な行動というのは状況によって変わるものですし、研究したい状況を待っていても自動的に出てくるものではありません。もちろん実験室などでその状況を作り出すこともやりますが、選挙のような公的なものであれば実験的に作り出すこともできません。しかしたとえば選挙結果などは、社会心理学における政治的行動に大きな影響を及ぼす（あるいは政治的行動の結果になる）ものですから、研究としては非常に重要なテーマになります。選挙になるまでは仕事ができない、というのでは

*1 もっとも最近では、web で調査の回答を得ることで、入力や誤入力のチェックなどの手間が大幅に軽減されています。

*2 この定義は Allport (1967) によるものです。

304 社会心理学者も困りますから、「あなたはどこに投票するつもりですか」と聞いて普段の政治的な態度を調べ
305 るという研究手法ができました。社会調査の始まりです。

306 そしてこれを心理学一般に応用しているのが、質問紙調査です。心理学の研究テーマは行動ですが、行動
307 を作り出せない場合や普段の状態を査定したいときに、「調査票に回答を求める」という方法を取るのです。
308 もちろん調査に対する回答は、嘘や見栄、間違いや勘違いなども入り込む可能性がありますので、設計は丁寧
309 に、分析は慎重に行う必要があります。調査票で過去のことを聞いたりすると記憶が歪んでいる可能性が
310 ありますし、未来のことを聞いたら嘘八百を答えられるかもしれません。今の気持ちを聞いても、自分の現在
311 の状況ははっきりわかっているかどうか、怪しいものです。それでもある程度の真実味はあるだろうと考
312 えて、こうした阻害要因を取り除いて本質を掴むために、いろいろな工夫がなされています。

313 尺度を使って感情や気分、考え方（認知傾向）や今後どう行動するつもりか（行動意図）を調べることがな
314 されていますが、本質的にはこれらすべてが、広い意味での態度です。態度の認知的側面、感情的側面、行
315 動的側面などと言われたりします（藤原, 2001）。この「態度」の説明や定義は、実はそれほど明確ではありま
316 せん。しかし、一般に特定の対象に対する、正負・量的な評価が可能なものと考えられています。たとえば「自
317 民党に対する態度」というのは、自民党という対象に対して、「好き」「嫌い」というポジティブ・ネガティブの評
318 価ができ、さらに「とても好き」「やや嫌い」のように量的に表現できるものでもある、と仮定されます。多くの
319 人に多くの項目で調査するのは、項目同士の誤差が相殺しあってこれが正規分布すると考えられるからであ
320 り、社会的な態度は極端な値が少なく平均的な値を持つひとが多いものとして、相対的に評価されます。

321 基本的に測定しようとしているものは、こうした特徴を持ったなにかである、ということをもまずは踏まえてお
322 きたいと思います*3。

323 2.1.2 3つの方法

324 心理尺度の作り方には大きく分けて3つのスタイルがあります。

325 1つ目はサーストーン法による尺度で、等現間隔法 (method of equal-appearing intervals) と呼ば
326 れるものです。2つ目はリッカート法 (Likert 法, ライカートと読む人もいる) と呼ばれるもので、5件法、7件
327 法など数段階のカテゴリラベルのもっとも近いところに丸をつけるという方法です。3つ目はSD (Semantic
328 Differential 法, 意味微分法と訳されることも) 法とよばれるものです。SD 法は態度測定というより、イメ
329 ージの測定を目的としたもので、対象を提示しつつペアになった形容詞を列挙して提示します。たとえば「専修
330 大学」という対象に対して、「激しい-落ち着いた」「慎重な-軽快な」といった形容詞対ではどちらの表現が近い
331 かを評定してもらいます。形容詞の対が作る軸上でいうと平均的にどのあたりに対象がプロットされるのか、
332 を見ることで対象ごとのイメージの違いを表現するのが基本的なアイデアです。SD 法は複数の対象に対す
333 るイメージの相対的比較ですから、スコアの点数化にはそれほど重きを置いていないのでここでは取り上げ
334 ません。

335 サーストーン法とリッカート法は、これを使って態度のスコアをつけることができます。小杉の自民党に対する
336 態度は4.8点だ、といったように、です。こうした数値化がどのような理屈でなされるのかを、今からみてい
337 うと思います。

*3 性格のように特定の対象を持たないものであっても、正規分布に従うと考えるのは自然ですから、この後説明する統計技法が適用できるものもあります。感情や気分といったものは、持続時間が短いので生理指標などで測定するべきであり、調査票によるアプローチは不向きですが、逆に言えばある程度一定の安定した心理状態であれば測定することができると考えられているのかもしれない。

2.2 サーストンの等現間隔法

サーストンの等現間隔法^{とうげんかんかくほう}は、社会的な態度について絶対評価を与える方法です。この方法で作成された尺度は、各項目（態度表明文と呼ばれます）に尺度値がついており、回答者は提示された項目に賛成であればその尺度値がその人の態度得点になります。この尺度値は事前に複数の評定者によって決めておく必要があります。すなわち、尺度を作る前の入念な準備が必要です。また、サーストンの尺度は 1 次元性を有していることが前提となります。

具体的な作成方法は次のような手順で行います。

1. 項目の収集
2. 評定者集団による評定
3. 尺度値の算出
4. 項目の選定

以下順に説明します。

■項目の収集 まずは測定したい社会的態度のテーマに沿って、項目を準備します。たとえば「自民党に対する態度」のように、誰でも思い描ける具体的な対象が良いでしょう。このようなテーマが決まれば、これに対する態度項目を色々考えます。「自民党のことを考えると夜も眠れない」とか「自民党に関係したニュースはなるべく見るようにしている」「近所の自民党員の事務所に行くことがある」というポジティブな態度もあるでしょうし、「自民党のニュースはなるべく聞きたくない」「自民党には投票しない」「自民党は不正まみれの悪い政党である」といったネガティブな態度もあるでしょう。こうした文言をなるべく多く、強い態度から弱い態度まで、ポジティブなものからネガティブなものまで網羅的に準備します。ニュートラルな項目も考えておく必要があります。

■評定者集団による評定 尺度値を決めるための事前準備です。まず評定者を無作為に集めます。少なくとも十数人は必要でしょう。評定者には事前に準備した項目が好意的－非好意的（または肯定的－否定的）の 1 次元にそって、7～11 段階ぐらいの多段階に分類してもらいます。「自民党のことを考えると夜も眠れない」というのは非常にポジティブなので 11 点、「自民党には投票しない」というのはかなりネガティブなので 2 点、といったようにです。

■尺度値の算出 このように各態度表明文を複数人で評価してもらいますから、その項目の平均値、中央値、分散などの記述統計量を計算できます。この中央値（あるいは平均値）をその項目の尺度値とします。ただし、ここで分散が大きい項目は、評定者によって評定の仕方がバラバラだということを意味しますよね。値が人によって定まらないというのは、その項目が刺激としてあまり好ましくないと考えられるので、項目候補から削除します。誰がみても 10 点とか誰がみても 3 点、といった分散が少ない項目が望ましいでしょう。

■項目の選出 さてこうして尺度値が計算できたら、それを順に並べていきます。「夜も眠れない」は 10.7 点、「事務所に行くことがある」は 9.5 点、「ニュースをなるべく見る」は 7.9 点・・・というようにしていくことができますね^{*4}。このとき、項目間の間隔が均等になるように項目を選別します。たとえば「夜も眠れない」と「事務所に行くことがある」の間隔は 1.2 点ですが、「事務所に行くことがある」と「ニュースをなるべく見る」の間隔は 1.6 点になっています。これでは等間隔と言えないので、1.2 点間隔すなわち 8.3 点ぐらいの項目を選

^{*4} 中央値なのになぜ小数点があるのだ、と思う人がいるかもしれません。平均値でもいいですし、評定者が偶数人の場合は中央値も両得点の平均や重みつき平均で小数点が出るからです。

373 出します。

374 このことからわかるように、サー斯顿法で尺度を作る場合は、事前に多くの項目を準備しておかないと
375 「ちょうどいいところの表明文がない」となってしまう恐れがあります。ですから最終的にできる尺度に含まれ
376 る項目の、5 倍から 10 倍ぐらいを事前に準備し、うまく等間隔に項目が選出できるようにしなければなりま
377 せん。

378 なぜ等間隔に選ぶのかというと、もうお分かりですね、これで得られる尺度値を**間隔尺度水準**として扱
379 たいからです。間隔が等しくなければ順序尺度にしかありませんが、間隔が等しいことがわかっていると、平
380 均や分散などの計算をし、相対的な比較をできるからです。またこの尺度を使うときは、すでに評定者集団に
381 よって尺度値がわかっていますから、回答者にずらりと並べられた尺度を見てもっとも自分の意見に近い項
382 目を選出してもらえば、その項目の尺度値がその人の態度得点だということができます。

383 評定者集団をなるべく偏りなく多く集めることで、事前に尺度の値を確定させておき、あとは本来研究対象
384 にしたかった人にその尺度を当てれば尺度値 (尺度得点) が求められる方法ですから、準備が大変だけど使
385 うときは確実に絶対的なスコアを与えることができるというのがこの方法の利点です。欠点はその準備コスト
386 の高さで、1 次元的な態度しか用いられないことでしょうか。また尺度構成の観点から重要なのは、選出プロ
387 セスによって項目間の尺度値が均等であることが保証されている点です。均等に選んだ後で、1, 2, 3, 4,
388 5 と数字を付け直しても構いません。大事なのは、こうしたプロセスのおかげで**間隔尺度水準**が維持され、以
389 後の分析に耐えうるスコアになっているという点です。

390 2.3 リッカートのシグマ法

391 次に紹介するのはリッカートの**シグマ法** (σ method) です。これはいわゆる 5 件法、7 件法と呼ばれる
392 採点方法で、「私は自民党の政治のやり方が好きだ」といった項目に対して、「まったく当てはまる」「かなり当
393 てはまる」「やや当てはまる」「どちらとも言えない」「やや当てはまらない」「あまり当てはまらない」「まったく当
394 てはまらない」といった順序づけられたカテゴリーにたいしてもっとも自分の考え・態度と近いところに丸をす
395 る、という方法で反応が得られます。この時の反応カテゴリーが、今回は 7 つありますから 7 件法 (7-points
396 scale) で回答を求めた、などと言います。5 段階なら 5 件法、4 段階なら 4 件法です。普通は「どちらとも言
397 えない」というところを用意するために奇数 (3,5,7,9) 件法を使いますが、日本人は「どちらとも言えない」を
398 選びやすいという中庸傾向があるとも言われていますので、意見をはっきりさせるために 4,6 件法も使われ
399 たりします。

400 項目はこれも複数あって、たくさん集められた項目を分析するために、もっとも当てはまるを 7、かなり当て
401 はまるを 6、以下同様にまったく当てはまらないを 1、とコード化し分析するのが一般的です。ただし注意
402 して欲しいのは、もっとも当てはまる = 7 としたのは**名義尺度水準**の数字の割り当て方と一緒に、このカテゴ
403 リーが 7 という尺度値を持っているわけではない点です。そもそもこの評定カテゴリーは、統計学的にはせい
404 ぜい順序尺度水準の性質しか持っていませんから (もっとも >かなり>やや)、カテゴリーに割り当てた数字
405 からそのまま平均や分散の計算をするのはおかしいはずなのです。もしこれらのカテゴリーの間隔が等し
406 かつとしても、3, 4 段階しかないようであればやはり**間隔尺度水準**の計算ができるほどの精度は持っていま
407 せん。数量的に分析するには (等間隔が担保された上で) 9 から 11 段階は必要と言われてます。

408 それではリッカート法ではどのようにして尺度値を決めるのでしょうか。リッカート法も測定しようとしてい
409 るのは態度であって、表に出てくる反応カテゴリーの背後には連続的な心理的態度というのが、と仮定し
410 ています。またこの (社会) 心理学的態度は、向きと大きさがあって正規分布を仮定できます。リッカート法も
411 正規分布に従う潜在的な連続変数があると仮定するのです。

412 さて、ある項目について、多くの人からデータを集めて「もっとも当てはまる」「かなり当てはまる」といったカ

413 テゴリーごとの集計ができたとしましょう。多くの人の態度も集積すれば正規分布に従いますから、きっとこの
 414 ヒストグラムも正規分布を反映したものになっているはずです。しかし我々が知りたいのはその背後にある連
 415 続体上の数字なわけです。

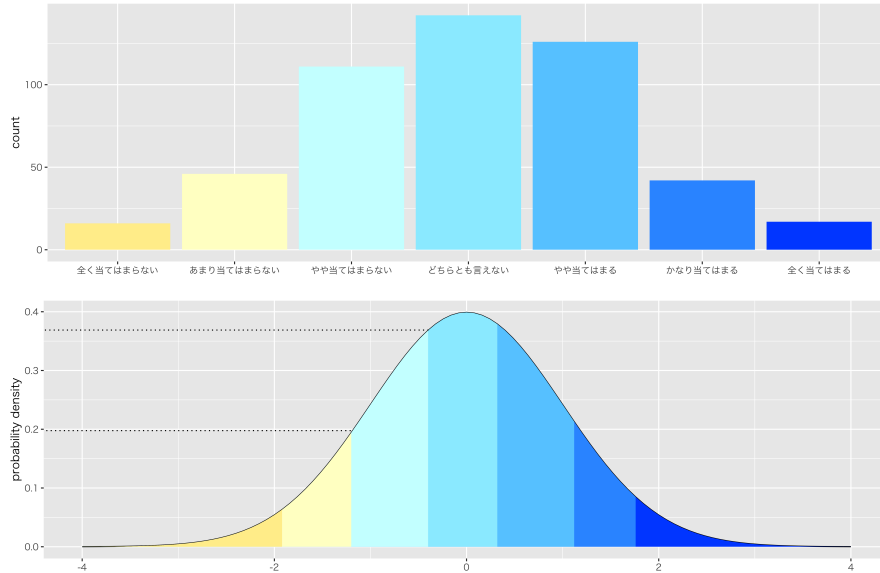


図 2.1 カテゴリ反応と背後の連続値

416 ここで図 2.1 を見てください。上段にあるのがある項目のヒストグラムの例です。しかし知りたいのは、下段
 417 にあるような正規分布の形をした連続体の変数のはずです。上段のヒストグラムは下段の状態を反映してい
 418 るはずですから、上段のカテゴリの相対頻度を元に、下段の正規分布を分割します。具体的な数字との対応
 419 は表 2.1 を見てください。出現度数を相対頻度にし、正規分布の面積を順に分割していくことになります。カ
 420 テゴリの下の方から順に分割するというので、表 2.1 の三段目には累積 (相対) 頻度を書いてあります。そ
 421 してこれを使って、標準正規分布の下から面積を考えます。統計環境 R では、qnorm 関数をつかうと累積
 422 確率の確率点が求められるのでした (表 2.1 の 4 段目)。またその時の確率密度も求めてあります (表の 5
 423 段目。R では dnorm 関数を使って求めます)。

424 さて、ではここからどのようにして尺度値を求めればいいのでしょうか。一般に C 件法で、下から
 425 $1, 2, 3, \dots, c, \dots, C$ とカテゴリ順に数字を割り振ったとして、第 c カテゴリの尺度値 Z_c を考えるとしま
 426 す。このカテゴリ c は標準正規分布において上限 z_c 、下限 z_{c-1} の確率点で挟まれる領域としていますから、
 427 この幅 ($[z_{c-1}, z_c]$) の平均を取ることを考えます。この点は、次の式で求められます (証明は付録 B, Pp.377
 428 参照)。

$$Z_c = \frac{(y_{z_{c-1}} - y_{z_c})}{p_c}$$

429 ここで y_c は z_c, z_{c-1} における確率密度、 p_c はカテゴリ c の相対頻度です。具体例でいきましょう。表 2.1
 430 の数字を使うと、「やや当てはまらない」($c = 3$) の尺度値は、 $\frac{0.20 - (0.37)}{0.22} = -0.7727273$ となります (分
 431 子は図 2.1 の点線部の差分、分母は該当箇所面積になります)^{*5}。このようにして計算された尺度値が、表
 432 2.1 の一番下の段にある数字です。

*5 表 2.1 は小数点下 2 桁までに丸めているので、正確な値ではありません。

表 2.1 カテゴリと数値の対応

反応カテゴリ	まったく当てはまらない	あまり当てはまらない	やや当てはまらない	どちらとも言えない	やや当てはまる	かなり当てはまる	まったく当てはまる
出現度数	16	46	111	142	126	42	17
相対度数	0.03	0.09	0.22	0.28	0.25	0.08	0.03
累積相対度数	0.03	0.12	0.35	0.63	0.88	0.97	1.00
累積相対度数の確率点	-1.85	-1.16	-0.40	0.33	1.19	1.83	∞
累積相対度数の確率密度	0.07	0.20	0.37	0.38	0.20	0.08	0.00
付与される尺度値	-2.33	-1.44	-0.77	-0.04	0.72	1.50	2.67

433 このように、累積度数をつかって尺度値を決めるリッカートの方法を**シグマ法**と言います。こうして作られた
 434 尺度値は連続体上の数字ですから間隔尺度水準になり、平均や分散をはじめとした数値計算に耐えうる値
 435 になっています。機械的に「まったく当てはまらない」から「まったく当てはまる」まで、1.0 刻みで数字を割り
 436 振っているのではないのです！

437 …とりたいところですが、今回の尺度値を眺めてみるとそこそこ等間隔（間隔は 0.8 ぐらいでしょうか）
 438 に並んでいますね。試しに各尺度値を 0.8 で割ってみると、 $-2.91, -1.81, -0.97, -0.04, 0.90, 1.88, 3.33$
 439 となります。四捨五入して小数点をなくしてみると、 $-3, -2, -1, 0, 1, 2, 3$ となりますね。そう、つまり**非常に**
 440 **ラフな近似値**でよければ、**機械的に 1,2,3… と数字を振っても同じこと**になります。ですから、実際の研
 441 究ではとくに深く考えずに 1,2,3… と割り振った数字をそのまま使われたりするのはです。大山鳴動して鼠一匹
 442 といいますが、苦労した割に得るものがないじゃないか、とお叱りを受けそうですが、少なくとも「なぜリッカ
 443 ト法は順序尺度ではなく間隔尺度のように扱って良いのか」という問いには答えられると思います。また、ここ
 444 にくるまでに、態度の 1 次元性や正規分布の仮定などが含まれていたことを改めて思い出してください。今回
 445 は数値例ですので、綺麗に七段階に分かれるようなものを用意しましたが、実際の調査では正規分布しない
 446 ものや、平均が低すぎるとか高すぎるものが結構みられます。それらに対して機械的につけた数字で分析す
 447 るのは決して適切な方法ではなく、シグマ法を始めその他の手法で適切な尺度値を付与すべきなのですが、
 448 人間は易きに流れるもので**ほとんど考慮されていないのが現状です***6。

449 2.4 尺度を評価する

450 このようにして作られた心理尺度は、それがきちんと測定したいものを測定できているか、評価する必要が
 451 あります。ここでは**信頼性 (reliability)** と **妥当性 (validity)** の 2 つの側面から説明します。

*6 この状況は決して良いものではなく、悪しき研究上の風習だと思われます。幸い、第 4 講で説明する**項目反応理論 (Item Response Theory)** の一種、**段階反応モデル (Ggraded Response Model)** を用いると、この問題点をカバーし
 つつ有用な情報が得られますので、みなさんは一足飛びにその手法を身につけた方が良いでしょう。

2.4.1 信頼性

信頼性は測定の安定性と言い換えてもいいかもしれません。すなわち、同じものを 2 回測っても同じ数字がつくことですね。測定するたびに数字が変わるようでは、その測定器（ここでは尺度ですが）は信用ならない、というわけです。

テスト理論の文脈では、テストのスコア X を本当に測りたいもののスコア t と誤差 e とに分割して考えます。古典的テスト理論のモデル式は次の通りです。

$$X = t + e$$

ここで、各項目についても同じことが言えると考え、 $X_i = t_i + e_i$ ということになります。複数項目で測定するのは、これの平均値を考えると $\bar{e} = 0$ (誤差の平均値はゼロになる) という仮定から、 $\bar{X} = \bar{t}$ となって真のスコアを得ることができる、ということがわかります*7。また、テストの分散 $Var(X)$ を考えると、テスト全体の分散は $Var(X) = Var(t) + Var(e)$ となり、真のスコアの分散と誤差の分散に分解できることがわかります。ここから、信頼性 Rel を次のように表現できます。

$$Rel = \frac{Var(t)}{Var(X)} = \frac{Var(X) - Var(e)}{Var(X)} = 1 - \frac{Var(e)}{Var(X)}$$

言葉で言えば、信頼性は全分散に占める真分散の割合であり、全体から誤差分散の割合を引き算したものであるとも言えます。

信頼性のない尺度があれば、その後の話は先に進みませんから、まずもって「尺度が信頼できるかどうか」を評価する必要があります。このことを信頼性は妥当性の上限、と表現したりします。信頼性を評価する方法は、測定値の安定の程度を評価できればいいのですから、複数の測定を持ってその相関係数を計算することでひとまず達成できます。しかし同じ尺度を何度も使うのは、調査回答者に要らぬ構えを持たせてしまいますから、普通は 1 回の尺度を分割してその特徴を見ることにします。尺度全体から計算される回答者の値は、項目の尺度値の合計であるのが普通です (テストの点数も正答数に対応していますね)。ですから、ある項目 j の尺度値は、 j を除いた残り $M - j$ 個の尺度値の和と高い相関をすることははずです。このように、各項目が尺度全体の値とどの程度相関しているかを見る **IT 相関 (Item-Total correlations)** は、尺度の信頼性を見る 1 つの指標になります。ある項目が、尺度全体と相関していなければ、それはその項目が尺度の中で目的と違うものを測定している可能性があり、それは必然的に尺度の安定を損ねる結果になるからです。

この考え方を発展させ、各項目が他の項目とどの程度相関しあっているか、つまり尺度の中でどの程度整合性がとれたもの＝一貫して同じものを測定しているのかを評価する指標として、 **α 係数 (alpha coefficient)** があります*8。これは、テストに含まれる項目数を M 、テスト全体の分散を V_t 、項目 j の分散を V_j と表すと、次の式で表されます。

$$\alpha = \frac{M}{M-1} \times \left(1 - \frac{\sum V_j}{V_t} \right)$$

この指標は、各項目が他の項目とどれくらい相関するかを総合的に表した指標で、とくに**内的整合性信頼性**と呼ばれます。このようにして、尺度の安定の程度である信頼性は数値化できますが、次にお話しする妥当

*7 この点については、心理学データ解析基礎の授業でも触れていますので、あっさりとした説明になっています。よくわからない人は「基礎」の方の資料を読み直して確認してください。

*8 クロンバックのアルファ (Cronbach's alpha) とも呼ばれます。

482 性については、数値化できないものです。

483 2.4.2 妥当性

484 **妥当性 (validity)** は、信頼性をその上限とした上で、それが何を測っているのかを改めて考える指標
485 です。

486 たとえば身長を測ろうとして、体重計を使うとします。成長に応じて、身長が伸びますが、それは体重とも関
487 係がありますので、身長の伸びに応じて体重も増えていくでしょう。体重計は安定した計測器で、信頼性は十
488 分あると思いますが、体重計で身長が測れていると言えるでしょうか。相関する変数ですので、部分的に Yes
489 といえそうですが、やはり身長は身長計で測ったほうが良いでしょう。身長計の方が、身長という特徴を的確
490 に捉え、より本質に近い測定をしているからです。このように、作ったものがしっかりとその本質を捉えている
491 かどうか、これが妥当性の基本的なポイントです。

492 妥当性はですから、そもそも概念としてその測定しようとしているものが適切かどうか (**構成概念妥
493 当性 (construct validity)**) とか、その文言でちゃんと質問できているか (**内容的妥当性 (content
494 validity)**)、理屈通りその測定値が結果と変動しているか (**基準関連妥当性 (criterion validity)**) と
495 言った面から検証されます。最後の基準関連妥当性については、基準値と尺度値をつかって数量的に検証で
496 きますが、構成概念妥当性や内容的妥当性は中身の問題であったり、言葉と概念の対応であったりするの
497 で、数理モデル的アプローチができるものではありません。数値化できないか大きな問題ではない、というの
498 はもちろん逆で、数値化できないところであるからこそ、専門的な観点、幅広い視野、批判的思考でもって検
499 証していかなければなりません。

500 量的に評価する方法としては、今後説明していく**因子分析 (Factor Analysis)** によって**因子的妥当
501 性 (factorial validity)** を見る方法ですとか、理論通りの因子に分離できているかを見る**弁別的妥当性
502 (distinctive validity)**、あるいは**収束的妥当性 (convergent validity)** などがあります。これらをまと
503 めて、**検証的因子分析 (confirmatory factor analysis)** によって理論通りの分類ができているかをモ
504 デル適合度の観点から評価する方法もあります。これらは次回以降お話ししていく、テスト理論の発展系のな
505 かで考えていくものですので、どうぞ楽しみに。

506 2.5 課題

507 **■リッカート法** シグマ法でなく機械的に数字を割り振るとどのような問題が生じるか、自分で数値例を
508 作って検証してみてください。

509 **■信頼性の記述と報告** 心理学の尺度作成に関する論文^{*9}を読み、信頼性についてどのように記述されて
510 いるかを確認してみましょう。

511 **■さまざまな妥当性** 妥当性にはさまざまなものがあります。Grimm and Yarnold (2001)などを参考
512 に、妥当性について自分なりに調べてみてください。

513 **■リッカートのシグマ法** 表 2.2 のように、「かなり当てはまる」や「まったく当てはまる」など尺度の右の方
514 に丸をつける人が多かった項目があったとします。この時の尺度値をリッカートのシグマ法に則って算出して
515 みてください。

^{*9} 日本心理学会が出している「心理学研究」という学会誌では、【資料】というカテゴリーで毎回のよう新しい尺度が報告されてい
ます。

表 2.2 天井効果の出た尺度

反応カテゴリ	まったく当てはまらない	あまり当てはまらない	やや当てはまらない	どちらとも言えない	やや当てはまる	かなり当てはまる	まったく当てはまる
出現度数	1.00	3.00	5.00	18.00	24.00	53.00	56.00
相対度数	0.01	0.02	0.03	0.11	0.15	0.33	0.35
累積相対度数	0.01	0.03	0.06	0.17	0.32	0.65	1.00
累積相対度数の確率点	-2.50	-1.96	-1.59	-0.96	-0.47	0.39	∞
累積相対度数の確率密度	0.02	0.06	0.11	0.25	0.36	0.37	0.00

第3章

テスト理論と因子分析

3.1 古典的テスト理論

前回の信頼性についての議論のなかで、古典的テスト理論について触れました。古典的テスト理論のモデルは $X = t + e$ で表されます。すなわち、テストの点数 X は真のスコア t と誤差 e に分割できるというものです。非常に単純なモデルですが、測定したものには誤差がついているという考え方が、言い換えると目に見えるものだけが真実ではないという考え方がしめされています。この考え方は、ソフトサイエンスの領域においては重要なことです。

またこの古典的テスト理論から、いくつかの重要な考えを読み取ることができます。ひとつは誤差についての考え方です。このモデルを $X_j = t_j + e_j$ のように、ある個人の変化しない属性について、 j 回測定したとします。この時、測定の平均は次のように計算できます*¹。

$$\begin{aligned}\bar{X} &= \frac{1}{n} \sum_{j=1}^n X_j \\ &= \frac{1}{n} \sum_{i=j}^n (t_j + e_j) && \text{定義より} \\ &= \frac{1}{n} \sum_{j=1}^n t_j + \frac{1}{n} \sum_{j=1}^n e_j && \text{分配して} \\ &= \bar{t} + \bar{e}\end{aligned}$$

さて古典的テスト理論では、誤差に関して次のことが仮定されます。

- 誤差の平均はゼロ。つまり誤差が出現するときは、「正に偏る」「負に偏る」といった一貫した傾向がないと考える。
- 真のスコアと誤差との相関はゼロ。つまり誤差は真のスコアに関係なく影響してくるもので、真のスコアと共変動するようであれば偶然によるものとはいえない。
- 異なる測定誤差間の相関はゼロ。誤差同士がなにか意味のある変動をしているのであれば、それはもう偶然によるものとはいえない。

これらはいずれも、誤差が「偶然によって現れる影響で、制御不可能なもの」という考え方からは自然

*¹ この式は、ある測定を多くの個人 i に対して行ったものとして、 $X_i = t_i + e_i$ と考えることもできますが、添字が異なるだけで式の展開に違いはありません。

535 な仮定だといえるでしょう。より詳しくいえば、誤差は測定に応じて毎回一定の傾向で生じる**系統誤差**
 536 (**systematic error**) と、全く傾向のつかめない**偶然誤差 (random error)** に分けて考えられますが、
 537 系統誤差は測定に際して工夫して取り除くべき問題であり、ここでは全くの偶然による誤差についての議論
 538 からです。

539 さて誤差の平均がゼロ、つまり $\bar{e} = 0$ ですから、 $\bar{X} = \bar{t}$ となって、いつかは誤差がなくなって真のスコアを
 540 得ることができるようになる、ということが示されます。

541 また平均は 0 ですが分散はゼロではありません*2。このテストの分散を考えると、次のようなことがわかり
 542 ます。

$$\begin{aligned}
 \text{Var}(X) &= \frac{1}{n} \sum_{i=j}^n (X_j - \bar{X})^2 && \text{定義より} \\
 &= \frac{1}{n} \sum_{j=1}^n ((t_j + e_j) - (\bar{t} + \bar{e}))^2 && X \text{ をテスト理論のモデルに} \\
 &= \frac{1}{n} \sum_{j=1}^n ((t_j - \bar{t}) + (e_j - \bar{e}))^2 && \text{同じ記号でまとめる} \\
 &= \frac{1}{n} \sum_{j=1}^n ((t_j - \bar{t})^2 + 2(t_j - \bar{t})(e_j - \bar{e}) + (e_j - \bar{e})^2) && \text{展開する} \\
 &= \frac{1}{n} \sum_{i=1}^n (t_j - \bar{t})^2 + \frac{1}{n} \sum_{j=1}^n 2(t_j - \bar{t})(e_j - \bar{e}) + \frac{1}{n} \sum_{j=1}^n (e_j - \bar{e})^2 && \sum \text{ を分配} \\
 &= \text{Var}(t) + 2\text{Cov}(te) + \text{Var}(e)
 \end{aligned}$$

543 ここで Cov とは共分散を表しています。第二項の $2\text{Cov}(te)$ は真のスコアと誤差との共分散 (を 2 倍したも
 544 の) を表していることになりましたが、共分散 (相関) がそもそも線形関係を表す指標であったことを思い出し
 545 てください。相関係数は共分散を標準化したものだったわけですが、そういう意味ではこの $\text{Cov}(te)$ という
 546 のは真のスコアと誤差との相関関係を表しているようなものです。さて、ここでも誤差の仮定から、相関はゼ
 547 ロです。すなわち、誤差とはどのような傾向もなく出現するもの、という考えられているのです。どのような傾
 548 向もないわけですから、当然なにかと相関関係にあるはずがない、すなわち $\text{Cov}(te) = 0$ であるとなります。

549 これを踏まえると、 $\text{Var}(X) = \text{Var}(t) + \text{Var}(e)$ のように、テストの分散が真のスコアの分散と誤差の
 550 分散に完全に分割されました。ここから、信頼性の定義は $\frac{\text{Var}(t)}{\text{Var}(X)} = \frac{\text{Var}(t)}{\text{Var}(t) + \text{Var}(e)}$ と表すことができ
 551 るのです。言葉で言えば、**信頼性**の定義は「全分散中にしめる真のスコアの分散の割合」ということになり
 552 ます。

553 ここまでは古典的テスト理論から示されることであり、これまでの復習になります。ここから、このテスト理論
 554 をより展開させていくことを考えましょう。

555 3.2 因子分析モデル

556 3.2.1 単因子モデル

557 今からお話するのは、**因子分析 (Factor Analysis)** というモデルです。因子分析モデルは古典的テス
 558 ト理論の拡張であり、もっとも簡単な 1 因子モデルは次のように表されます。

*2 ガウスの考えた誤差論から、誤差は確率**正規分布 (Gaussian Curve)** に従うと考えられます。

$$z_{ij} = a_j f_i + e_{ij}$$

559 ここで Z_{ij} は個人 i の項目 j に対する反応を標準得点で表したもの*3, a_j は項目 j の**因子負荷量**
 560 (**factor loading**), f_i は個人 i の**因子得点 (factor score)**, e_{ij} は個人 i と項目 j の組み合わさった時
 561 に生じた誤差と呼ばれます。

562 **因子負荷量 (factor loading)** とは, 因子というこのテストで測定したい特性と, 項目との関係の強さを
 563 表しているものです。**因子得点 (factor score)** とは, 因子というこのテストで測定したい特性と, 個人との
 564 関係の強さを表しているもので, その人のスコアだということができます。

565 記号について添字に注目してください。添字 i は個人を, 添字 j は項目を表していますが, 因子負荷
 566 量は a_j と表されています。つまり項目によって変わる変量です。因子得点は f_i と表されています。つまり人
 567 によって変わる変量です。古典的テスト理論をこの添字を使って表現するならば, $X_{ij} = t_i + e_{ij}$ となります
 568 が, これと比べてみると t_i が $a_j f_i$ に変わったのが因子分析モデルだということになります。 t_i は個人につい
 569 ての真のスコアなのですが, 古典的テスト理論の場合, テストの点数は個人のこの能力だけを反映していると
 570 考えられていたこととなります。もしテストの問題が難しすぎて, まったく答えることができなければ, その人の
 571 能力はゼロということになるわけです。しかし中には悪い問題というものもあって, たとえば小学生に高校で習
 572 う知識が必要な問題を解かせるような問題があれば, 誰だって解けないかもしれません。解けない問題を出
 573 して「学力が低い」と結論づけるのはやや暴力的ですらありますね。このように, 古典的テスト理論は項目の良
 574 し悪しといったものが評価できないモデルだったのです。

575 因子分析モデルはこれを改良し, $a_j f_i$ としました。すなわち, ある項目に対する反応 z_{ij} は, その項目が測
 576 定したい特徴を十分に反映しているかどうか (a_j) と, その人がその特徴を有しているかどうか (f_i) の両方
 577 が成立している必要があるわけです。掛け算ですから, 一方がゼロであれば結果もゼロになります。すなわち
 578 測定したい特徴を反映していない項目 ($a_j = 0$) であれば, どれほどそれについての能力 (f_i) が高くても反
 579 応できないのです。たとえば美的センスに非常に秀でた人がいても, 数学のテストでその能力を反映させるこ
 580 とはできませんよね。これは数学のテストというのが数学力を測定するものであって, 美的センスを測定する
 581 ものではないからです。

582 因子分析は知能検査や性格検査の文脈から生まれてきたものです。心理学において「知能」とは, 何にで
 583 も応用可能な一般的な知能というのがあるのか, あるいは語彙力や計算力といった複数の個別の能力があ
 584 るのか, という議論がありました。知能検査としていろいろなものが考えられますが, それらがきちんと当該能
 585 力を測定する検査法だったかどうかはわからないわけで, 因子分析モデルはそこを評価できるようにした, と
 586 も言えます。性格検査についても同様で, 特性論的に考えるならば人間の性質というのは複数のもの, たとえ
 587 ば外向性, 神経症傾向, 開放性, 協調性, 勤勉性*4などがあり, ある性格検査の項目は協調性を測定するの
 588 には向いているけれども, 神経症傾向を測定するには向いていないということがあるわけです。このように,
 589 心理学と因子分析モデルは深い関係があります。

590 3.2.2 多因子モデル

591 さて先ほどは一因子, あるいは単因子ともいいますが, 測定したい特徴が1つだけのモデルでした。学力
 592 テストなどは一因子で問題ありません。国語のテストは国語の能力を, 数学のテストは数学の能力を測定して

*3 標準得点 (Standard Score) とは, 素点 X_j を $Z_j = \frac{X_j - \bar{X}}{\sigma_j}$ と変換したものであることを思い出してください。標準化さ
 れたスコアは平均が0, 分散が1になりますので, 単位の異なるもの同士であっても標準得点を使うと相互に比較可能になるの
 でした。

*4 小塩 (2020) のビッグファイブについての解説に基づいています。

593 いれば良いのであって、数学のテストを解くのに美的センス (真美的能力) が必要というのは、むしろ困った
 594 状況ですらありますね。しかし、知能検査や性格検査の場合はそうではありません。ある行動傾向, ある形容
 595 詞, ある検査がたった 1 つの能力・性質・心理的要因だけを反映しているとは限りません。たとえば人に優しく
 596 振る舞うといっても、その背後に外向性があるのか, あるいはそうすると自分がよく見られるからという利己
 597 的な性格があるのか等々が考えられます。1 つの項目に複数の要素 (因子) が複合的に影響していることを
 598 考えるべきです。そこで因子の数が 1 つではなく、複数ある多因子モデルを考えることにします。多因子モデ
 599 ルは次のように表現されます。

$$z_{ij} = a_{j1}f_{i1} + a_{j2}f_{i2} + a_{j3}f_{i3} + \cdots + a_{jm}f_{im} + d_j u_{ij} \quad (3.1)$$

600 記号の意味は単因子モデルと基本的には同じです。 z_{ij} は個人 i の項目 j に対する反応を標準得点で表
 601 したものであり, a_{jm} は第 m 因子の**因子負荷量**, f_{im} は第 m 因子の**因子得点**を表しています。因子の数
 602 が複数あるモデルですから, $a_{j1}, a_{j2}, \dots, a_{jm}$ とか $f_{i1}, f_{i2}, \dots, f_{im}$ のように因子の番号と項目・個人の
 603 添字の組み合わせになっていることを確認してください。最後の e_{ij} が $d_j u_{ij}$ となっていますが, これは誤差
 604 についても他の因子と形を同じくし, 項目に依存するものとそれ以外に区別しているだけです。

605 さて, ここでは因子を m 個あるとしています。この因子はどの項目にも共通して働くので, **共通因子**
 606 (**common factor**) と呼ばれます。共通因子がいくつあるかは事前にわかりませんが, 一般的にその数は
 607 数個～十数個になります。性格心理学は長い研究の中で, 性格を表す言葉に共通する因子はおおよそ 5 つ
 608 ぐらいであろう, という答えを得るに至りましたが, それ以外の領域では領域ごとの見解があるでしょう。知能
 609 が何種類の因子に分かれるのか, あるいはとある心の状態がどういう構造をしているのか, というのは心理
 610 学的にみても十分興味のある考え方です。もちろん因子分析によって得られる因子が, 人間の潜在的な知能
 611 や概念に直接対応しているとは言えないのですが*5, それでも因子がどのような構造 (しくみ) をしているの
 612 かについての一定の情報を与えてくれます。多因子モデルが心理学一般で広まったのは, こうした心の「構
 613 造」に注目する学問との相性が良かったということでしょう。

614 3.3 因子分析の定理

615 3.3.1 因子分析モデルの展開

616 因子分析モデルも古典的テスト理論のように, 式の展開から何が見えてくるか考えてみましょう*6。

617 左辺の z_{ij} は観測されたデータから算出できるものですが, 右辺の因子負荷量, 因子得点はいずれも未知
 618 数です。データに対して未知数が多すぎるようで, これではどのようにして答えを見つけ出せば良いのかわか
 619 らないかもしれません。たとえばある人のある項目に対する標準得点が 0.12 であるとして, それが 0.4×0.3
 620 で得られるのか, 0.2×0.6 で得られるのか, はたまた他の数値の組み合わせで得られるのか, を解く数学的
 621 技術は存在しません。この方程式はこのままでは解けないのです。

622 そこで, この未知数だらけの方程式を解くために, 因子について以下のような条件を置きます。

- 623 • 共通因子の因子得点, 独自因子の因子得点は, 標準化されている。すなわち, いずれの因子得点も平
 624 均点は 0 であり, 分散は 1 である。
- 625 • 独自因子は共通因子, 他の独自因子と相関しない。

*5 むしろテスト項目や調査票などに対する反応パターンが因子として出てくるだけで, 性格や知能が数次元あるというより, 我々は性格や知能を数次元で捉えることしかできない, という言い方の方が正しいでしょう。

*6 以下このセクションは小杉 (2018) の原稿を再構成したものです。

626 この他に、状況に応じて因子同士の間に関連を仮定します。

- 627 • 共通因子同士の相関を認めないのを「直交因子モデル」、認めるのを「斜交因子モデル」と呼ぶ。

628 このような仮定を置いたら問題が解けるようになるのでしょうか？ 実はこの問題を解く鍵は、多変量デー
629 タであればなんとかなるのです！

630 2つの変数、 j と k の標準得点から、

$$r_{jk} = \frac{1}{N} \sum_{i=1}^N z_{ij} z_{ik} \tag{3.2}$$

631 のように、相関係数が算出されることを思い出してください。先ほどの因子分析の基本代数式 (式 3.1) を
632 この式に代入してみましょう。

$$\begin{aligned} r_{jk} &= \frac{1}{N} \sum_{i=1}^N z_{ij} z_{ik} \\ &= \frac{1}{N} \sum_{i=1}^N (a_{j1}f_{i1} + a_{j2}f_{i2} + \dots + a_{jm}f_{im} + d_j u_{ij})(a_{k1}f_{i1} + a_{k2}f_{i2} + \dots + a_{km}f_{im} + d_k u_{ik}) \end{aligned} \tag{3.3}$$

633 これは代数の計算としてやっていくと、非常に煩雑で間違いが起きやすそうです。そこで、少し視覚化して
634 わかりやすくしてみましょう。多項式の掛け算は、各項目の総当たり戦ですので、列方向に z_{ij} 、行方向に z_{ik}
635 の各要素を置いて、要素同士の組み合わせ表を作ります (図 3.1)。

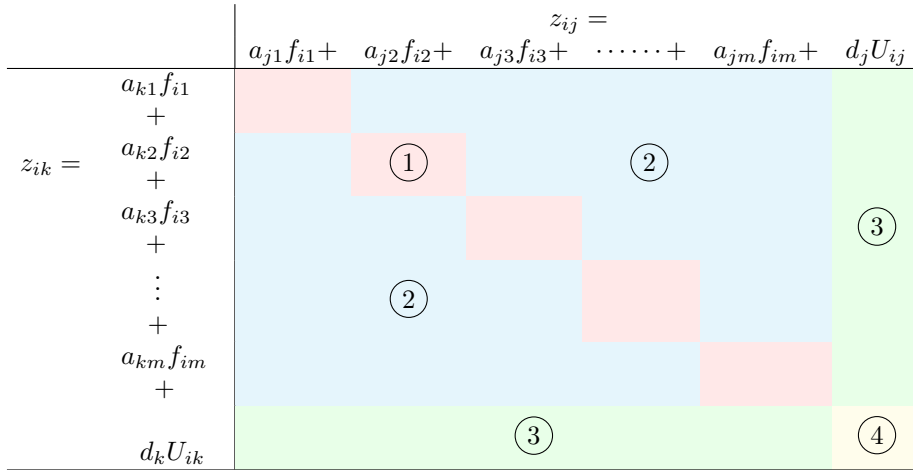


図 3.1 項目同士の総当たりを考える

636 図 3.1 に示されたのは個人 i についてのものであり、これが人数分ある、すなわち $\sum_{i=1}^N$ をつけないといけ
637 ないことに注意してください。さて図を軽く色分けしてあるのですが、ここにあるように計算すべき領域を 4 つ
638 に分けて考えていきましょう。

- 639 ①の領域 因子 p と p の積和部分 (同じ因子同士の掛け合わせ)
- 640 ②の領域 因子 p と q の積和部分 (異なる因子同士の掛け合わせ)
- 641 ③の領域 因子 $p(q)$ と独自因子の積和部分
- 642 ④の領域 独自因子同士の積和部分

643 この各パートを順に計算していきましょう。まず①ですが、たとえば第一因子については

$$\frac{1}{N} \sum_{i=1}^N a_{j1} a_{k1} F_{i1}^2 \quad (3.4)$$

644 となるのがわかります。ここで、 a_{j1} と a_{k1} は N には関係がない (\sum は i が 1 から N まで変化すること
645 を意味しているが、係数 i はどちらにも入っていない) ので、総和して割る意味がないことに気づきます。そう
646 になると、必然的にこの式は、

$$\frac{1}{N} \sum_{i=1}^N a_{j1} a_{k1} F_{i1}^2 = a_{j1} a_{k1} \frac{1}{N} \sum F_{i1}^2 \quad (3.5)$$

647 となります。この F_{i1} は因子得点であり、上の仮定より標準化されたものということになります。さて、標準
648 得点と標準得点の積和平均は相関係数になることをもう一度思い出してください！そうすると、これは自分
649 自身との相関係数を表していることとなりますから、当然 $F_{i1}^2 = 1.0$ であることがわかります。

650 結局、①のエリアは

$$\frac{1}{N} \sum_{i=1}^N a_{j1} a_{k1} F_{i1}^2 = a_{j1} a_{k1} \frac{1}{N} \sum F_{i1}^2 = a_{j1} a_{k1} \quad (3.6)$$

651 と、とてもあっさり書き下すことができます。

652 つづいて②を見てみましょう。ここは異なる因子がかけ合わさった部分ですね。落ち着いて、第一因子と第
653 二因子を例にして考えてみましょう。この箇所で見られるのは、

$$\frac{1}{N} \sum a_{j1} F_{i1} a_{k2} F_{i2} = a_{j1} a_{k2} \frac{1}{N} \sum F_{i1} F_{i2} \quad (3.7)$$

654 ということになります。ここで、 F_{i1} および F_{i2} はそれぞれ第一、第二因子における個人 i の因子得点を意味
655 しています。因子得点は標準化されていることをもう一度思い出すと、これは第一因子と第二因子の相関係
656 数になります。ここで、この因子分析が直交因子モデルだと考えますと、因子同士に相関がないわけですから、
657 数字としては 0.0 で消えてしまいます。するとこの部分は、

$$\frac{1}{N} \sum a_{j1} F_{i1} a_{k2} F_{i2} = a_{j1} a_{k2} \frac{1}{N} \sum F_{i1} F_{i2} = 0 \quad (3.8)$$

658 となるのがわかりました。つまり、この領域②は、すべて 0 になってしまうのです。

659 続いて③の部分について考えてみましょう。これはある共通因子と独自因子の積和部分です。例によって
660 標準得点同士の関係から、相関係数を算出することになりますが、独自因子は共通因子と無相関であること
661 を考えると、

$$\frac{1}{N} \sum a_{i1} F_{i1} d_j U_{ij} = a_{ik} d_j \frac{1}{N} \sum U_{ij} F_{i1} = 0 \quad (3.9)$$

662 とこのように、この箇所もすべて 0 になってしまいます。

663 最後の④に至っては、独自因子と独自因子の積和ですから、これも

$$\frac{1}{N} \sum d_j d_k U_{ij} U_{ik} = d_j d_k \frac{1}{N} \sum U_{ij} U_{ik} = 0 \quad (3.10)$$

664 のように 0 になります。結局、消えて無くなるのがほとんどで、残るのは①の部分だけであり、 r_{jk} を考えると
665 きはそこだけ考慮すればよいこととなります。

666 整理すると、

$$r_{jk} = a_{j1} a_{k1} + a_{j2} a_{k2} + \cdots + a_{jm} a_{km} \quad (3.11)$$

667 ということがわかります。つまり、項目 j と項目 k の相関係数は、項目 j の因子負荷量と項目 k の因子負荷
 668 量を、すべての因子について総和したものであるということです。因子分析の基本モデルから導出されるこの
 669 定理を、とくに**因子分析の第二定理**と呼びます。

670 ここで同じ項目同士の相関を考えてみましょう。項目 j と項目 j の相関係数は、もちろん 1.0 になります
 671 ね。これを因子分析の基本式で表すと、次のように表現できます。

$$r_{jj} = a_{j1}^2 + a_{j2}^2 + \cdots + a_{jm}^2 + d_j^2 = 1.0 \quad (3.12)$$

672 さて、この式が意味するのはなんでしょう。意味を考えてみると、ある項目それ自身との相関係数は、因
 673 子負荷の二乗和からなっている、ということがわかります。これこそ**因子分析の第一定理**と呼ばれるものであ
 674 り、解けるはずのなかった方程式を解くための鍵となる式なのです。

675 3.3.2 因子分析の定理

676 数式の展開はいったんここまでにして、第一定理は次のような形をしているのでした。

$$a_{j1}^2 + a_{j2}^2 + \cdots + a_{jm}^2 + d_j^2 = 1.0$$

677 ここで共通因子部分を、 $a_{j1}^2 + a_{j2}^2 + \cdots + a_{jm}^2 = h_j^2$ のようにすると、この式は単に $h_j^2 + d_j^2 = 1.0$ となり
 678 ます。この h_j^2 のことをとくに**共通性 (communality)** といいます。この式は共通性と独自因子の二乗和
 679 が 1.0 になることを意味しています。言い換えると、全体を 100% とした比率で共通性と誤差を比較できると
 680 いうことです。共通性は因子負荷量の二乗和で、共通因子はそのテストの背後にある共通の要因、すなわち
 681 テストで測定したいものだったわけです。古典的テスト理論では、モデル式の展開から信頼性を全分散中に
 682 示る真のスコアの割合と定義しましたが、因子分析モデルはこのように 1 つの項目 j における共通因子の割
 683 合を算出し、項目の信頼性を考えることができます。因子分析モデルにおける信頼性は、1 項目の中に含ま
 684 れる共通因子の大きさだとも言えるわけです。逆に $d_j^2 = 1 - h_j^2$ は**独自性 (uniqueness)** は、当該項目が
 685 そのテストで測っていないものの大きさを表しており、この割合があまりにも大きいと「この項目は全然関係
 686 ないものを測っちゃってるんじゃないか？」と疑われることとなります。多因子モデルにおいては、多角的に
 687 対象を切り分けるために多くの質問を投げかけるわけですが、独自性の高い項目は回答者に負担をかける
 688 だけの邪魔なものですから、実践上はこうした項目を除外することが少なくありません。因子分析には**単純構
 689 造の原則 (principle of simple structure)** と呼ばれるものがあり、項目は該当する因子を適切に反映
 690 し、かつ、他の因子と関係ないことが美しいとされます。尺度構成段階では、共通性 (独自性) をみて項目の
 691 良し悪しが判断されるのです。

692 次に第二定理を見てみましょう。第二定理は次のような形をしているのでした。

$$r_{jk} = a_{j1}a_{k1} + a_{j2}a_{k2} + \cdots + a_{jm}a_{km}$$

693 2 つの項目 j と k の相関係数は、それぞれの因子負荷量の積和の形で表される、というものです。ここに誤
 694 差の話は入ってこず、共通因子だけで話ができます。

695 左辺の相関係数は、2 つの項目がどれほど同じものを測定しているかの指標です。相関係数 (の絶対値)
 696 が高ければ高いほど、2 つの項目は同じものを指し示しているわけです。逆に相関係数が低いということは、
 697 2 つの項目に関係がないことを表します。ここで右辺に目をやりますと、右辺の各項目は因子負荷量の積の
 698 形になっています。左辺の値が小さくなる 1 つの理由は、ある共通因子 m が項目 j, k に対して、異なる方向
 699 で寄与しているからだと考えられるでしょう。そしてそのパターンが一貫していないという状況です。そもそも
 700 相関係数が小さいところからは因子を見つけ出すのは難しいのですが、そうした状況があるのはある項目ペ
 701 アについて因子同士の向きがバラバラに影響しているからだとと言えます。そのような状況は、測定がきちんと

702 できているかどうか怪しいですね。**測定の一義性**とも言われますが、そのような尺度は妥当性が低いといえ
703 るでしょう。

704 このように、因子分析モデルは第一定理で信頼性を、第二定理で妥当性をあらわすものになっているの
705 です。

706 3.4 課題

707 ■**テスト理論と因子分析** 因子分析モデルは古典的テスト理論をどのように発展させたのか、添字に注意
708 しながら数式で表現してみよう。

709 ■**因子分析の定理の導出** 因子分析の定理の導出を自分でできるようになろう。

710 ■**因子分析の定理の意味** 因子分析の定理が何を意味しているのか、自分なりの言葉で説明してみよう。

711 第4章

712 現代テスト理論

713 4.1 因子分析とテスト理論

714 ここまで、古典的テスト理論と因子分析の話を見てきました。古典的テスト理論から、テストの信頼性と妥当
715 性の話を導きました。次に因子分析モデルによって、古典的テスト理論が多因子 (多次元) モデルに展開され
716 るのでした。

717 テストの理論も心理学の研究も、目に見えない「学力」や「性格」といったものを測定するという意味で、
718 ツールとしては同じものを使うわけです。一方ではテスト、他方では質問紙とか尺度と呼ばれますが、狙いは
719 回答者の反応パターンから潜在的な性質を見出そうとするものです。ここで、改めてテストの理論に戻りたい
720 と思います。ただしこれまで古典的テスト理論と呼ばれていたものは、その名の通り古典的であって、現代的
721 なテスト理論はどうなっているのか、というところを考えてみたいと思います。

722 現代的なテスト理論、新しいテスト理論ともいわれますが、それは**項目反応理論 (Item Response**
723 **Theory)**、あるいは項目応答理論とよばれます。略して **IRT** と表現されることも多いですね。この理論はい
724 わゆる「学力テスト」などの要請から展開してきたものです。心理学的なアプローチからは、従属変数が連続
725 的であったり^{*1}、心理的な構造が知りたいために多因子であったりするのが、自然な発想でした。これに対し
726 て学力テストのようなものは、各項目の結果が「正答/誤答」の二種類しかありません^{*2}。数値としては 1/0 の
727 二値、バイナリデータであり、尺度水準は名義になります。また、測定したいものは一因子です。国語のテスト
728 は国語の能力を、算数のテストは算数の能力を測定するべきだと言えるからです。このように、項目反応理論
729 は因子分析の特殊系だということが出来ます。

730 受講生のみなさんは心理学での応用例の方が興味があるかと思いますが、後ほどこの項目反応理論のモ
731 デルが展開し、再び因子分析モデルに戻ってきますのでお楽しみに。それまではひとまず、テストの項目を分
732 析するというのはどのようなことがなされているのかをみていきたいと思います。

733 4.2 通過率と累積正規分布

734 みなさんは大学入学共通テスト (旧センター試験、さらにその前は共通一次試験と言いました) や、学内の
735 定期テスト、模試など色々なシーンでテストを受けてきたことと思います。模試などでは偏差値が明らかにな
736 り、自分の実力が相対的にどのあたりにあるのかがわかるようになっていたかと思います。大学共通テストな
737 どは 50 万人ぐらいが一度に受験しますから、さまざまな学力の人がそこには含まれるのですが、成績を図に

*1 因子分析モデルの式 3.1 が Z_{ij} から始まっていたことを思い出してください。標準化されているということは、平均や標準偏差が求められているということであり、**間隔尺度水準**以上の数字が前提とされています。

*2 部分点というのがあるじゃないか、と思うかもしれませんが、それはひとまず横に置いてください。

738 するととても綺麗な正規分布になることが知られています*3。正規分布は誤差の分布でもあります。多くの
739 要因が考えられる際の集積的データも、自然とこの形になることがわかります。

740 さて、学力のような潜在変数が標準正規分布に従うと仮定しましょう。この分布の形はどこの確率点がどれ
741 ぐらいの確率密度を持っているか、あるいはある確率点以上・以下の面積が全体の何 % を表すものです
742 が、縦軸をある点以下の累積確率に書き直してみましょう (図 4.1)。図 4.1 の上の図がいわゆる正規分布の
743 分布の形、確率密度関数です。下の図はこれを累積確率に書き換えたものになっています。累積確率は 0%
から始まって、最終的に 100% にまでどのように増えていくかを示した図になります。

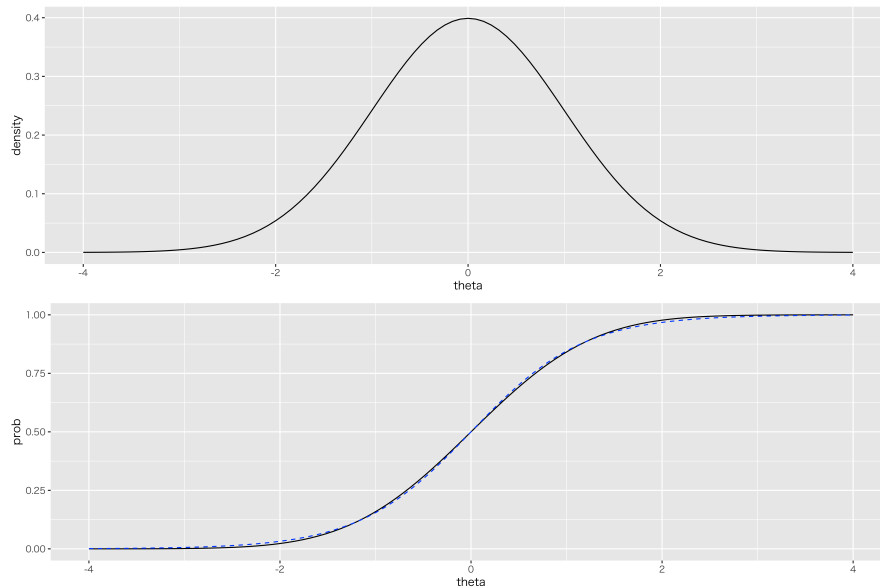


図 4.1 正規分布の確率密度関数 (上) と累積確率関数 (下)

744

745 累積正規確率は、テスト理論と密接な関係があります。というのも、学力が正規分布すると考えるなら、累
746 積正規分布の形はあるテスト項目の**通過率 (pass ratio)**と同じ形になると考えられるからです。

747 通過率とは、あるテスト項目に正答する人の割合のことです。ここで複数の項目からなる、あるテストをした
748 としましょう。正答数を数えてその人の成績とすると、よくできたテストであれば成績は正規分布に従います。
749 さらに、ある項目と成績との相関 (**IT 相関 (Item-Total correlations)**) は高いはずですね。つまりその
750 項目に正答することが、テスト全体の成績と高く関係しているのです、その項目はテスト全体が測ろうとしている
751 ものを反映していると考えられるからです。また、成績をもとに被験者集団を 5 群に分けたとしましょう。「成
752 績上位群 (HH)」、「成績やや上位 (MH)」、「成績中程度 (M)」、「成績やや下位 (ML)」、「成績下位 (LL)」
753 です。このとき、各群の平均通過率を考えると、図 4.2 左上図のようになるのが理想的です。つまり、成績が
754 高い人たちの通過率は高く、成績が低い人達の通過率は低くなるはずですね。同じ図の右上は、LL 群でも
755 半分ぐらいが通過し、その後の群は過半数、ほとんどが通過するようになっています。これは、この項目が簡
756 単すぎたことを意味しています。簡単すぎる問題は、それはそれで被験者の特徴が弁別できないという意味
757 で悪い項目です。逆に図の左下にあるのは、HH 群でも半分以下の通過率しかありません。つまり難しすぎ
758 る問題です。ほとんどの人が間違えてしまうわけですから、これも良い試験問題とは言えないでしょう。右下
759 に至っては逆転していて、どうやったらこういう項目が作れるのか却ってわからないほどですが、学力の低い

*3 山内 (2010) の見返し (表紙を開いた最初の内側のページ) に、センター試験の成績分布が載っており、綺麗な正規分布であることが示されています。

760 人だけが正答できて、学力の高い人は正答しない項目ということになります。もちろんこんな項目はよくない
 761 わけで、IT 関連で負の相関が出ているわけですから、テストの文脈でいうなら「そのテストで測っていない
 762 何か別の能力を測っている」と考えるしかありません。

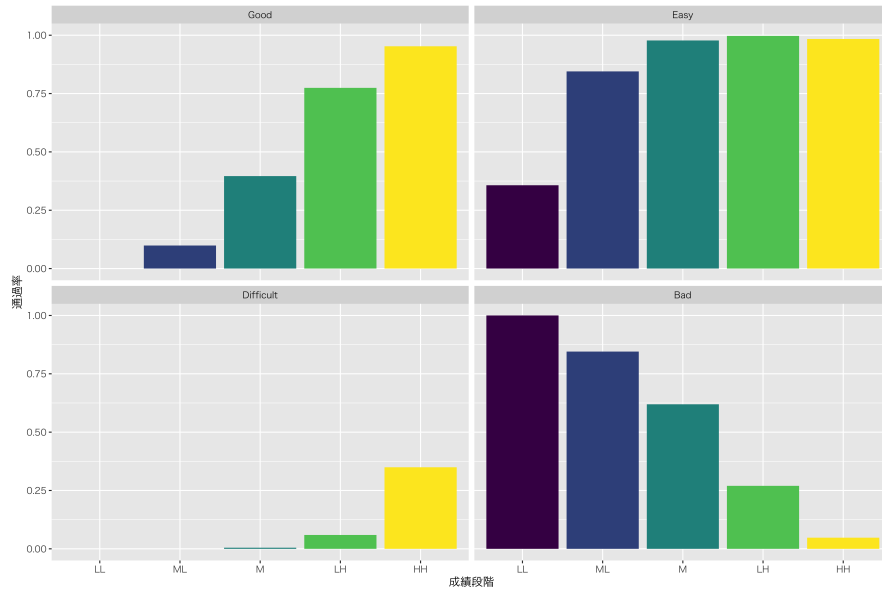


図 4.2 群ごとの平均通過率。左上が良いパターン。右上は簡単すぎる、左下は難しすぎる項目。右下は逆転していて良くない項目。

763 ともあれ、このようなやり方で項目の良し悪しを見ていくことができます。また、図 4.2 は 5 段階ですが、7
 764 段階、9 段階とよりきめ細かくしていくと、理想的な形は累積正規分布により近づいていきます。新しいテスト
 765 理論による項目分析はこの累積正規分布の形を基本とし、これを拡張することで各項目の特徴を描いていく
 766 ことになります。

767 ところで 1 つ前の図 4.1 の下の図には、実線と点線の 2 つの線が絡んでいることにお気づきでしょうか。実
 768 線の方は確率分布関数から累積確率を出して描いたものですが*4、点線のほうは次の関数を使って描いて
 769 います*5。

$$f(x) = \frac{1}{1 + \exp(-1.7x)}$$

770 この関数、図から明らかなように累積正規分布とほとんど同じですよね。累積正規分布の関数を直接使う
 771 と、積分計算 (\int を使うやつ) が入ってくるのでちょっと計算が面倒ですから、こちらの関数の方を近似関数
 772 として用います。この関数のことを**ロジスティック関数 (logistic function)** と言います。ロジスティック関
 773 数そのものは、先の式から 1.7 という係数を除いた $\frac{1}{1 + \exp(-x)}$ で表されるもので、 $-\infty$ から $+\infty$ まで
 774 のどんな数字が与えられても、答えを 0 から 1 の範囲に変換してしまう関数です。この特徴はとても便利で
 775 す。というのも、結果が 0 と 1 の間に入るといことは、比率を表していると考えられるからです。0/1 のバイ
 776 ナリデータが従属変数のときに、独立変数をこのロジスティック関数で変換してやれば 0 か 1 のどちらに近い
 777 か、どれぐらいの比率で 1 の目が出るかがわかります。項目反応理論も結果が 0/1 (誤答/正答) ですから、

*4 R では `pnorm` 関数を使って描きます。数式でいうなら、 $\int_{-\infty}^p \frac{1}{\sqrt{2\pi}} \exp(-\frac{x^2}{2})$ となります。

*5 `exp` というのは指数関数で、 $\exp(x) = e^x$ のことです。ここで e は数学定数で、 $e = 2.718282\dots$ という数字です。

778 「学力」のような目に見えない能力をロジスティック変換してやれば、結果が正答率になるというのは大変便
779 利なのですね。

780 それではこのロジスティック関数をつかった項目分析の話に進んでいきましょう。

781 4.3 項目母数の特徴

782 ロジスティック曲線が累積正規分布の近似関数になっていること、テスト項目の分析には通過率を使って
783 考えることを見てきました。とくに通過率の分析 (図 4.2) では、その項目が難しい設問だったのか、簡単なも
784 のだったのかを見ることができました。ロジスティック曲線もこの「項目の難しさ」を表現できるように、次のよ
785 うに拡張できます。

$$p(\theta) = \frac{1}{1 + \exp(-1.7(\theta - b))}$$

786 左辺の $p(\theta)$ に含まれる θ は、潜在変数のスコア、**因子得点**であり、ここでは標準化された学力ですから、
787 偏差値のようなものだと思ってください*6。 $p(\theta)$ は能力 θ の人がこの項目に正答する確率=(通過率)です。

788 ここで b という変数が入ってきました。これが**困難度 (difficulty)**を表す指標です。 $b = 0$ のときは
789 標準正規分布と同じ形になりますが、 $b = 1$ ならばこの式は右に、 $b = -1$ ならば左に動きます。つまり
790 $b > 0$ であれば難しく、 $b < 0$ であれば簡単であることを表現していることとなります。図 4.3 に困難度
が $b = -1, 0, +1$ の時の曲線を書いてみましたので確認してください。このように困難度を表現するパラ

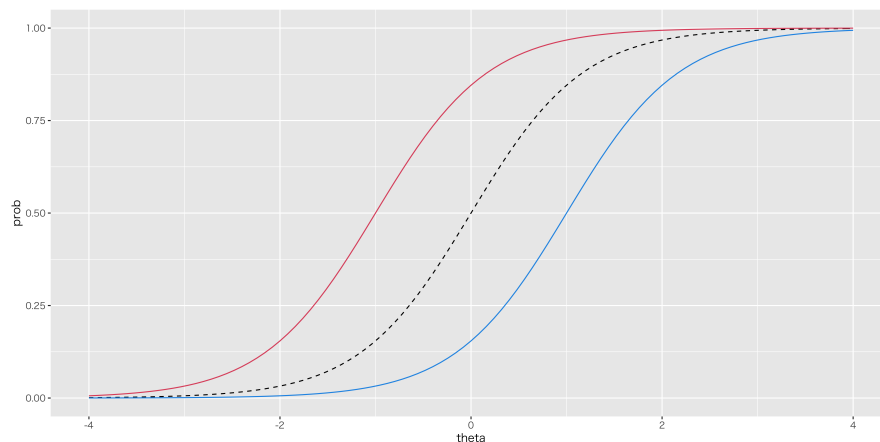


図 4.3 困難度母数の入ったロジスティック曲線。1PL ロジスティックモデルともいう。点線が $b = 0$ の標準的曲線。赤が $b = -1$ 、青が $b = +1$ の例

791
792 メータを追加したモデルを **1 パラメータ・ロジスティックモデル (One Parameter Logistic model)**
793 と言います。実際のテストの回答パターンにたいし、各項目にこの曲線を当てはめて困難度を推定するこ
794 とで、項目を評価できるようになります。このように項目の特徴を描く曲線のことを**項目特性曲線 (Item**
795 **Characteristic Curve, ICC)** と言います。

796 さらにパラメータを追加して、次のようにすると **2 パラメータ・ロジスティックモデル (Two Parameter**
797 **Logistic model)** になります。

$$p(\theta) = \frac{1}{1 + \exp(-1.7a(\theta - b))}$$

*6 偏差値は標準化スコア z_i を $10z_i + 50$ と変換したものを指します。ここはその変換前の z_i と同じです。

798 ここで a という母数 (パラメータ) が入ってきました。これは $a = 1$ だともとのモデルのままなのですが、これ
 799 が小さくなると曲線が傾き、大きくなると曲線のカーブが強くなります (図 4.4)。曲線が緩やかになると (図
 800 4.3 の赤線)、 θ の違いに対して通過率の変化が乏しくなります。言い換えると感度が悪くなるわけです。逆に
 801 曲線の立ち上がりが強くなると (図 4.4 の青線)、 θ がある一定のところを超えるかどうかで正答率がグッと
 802 変化することになります。つまりこのパラメータは、回答者の能力 θ を分類する力の強さを表しているのです。
 803 このパラメータのことをとくに**識別力 (discriminant)** と呼びます。

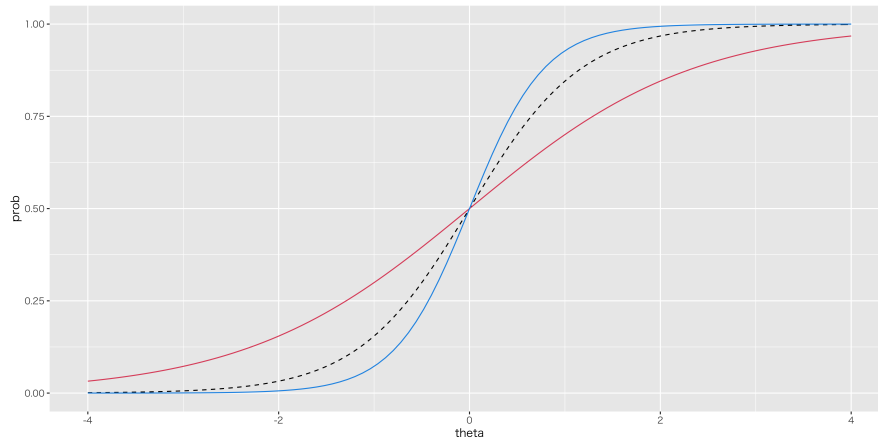


図 4.4 識別力母数の入ったロジスティック曲線。2PL ロジスティックモデルともいう。点線が $a = 1$ の標準的曲線。赤が $a = 0.5$, 青が $a = 1.5$ の例

804 一般にはここまで紹介した 2PL モデルがよく用いられます。他にも 3, 4, 5 つとパラメータが増えたモデル
 805 もありますが^{*7}, もとのロジスティック曲線に特徴を付け足していったもので、曲線をずらしたり、曲げたり
 806 しながらもとのデータにうまく当てはまるようにしつつ、その特徴を解釈できるように工夫しています。

807 いずれにせよ、テストの結果から各項目の特徴を記述します。少し例示したほうがわかりやすいでしょう。
 808 図 4.5 には心理学データ解析基礎で行った試験の結果から、2PL ロジスティックモデルを当てはめて項目
 809 分析をした例です^{*8}。同じデータの項目母数を表 4.1 にも示しました。表 4.1 の数字と図 4.5 の曲線の対応
 810 をよく確認してください。

811 たとえば、項目 I0022 は困難度が-1.92、識別力が1.30 です。困難度がマイナスですので、これはかなり
 812 簡単な問題だということになります。具体的には、偏差値 50 すなわち $\theta = 0$ の人であっても 98.58% 正解
 813 するわけですから、ほとんどの人にとって容易い問題であることがわかります。ちなみにこの問題、具体的
 814 は帰無仮説検定に関する問いで、「差がない」「偏りがない」といった仮説は何と呼ばれるか」というもので
 815 した^{*9}。

816 一方、困難度が 0 近いところの例として項目 M0605 をあげますが、これが平均的な難易度の質問になって
 817 います。困難度母数 $b = 0.0$ であれば偏差値 50, すなわち $\theta = 0$ の人が正答する確率が 50% の質問とい
 818 うことになりますが、今回の M0605 はそれより少し難しいので、偏差値 50 の人で 20.70% の割合で正答で
 819 きます。この問いについて偏差値が 70 ($\theta = 2$) であれば、92.96% の確率で正答できることになりまし、偏

*7 詳しくは <http://antlers.rd.dnc.ac.jp/~shojima/exmk/jindex.htm> を参照してください。3 つめのパラメータはあて推量母数、4 つ目は上方漸近母数、5 つ目は非対称母数と呼ばれています。

*8 心理学データ解析基礎の授業では、過去の受講生のデータとさまざまな質問をプールしたデータを貯めてあります。みなさんが受けたテストには入っていない項目かもしれませんが、これまでどこかで出題され、回答パターンが得られている実際のテストです。

*9 いうまでもないですが、答えは「帰無仮説」です。

820 差値 $30(\theta = -2)$ であれば 0.51%, つまりほとんど正答は望めないということになります。ちなみにこの問題
 821 は「重回帰分析において、標準化されたデータを使って分析をすることで、モデルの適合度を上げることがで
 822 きる」を Yes か No かで判断させるという質問でした。

823 他にも項目 I0017 は困難度が 2.17, 識別力は 0.69 です。困難度が最も高いグラフで、図の曲線が一番
 824 右にあるラインになっています。ただ識別力がやや低いので、曲線はよりなだらかになっていますね。困難度
 825 が高いので、 $\theta = 0$ でも 7.24% しか正答できません。難しい！ $\theta = 2$ で 45.10% ですから、かなり能力の
 826 高い人でも半分は間違えるような問題です。ちなみにこれは標本分散の期待値が母分散からどれぐらいずれ
 827 るのかを計算する問題でした。

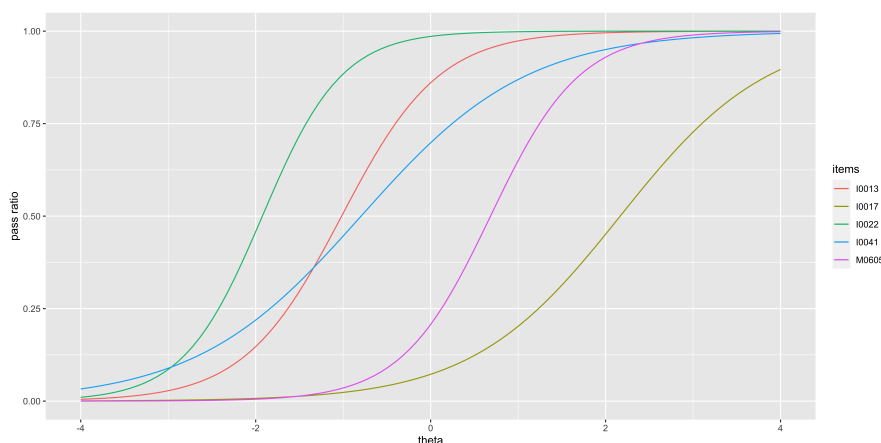


図 4.5 実際のテストに 2PL モデルを適用した例

表 4.1 各項目の困難度と識別力

	項目 ID	困難度	識別力
1	I0013	-1.02	1.05
2	I0017	2.17	0.69
3	I0022	-1.92	1.30
4	I0041	-0.79	0.62
5	M0605	0.68	1.15

828 4.4 被験者母数の推定

829 項目反応理論における因子得点の推定は、項目の特徴を表す項目母数に対して**被験者母数**と呼ばれ、上
 830 述の項目特性に基づいて行われます。先ほどの図 4.5 をもとに説明します。ある回答者が項目 I0013 に正
 831 答したとしましょう。その人の能力値 (因子得点, θ) はどの辺りにあるかといえば、確率の曲線に沿った下の領
 832 域のどこかということになります (図 4.6)。この曲線、ICC は項目の特徴を表したもので、ある θ の値の能力
 833 があればどの程度の確率で正答できるかを表した項目の特徴でもあります。逆にある人の θ がどのあたり
 834 にありそうかを示しているともいえます。たとえばこの ICC において、 θ が 0.5 のときの通過率は 86.05%
 835 ですが、言い換えればこの項目に正解した人が $\theta = 0.5$ である可能性も高そうです。 $\theta = -2$ の通過率は
 836 14.71% ぐらいですから、能力がこんなに低いとは思えませんし、1 つの項目の話でしかないですが、希望的

837 観測をするなら θ がもっと高い可能性もあるでしょう*10。

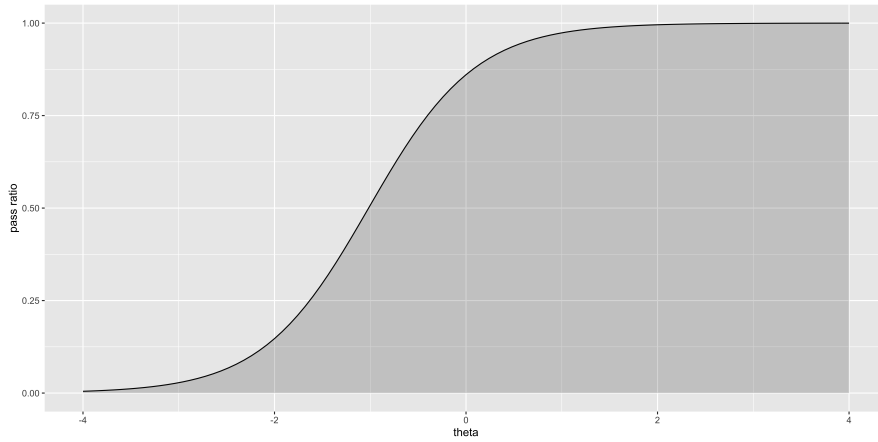


図 4.6 実際のテスト項目 I0013 に正答した人の能力がありそうな領域

838 次に、困難度のより高い項目である I0017 には誤答したとしましょう。その人は、I0017 の ICC の下の
839 領域には**ない**はずです。項目 I0013 の ICC の下で、かつ、項目 I0017 の ICC の上にあるはずなので、図
4.7 のように塗りつぶされた領域の中に入ります。

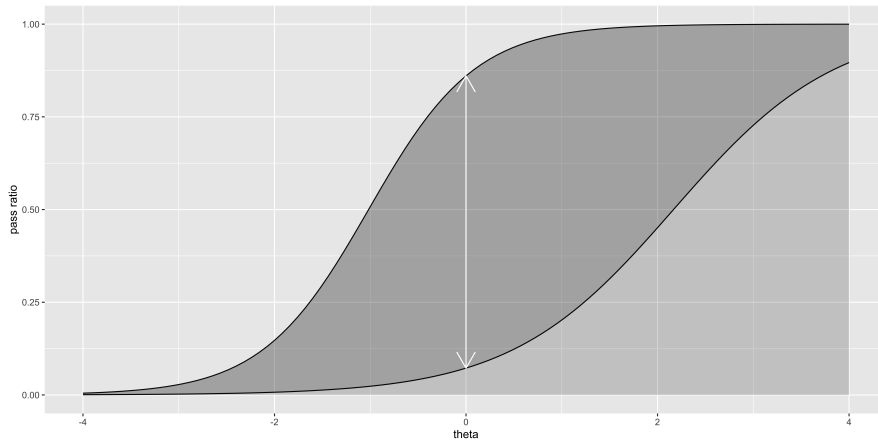


図 4.7 つぎのテスト項目 I0013 には誤答した人の能力がありそうな領域

840

841 この2つの曲線の間にある、濃く彩られた領域の幅が、回答者の能力値がありそうな程度を表しているの
842 です。図 4.7 には $\theta = 0$ の可能性の大きさを矢印で示してありますが、 θ の値はここだけに限らずこの曲線
843 の幅のどこかです。ただ $\theta = 2$ や $\theta = -2$ より $\theta = 0$ のほうが、より「ありそう」な値だということがわかり
844 ます。

845 ここでさらに同じ人に、項目 M0605 の質問をして、この人がそれにも正解したとしましょう。この人の能力 θ
846 のありそうな範囲はさらに絞り込むことができます (図 4.8)。項目 I0013 より難しい質問に正解したわけ
847 ですから、 $\theta = 0$ の可能性はグッと小さくなり、それよりも $\theta = 2$ ぐらいの方がありそうだ、ということになっ

*10 θ のありそうな「確率」とは言ってないことに注意してください。すぐにわかることですが、この ICC の下の面積を積分しても 1.0 にはなりませんのでこれは確率ではなく、尤度 (likelihood) のほうなのです！

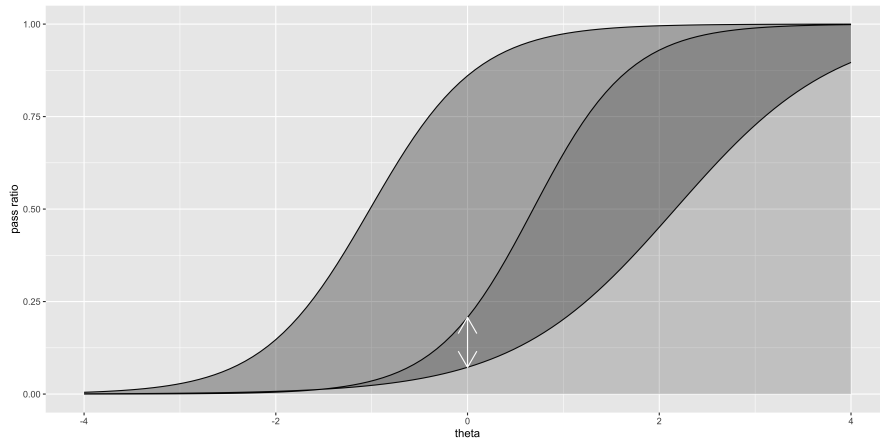


図 4.8 さらにテスト項目 M0605 には正答した人の能力がありそうな領域

848 できます。このように、1つ1つの項目からこの人の能力 θ がどのあたりにありそうか、というのを絞ってい
849 き、テストに含まれる項目全部を使えば、かなり狭い領域で「この辺りにあるはずだ」と推定できるでしょう。

850 この IRT を用いた推定方法は、このようにテストの項目ごとにその特徴を分析できるのが特徴です。テスト
851 の中でも良い項目、悪い項目というはあるでしょうが、どのように悪いのかを困難度や識別力といった項目
852 の特徴を使って表現できます。これらの数字は、テストに含まれる項目ごとの難易度を相対的に比較してい
853 く中で作られるものであり、回答者の能力に依存するものではありません。テスト理論が被験者の特徴と項目
854 の特徴を分離するところから始まったことを、改めて思い出してください。

855 また、このような項目同士の比較から求められた項目の特徴をつかって、被験者の能力（因子得点）を推
856 定する方法についても紹介しました。その過程の中で気づいたと思いますが、複数の回答を通じてある回
857 答者の能力 θ のありそうな領域が狭められていく中で、明らかにその人にとって簡単すぎる問題、難しすぎ
858 る問題は意味を成しません。たとえば図 4.8 の段階まで絞り込まれているときに、項目 I0013 よりも簡単な
859 I0022 を出題しても、おそらくほぼ確実に正答し、そのことは「領域を狭める」ことにはなんの貢献もしないで
860 しょう。回答者の能力に相応しい質問を選んで出すことができれば、とても効率よくその領域を絞り込んでい
861 くことができます。こうした考え方は、次回お話しするコンピュータ適応型テスト (Computer Adapted
862 Test) として実装されます。

863 今回はテスト理論について紹介してきましたが、この考え方は性格テストなどで用いられているリッカート
864 形式の尺度に応用できるように、発展していきます。次回はその辺りを解説していこうと思います。

865 4.5 課題

866 ■テスト理論と因子分析 項目反応理論のロジスティックモデルについて、項目母数が何を意味している
867 か、項目母数が変わると ICC がどのように変化するかを自分で説明できるようになろう。

868 ■項目反応理論の利点 項目反応理論を用いた採点方法を使うと、どういう長短所があるだろうか。次回の
869 内容に先駆けて、自分なりに考えてみよう。

第 5 章

現代テスト理論その 2

5.1 現代テスト理論の特徴

前は現代テスト理論として項目反応理論 (IRT) を紹介しました。ロジスティック曲線を応用して項目の特徴を描画し、それを使って被験者母数を推定する方法についても解説しました。この一連の手続きに基づき、現代テスト理論の利点を考えてみたいと思います。

5.1.1 現代テスト理論の利点 1: 項目母数と被験者母数の分離

古典的テスト理論からの発展として、現代テスト理論では被験者母数 θ_i と項目母数 a_j, b_j を区分して考えるようになりました。項目母数は通過率のアイデアを精緻にしたものですが、この通過率は項目群の総和を元に考えられていたことを思い出してください。すなわち、項目母数の計算には項目の相対的な困難度だけをを用いています。イメージとしては鋳物の硬度検査のようなものです。2つの異なる硬さの物をぶつけて崩れた方が負け＝より硬度が低いと考えるように、2つの異なる項目を被験者に与えて、より正答者数が少ない方がより難しいと考えるのです。これはつまり、回答者の学力が高かろうが低かろうが、困難度が $b_x < b_y$ であるという関係に違いはないという考え方です。

これまでのテストや心理尺度の作り方に比べると、この点が大きく違います。学校などのテストは教員が作っていますから、教員が自分の感覚で「こちらの方がより発展的な内容だ」「こちらの方が難しいだろう」という問いに大きな配点がなされたりするでしょう。その後テストの平均点をみて「今回のテストは簡単にすぎたか」という判断をしたりするでしょう。しかし IRT では項目それ自体に困難度を決めさせますから、そこに作成者や回答者の意図は含まれません。平均正答率が高いからと言って簡単な問題なのではなく、項目の特徴として困難度が決まるのです。

たとえばサー斯顿尺度の作り方を思い出してください (セクション 2.2, Pp.27 参照)。サー斯顿尺度では尺度適用前に評定者集団によって尺度値を決めます。この評定者集団が偏った思想の持ち主だけで固められていた場合、尺度の点数は極端なものになり、普通の人がある尺度に回答すると極めて低い点数、高い点数になってしまうかもしれません。あるいはリッカート尺度の作り方を思い出してください。(セクション 2.3, Pp.28 参照)。リッカート尺度では回答者の累積度数から尺度値を算出します。先ほど同様、回答者集団の態度に偏りがあれば、尺度の点数は標準的な物ではなくなるでしょう。つまり「誰を対象に測定するか」によって目盛りが変わるようなものです。これでは測定結果の一般化は難しいでしょう。たとえば本学で作られた尺度を、他大学でやってみると違った尺度値になるのですから、研究結果はせいぜい「その大学ではそうなんだろう」となり一般化できなくなります。従来方法は、回答者と項目の特徴が関連しすぎていたのです。

これに対し、IRT を使って項目の特徴を計算する場合は、相対的な難易度に違いはありませんから、どの

900 大学で作った尺度であっても統一的な解釈が可能です。尺度作成時に幅広くデータを集め、項目母数を確
 901 定させてしまえばどこでも統一的な評価ができます。テストなど学力を測定する際に大学間での違いが見
 902 られたとしても、その難易度を調整するのも簡単で、共通する項目を入れておけばそこを基準に相対的な困
 903 難度調整ができます*1。良問と悪問の評価と、回答者の評価を分けることは重要なポイントなのです。

904 5.1.2 現代テスト理論の利点 2: 被験者母数の推定の利点

905 被験者母数と項目母数の分離は、さらに別の利点も生み出します。それはデータが一部欠落した場合の補
 906 完に関係します。

907 リッカート尺度では回答者全体の相対頻度から、カテゴリの尺度値を決定するのです。ここである人が特
 908 定の項目にだけ回答をし忘れたとします (調査研究ではよくあることで、同じような目盛りが並んでいると一
 909 行飛ばして丸をつけちゃうようなことはよくあります)。そうすると、その項目だけ合計人数が変わりますから、
 910 計算が面倒です。また相関係数を計算する時にも、その人のその項目については計算できなくなります。因子
 911 分析は相関係数から計算を始めますから、一箇所でも欠損値があるとその人のデータを抜いてしまうか*2、
 912 他の値を代入して補完するか*3、手間でも計算の時にそこだけ外して計算するか*4といった工夫が必要で
 913 す。因子分析結果に大きな違いは出なくても、その人の因子得点は計算できないことに違いはありません。

914 それに対して、IRT の被験者母数の推定は、項目母数をつかった ICC をもとに一人ずつ絞り込んでいく
 915 というものでした。もしある人が回答していないことがあっても、計算ができなくなることはありません。その項
 916 目の情報が得られないので絞り込み精度は上がりませんが、ヒントが減っただけで回答できないわけではない
 917 のです。このように、得られた情報すべてを使ってその人の被験者母数 (因子得点) を推定する方法のこと
 918 を、**完全情報最尤推定 (Full Information Maximum Likelihood)** と言います。このように IRT を
 919 使うと、必ずしも全問に回答していなくてもスコアは計算できるということになります。欠損値があるからその
 920 人のデータは使えない、ということがないのでいいですね！

921 もっと言うと、IRT では全員が全員同じ問題に回答する必要はありません。たとえば能力値が $\theta_i = -0.3$
 922 ぐらいにありそうだ、と絞り込めている段階で、次の問題の困難度母数が $b_j = 2.5$ であれば、おそらくほぼ
 923 確実にその人は回答できないでしょう*5。その人にいくら困難度母数の高い質問を繰り返しても、ほとんど
 924 誤答がつづくだけで、とくにその人の θ がありそうな領域を狭めるヒントにはなりません。むしろ $b_j = -1$
 925 とか $b_j = -0.5$ のような簡単な問題を出して*6、それらに正答できるかどうかを見極め、絞り込んで行っ
 926 た方が効率的です。紙に印刷されているテスト (Paper Based Test) であれば、印刷された問題は変えよ
 927 うがありませんから、困難度順に問題を並べると、あるところから先はずっと不正解が続く人が続出しま
 928 す。ずっと不正解なところの問題はいくら良問でも、その人の能力を測るのには役立ちません。ヒントが
 929 増えないからです。であれば、回答者の学力に合った問題を、その都度その都度出題した方がいいですよ
 930 ね。コンピュータを使って回答者に相応しい質問をダイナミックに組み替える、コンピュータに基づいたテスト
 931 (Computer Based Test), 別名**コンピュータ適応型テスト (Computer Adaptive Test)** というのがそ
 932 れです。CAT になると、回答者ごとに問題が変わりますからカンニング対策の必要も無くなって、とても便利

*1 こうしたテスト間の困難度調整のことをテストの**等価 (equation)** と言います。

*2 リストワイズ削除と言います。

*3 欠損値補完については、平均値や中央値を代入したり、同じようなパターンで回答している人の値を使い回したり、回帰分析でその項目の値の推定値を入れたり、とさまざまな方法が考えられてきました。欠損値発生メカニズムにもよりますが、いずれもある程度バイアスのかかった値になってしまいます。統計的によりバイアスの少ない適切な代入法が考えられてはいますが、そもそも欠損値がないのが最も望ましいことに変わりはありません。

*4 ペアワイズ削除と言います。

*5 2PL モデルで $a = 1, b = 2.5$ のとき、 $\theta = -0.3$ が正解する確率は 0.8492% です。

*6 2PL モデルで $a = 1, b = -1$ のとき、 $\theta = -0.3$ が正解する確率は 76.674%, $b = -0.5$ のであれば 58.419% です。

933 になること間違いなしです。

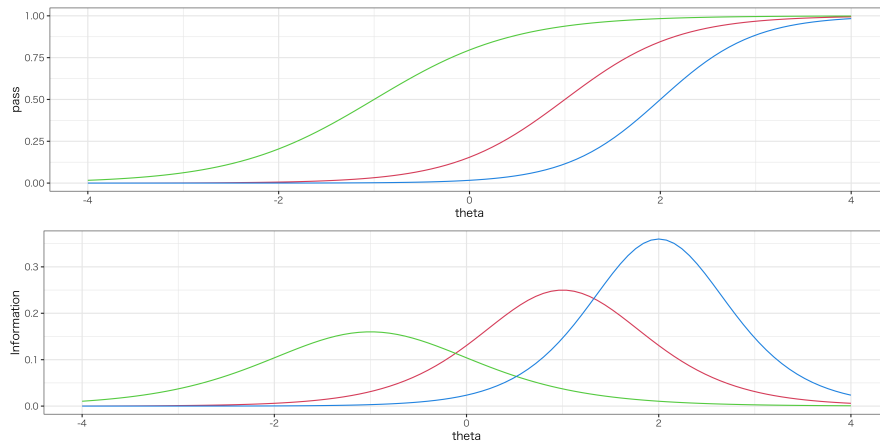
934 5.1.3 現代テスト理論の利点 3: 信頼性についての考え方

935 IRT の考え方からいうと、どんな項目でも何らかの情報を提供してくれるはずで
936 す。たとえば $b_j = 3$ のようなとても難しい項目が合ったとします。これがどれぐ
937 らい難しいかという、偏差値 70 の人でも 15% しか正解できないレベ
938 ルです。ほとんどの人にとっては誤答にしかならない難しすぎる悪問だ、と
939 言いたくなるかもしれませんが、学力が高い人がどれほど高いレベルでや
940 れるのかを検証するためには必要な問題です。偏差値が 70 なのか、75
941 なのか、80 までいけるのか、といったことを見極めるためにはこの問題
942 でない情報が必要なものなのです。

942 つまり、どの項目にもその項目が得意とする領域があるはずで
943 す。この項目はこの領域の回答者を絞り込む時に、最も有用な情報をもた
944 せしてくれるはず、という θ の場所があるはずなのです。これを表現する
945 のが項目情報曲線 (Item Information Curve) といい、次の式で表される項目
946 情報関数で描くことができます。

$$I(\theta) = a_j^2 p_j(\theta) q_j(\theta)$$

946 ここで $p_j(\theta)$ はその項目 j の θ における正答率 (通過率)、 $q_j(\theta)$ は誤
947 答率を表しています。能力が平均的、すなわち $\theta = 0$ のときは、 0.5×0.5 に
948 なる平均的な困難度 ($b = 0$) の設問が最も大きな値になる、というわけ
949 ですね。図 5.1 にいくつかの ICC とそれに対応する IIC を描きましたので、
950 確認してください。IIC の



948 図 5.1 項目特徴曲線 (ICC, 上図) と項目情報曲線 (IIC, 下図) の例。左から順に $a_1 = 0.8, b_1 = -1$
949 の識別力が弱く簡単な項目, $a_2 = 1, b_2 = 1$ のやや困難な項目, $a_3 = 1.2, b_3 = 2$ の識別力が強く困難
950 な項目。

948 ピークは、対応する ICC の $\theta = 0.5$ のところにあること、識別力はピークの尖り具合に関わっていることを確
949 認してください。

951 さて、IIC はその項目がどこで情報をもたらしてくれるか、ということを表しているのです。言い方を変え
952 ると、IIC のピークはその項目の最も信頼できるところであるとも言えます。つまり IRT において信頼性は
953 潜在特性の関数になっているのです。古典的テスト理論ではテスト全体の分散に占める真分散の割合のこと
954 を信頼性というのです。因子分析理論では項目の中の共通因子負荷量の二乗和、共通性とその項目の信

955 頼性を表しているのです。信頼性を見る水準がテスト全体から項目別の評価に発展したわけですが、IRT
956 ではさらにその項目の最も感度の良いところを探る関数として、その信頼性を評価できるようになったといえ
957 るでしょう。

958 また IIC はある項目から得られる情報のことを意味しますが、テストに含まれているすべての情報関数を
959 足し合わせることで、そのテストから得られる情報の大きさを関数として評価できます。すなわちテスト全体
960 の情報量 $I_T(\theta)$ は、 $I_T(\theta) = \sum_{j=1}^M I_j(\theta)$ であり、この関数のことを**テスト情報関数 (Test Information Curve)** といいます。図 5.1 の 3 つの項目からなる TIC を示したのが図 5.2 です。これを見ると、この 3 つ

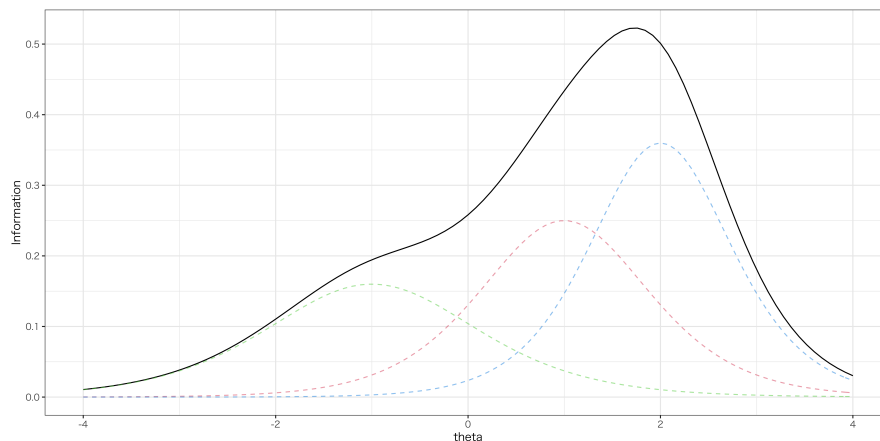


図 5.2 テスト情報関数 (黒の実線部)。理解を進めるために各 IIC を薄い点線で表現した。

961 の項目からなるテストは $\theta = 2$ より少し低いレベルを測定するとき最も鋭敏に働くということがわかります。
962 このように、項目母数がわかっていれば事前にテストの特徴をどのあたりに持ってくるかを決定でき、自由自
963 在にテストをデザインできるようになるわけです。

964 テストの例で話をしていますが、心理学的な領域でももちろん便利な手法です。たとえば高い認知能力レ
965ベルの人をとくに選出したいとか、精神的な健康度がごく低い人をしっかりと検出したいといった目的があれ
966ば、そのあたりにピークが来るような項目からなる質問項目からなる調査票を構成すれば良いのです*7。
967

968 5.1.4 現代テスト理論の問題点

969 ここまで見てきたように、IRT はさまざまな利点があります。しかし欠点がないわけではありません。

970 ここまでの話はすべて、項目母数が既にわかっているという前提付きで進めてきました。では項目母数
971 はどのようにして定めるのでしょうか？ これはもちろん得られたデータから算出できるのですが、そのためには
972 事前に多くの被験者から回答を集め、項目母数の値はほぼ間違いなくこれぐらいだろう、と言えるほど安定し
973 たものである必要があります。テストの場合、回答が 0/1 というバイナリデータで得られますから、そもそもそ
974 れほど情報がある反応ではありません。バイナリデータからその項目の特徴を安定して推定するためには、
975 かなり多くの被験者 (数千から数万単位) を集めて項目に回答させておく必要があります。テスト項目は一度
976 使ってみないと、項目母数がどうなるかわからないというのもポイントです。

977 また、CAT など項目をダイナミックに組み合わせるためには、選べるぐらさまざまな項目を準備しておか

*7 具体例として小杉 (2014) をあげておきます。学校適応感を測定するため、テストのピークがやや低いところに来るようになって
います。

978 なければなりません。項目を集めたものを**項目プール (Pool of Items)** といいます。これも数千から数
 979 万の単位で用意しておく必要があります。なぜなら、テスト項目は事前に 1 回は使っているわけですから、数
 980 えるほどしか項目がなければ受験生が正答を事前に丸暗記できてしまうからです。もちろん項目プールが数
 981 千から数万あっても一度どこかで使われていますから、過去問をすべて完全に丸暗記すればその人は満点
 982 が取れてしまいます。もっともそれだけのものを覚えられるのは、ある意味学力が高いといっても差し支えな
 983 いと思いますが。

984 日本でおこなわれる大学入学共通試験をはじめ、試験問題というのは極めて厳重な管理下に置かれ、事
 985 前にその情報が漏れることは公平性の観点から言って不適切であるとされています。しかし IRT で分析す
 986 るためには、事前に項目の特徴を知っていなければならないのです。CAT をつかって入学試験などができ
 987 れば、カンニング対策にもなりますし、受験生は何度でもチャレンジできるので利点も多いのですが、「公平
 988 性のために新しいテストでなければならない」となるとなかなか実用化できないところがあるというのも事実
 989 です*8。

990 ところで、心理学の場合は学力テストのように正答・誤答があるものではなく、「当てはまる」から「当てはま
 991 らない」といった軸上で多段階の反応を求めることが一般的です。テスト理論を多段階のモデルに応用できる
 992 のかと言えば、幸いにしてその答えは YES です。

993 5.2 段階反応モデル

994 リッカート法などで作られる尺度は、一般に 5, 7 段階のものが多くあります。少ないものでは 3 件法*9で
 995 あったり、ものによっては 4, 6 件法であることもあります*10。しかしこれらの段階はいずれも順序尺度水準
 996 の情報しか持っておらず、そのままでは尺度値として使うことができません。シグマ法などで数値化すれば良
 997 いのですが、その手間を省いて分析する悪い習慣もあることは既に指摘した通りです。

998 IRT の多段階版はこうした問題に対応できる方法です。IRT の多段階モデルは大きくわけて 2 つあり、
 999 1 つは**段階反応モデル (Graded Response Model; GRM)**(Samejima, 1997)、もうひとつは**多段階採
 1000 点モデル (Partial Credit Model)**(Muraki, 1992)と呼ばれています。どちらも発想は似たようなところ
 1001 があり、ここでは GRM について解説します。興味がある人は、豊田 (2012) や加藤・山田・川端 (2014) な
 1002 ど専門書を参考にしてください。

1003 GRM の考え方の基本は、段階反応の背後には正規分布する連続的な潜在特性 θ がある、と仮定すると
 1004 ころです。心理的な能力、学力、性質などは連続的なのですが、それが表に出てくる時は離散的 (順序的) だ
 1005 というわけです。図 5.3 に図示されているように、正規分布がある**閾値 (threshold)**(これを b_k と表します
 1006 が) を超えると出現する時は次のカテゴリになる、ということを考えます。図は三段階の例ですが、図から明ら
 1007 かなように k 段階であれば閾値の数は $k - 1$ 個あることとなります。横軸 θ は心理的な態度や性質の強さだ
 1008 と思ってください。さてそうすると、 $\theta = 2.0$ ぐらいであれば、ほぼ間違いなく「当てはまる」に回答すること
 1009 になりますし、 $\theta = -2$ であれば「当てはまらない」に回答するようになるはずで。このように θ が大きくなれ
 1010 ばなるほど最後のカテゴリに反応する確率は上がっていきますから、ここは 2PL モデルの時のようにロジス
 1011 ティック曲線で「当てはまる」に回答する確率を表現できるでしょう。問題は、それより下の段階に反応する確
 1012 率をどのように表現するか、です。

*8 令和 2 年度に大学入試センター試験から大学入学共通試験に変わりましたが、改革前の計画では CAT 化することが盛り込まれていました。しかし実際には、受験生のためのコンピュータやタブレットを準備したり、安定した通信網が必要であったり、というハード的な問題もあって見送られてしまいました。

*9 たとえば YG 性格検査は 3 段階です。

*10 偶数の段階にすることで、必ずどちらかの極に寄るようにして集計できます。日本人は「どちらでもない」に回答しがちな中点集中傾向があるとも言われているので、わざと肯定・否定のどちらかに寄せようという考え方です。

1013 ここである段階に反応する確率を考えるために、少し表現を改めます。すなわち、個人 i の項目 j に
 1014 対する反応が、カテゴリ k 以上になる確率として、 $P_{jk}^+(\theta) = P(x_{ij} \geq k|\theta)$ をまず考えます。ここで
 1015 $k = 0, 1, 2, \dots, K$ とします。先ほど示したように、 $k = K$ 、すなわち一番上のカテゴリ (ここでは「当てはま
 1016 る」) に回答する確率は、2PL ロジスティック関数と同じ形をしていますから、次のように表現できます。

$$P_{jK}(\theta) = P_{jK}^+(\theta) = \frac{1}{1 + \exp(-a_j(\theta - b_{jK}))}$$

1017 ここで右辺の a_j は**識別力**、 b_{jK} はカテゴリ K の**困難度**を表しています。左辺の $P_{jK}(\theta)$ は θ の人が項
 1018 目 j のカテゴリ K に反応する確率で、それは k 以上に反応する確率 $P_{jk}^+(\theta)$ と一致していることを表してい
 1019 ます。

1020 次に、もっとも低い段階に回答する確率を考えましょう。これは θ が大きくなればなるほど減っていくはず
 1021 で、いわばロジスティック曲線の逆のような形をするはずで、反応確率は最大でも 1.0 ですから、逆という
 1022 ことは 1.0 から引いてやればよいでしょう。

$$P_{j0}^+(\theta) = 1.0 - \frac{1}{1 + \exp(-a_j(\theta - b_{j0}))}$$

1023 問題は「どちらでもない」に反応する確率です。これは引き算で考えることができます。すなわち「どちらでも
 1024 ない」以上に反応する確率から、「当てはまる」以上に反応する確率を引いてやれば良いのです。

$$P_{jk}(\theta) = P_{jk}^+(\theta) - P_{j,k+1}^+(\theta)$$

1025 ここにあるように、段階数が増えたとしても k 番目のカテゴリ以上に反応する確率から、 $k + 1$ 番目に
 1026 反応する確率を引いてやれば、 k 番目のカテゴリに反応する確率が得られます。この計算をして描かれ
 1027 る曲線のことを**項目反応カテゴリ特性曲線 (Item Response Category Characteristic Curve; IRCCC)**、あるいは単に**カテゴリ確率曲線 (Category Probability Curve)**と呼ばれます。IRCCC
 1028 は図 5.3 の下段に示されています。

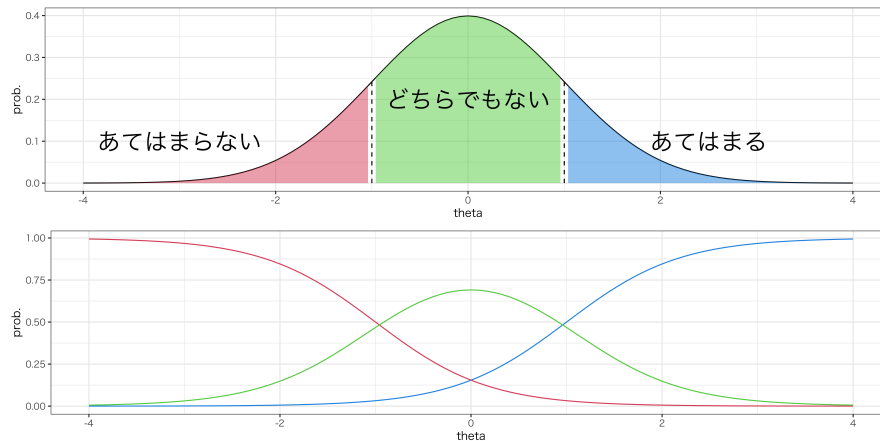


図 5.3 正規分布と閾値 (上図) と IRCCC(下図)

1030 IRCCC を、 θ の値がマイナスからプラスの方向に動かしながら見ていってください。最初は当然「当ては
 1031 まらない」に反応する確率が一番高いのですが、それが徐々に下がっていきます。「どちらでもない」の反応確
 1032 率は徐々に増えていき、閾値 b_{j1} で「当てはまらない」と逆転しピークを迎えることとなります。その頃「当ては

1033 まる」も徐々に増えていき、閾値 b_{j2} でピークが逆転する、というようになります。ピークは逆転されても、他の確
 1034 率が0になっているわけではなく、可能性は残っています。また IRCCC も ICC 同様に変換して、情報曲線
 1035 に帰することができます。すなわちこのあたりで鋭敏に情報を検出できるかを表現する項目反応カテゴリ情報
 1036 曲線を描くこともできます。このようにして、段階反応でもその項目の特徴をデザインできるのです。

1037 5.2.1 適切な反応段階を考える

1038 実際の調査研究をすると、図 5.4 の上の段のような IRCCC が描かれることがあります。何かおかしいところ
 1039 がありますか？これは 5 段階の反応モデルですが、4 番目の反応カテゴリがずいぶん低く、そのピーク
 1040 が 3 番目と 5 番目の反応カテゴリに潰されてしまっていますね。つまり、4 番目の反応がもっとも際立つシー
 1041 ンがないということです。言い換えるならば、これは尺度作成側が 5 段階だと思っていたにもかかわらず、回
 1042 答者はどういう時に「やや当てはまる」と答えるのかがはっきりせず、実質 4 段階でしか反応していないことを
 1043 表しています。

1044 このような IRCCC が描かれてしまう場合は、 $k = 3$ の反応と $k = 4$ の反応を合わせてひとつにしてしま
 1045 うなど、段階の修正を考えると良いでしょう。具体的にはデータで 4 と入力していたものを、3 に置き換えたり
 1046 します*11。修正したのが下の図になります。このように修正しても、情報関数は変わりません。同じデータから
 1047 得られる情報は同じだからです。

1048 このように 5 段階、7 段階を設定して回答者に無理やり回答を求めても、分析するとカテゴリのピークが潰
 1049 れていることがあります。回答者の反応しやすいカテゴリ数を丁寧に設計してやることが重要です。もちろん
 無分別に尺度値をつけて、そのまま分析するのはもっとも不適切な方法です。

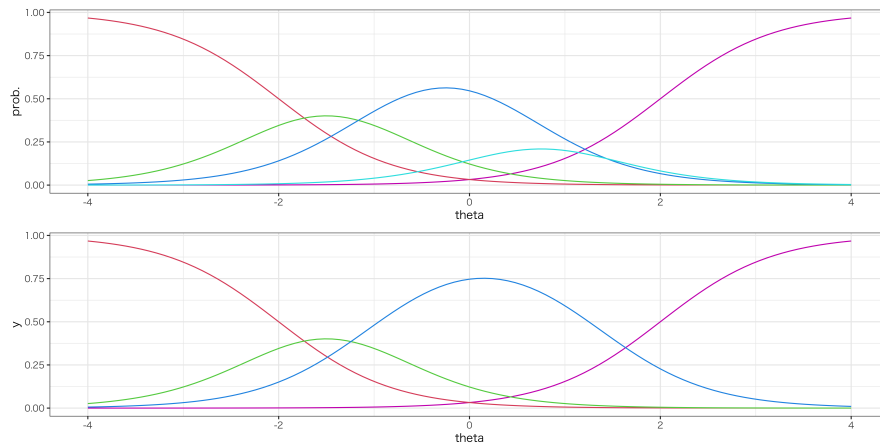


図 5.4 適切な反応段階をデザインする

1050

1051 5.3 因子分析の歴史と展開

1052 ところで、因子分析モデルもテスト理論も、潜在変数モデルとしては同じでいずれも古典的テスト理論から
 1053 の発展系なのでした。因子分析モデルは多段階反応が一般的で、多因子モデルで「潜在的な (心理学的な)
 1054 構造はどうか」ということを問題にします。ここでの目的は全体に共通する要素やその構造であり、何種類に
 1055 別れて要素間の関係はどうなっているのか、というところが中心的関心事になります。一方、項目反応理論は

*11 4 を 5 に書き換えても構いません。ヒストグラムを見て、より正規分布っぽくなるようにすると良いでしょう

1056 バイナリ反応が一般的で、因子数はひとつです。学力テストはその学力が測定できていることが重要で、因子
1057 の構造よりも因子得点をより精緻に推定できることの方が重要だからです。因子分析の言葉で言えば、因子
1058 得点をより精緻に表現しようとしているわけです。

1059 さて、GRMは、項目反応理論の多段階モデルでした。実はGRMは因子分析モデルの発展系でもある
1060 のです。因子分析は相関係数のモデルであったことを思い出してください。因子分析モデル自体は z_{ij} から
1061 始まっていましたが、変数同士の共分散 r_{ij} を考えるといくつかの仮定から**因子負荷量**だけのモデルに簡略
1062 化され、推定できるようになるのです*12。この標準化された共分散、すなわち相関係数はピアソンの積率相
1063 関係数とも言われ、間隔尺度水準以上の数値を使って計算されます。多段階の反応は順序尺度水準です
1064 から、相関係数を計算するのは不適切で、そのまま因子分析することはできません*13。では**順序尺度水準**の相
1065 関係数がないのかといわれると、あります。

1066 順序尺度水準の変数 × 順序尺度水準の変数の相関は**ポリコリック相関係数 (polychoric correlation)**
1067 といいます。順序尺度水準の変数 × 間隔尺度水準の変数の相関は**ポリシリアル相関係数**
1068 **(polyserial correlation)** といいます。ついでにバイナリ変数 × バイナリ変数の相関係数は**テトラコリック**
1069 **相関係数 (tetracholic correlation)** といいます。

1070 これらの相関係数はいずれも、順序(あるいはバイナリ)変数の背後に連続体があると考え、潜在的な連
続体 × 潜在的な連続体の相関係数を連続体のカテゴリが変わる閾値とともに推定するのです(図5.5)。

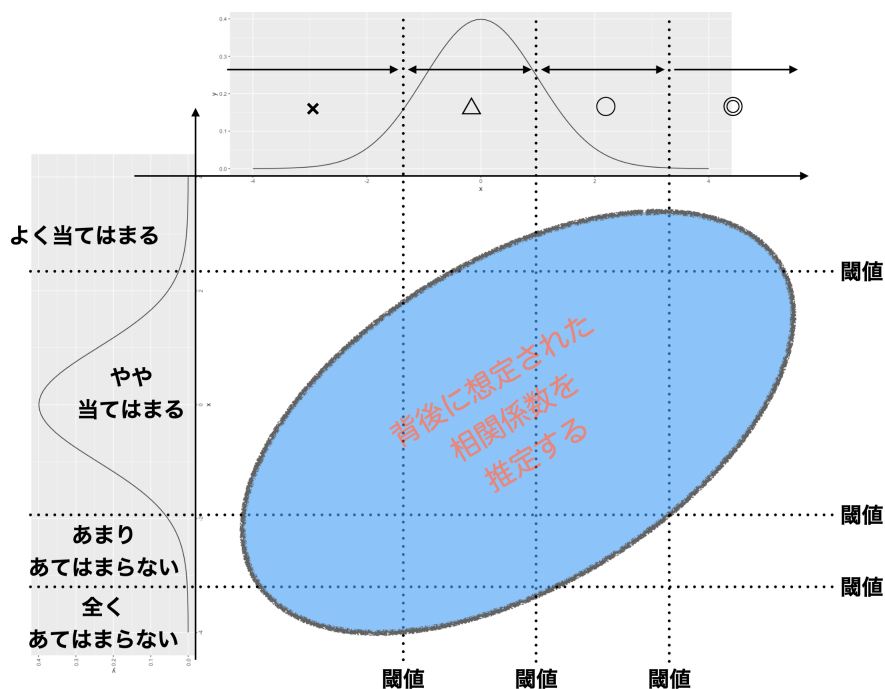


図 5.5 ポリコリック相関係数のイメージ

1071 こうして推定された相関係数を使って因子分析を行うと、順序尺度水準に適した因子分析を行うことが
1072 できます。この方法をとくに**カテゴリカル因子分析 (categorical factor analysis)** といいます。もっともこ
1073 の名前を覚えておく必要はありません。**カテゴリカル因子分析**は**段階反応モデル**と**数学的に等価**であるこ
1074

*12 具体的にどうやって因子負荷量を算出するかは次回以降のお楽しみです。

*13 できないのですが、尺度値に変換することもなく機械的に分析してしまう悪い習慣が蔓延しているのは何度も指摘している通りです。くどいと思われるかもしれませんが、私は憤っているのです。

1075 とがわかっています。GRM をすることはカテゴリカル因子分析をしていることと同じ、なのです。

1076 もっとも GRM は項目反応理論の系列ですから、単因子構造を仮定しています。複数の因子を想定す
1077 るカテゴリカル因子分析に対応するのは、正確には**多次元項目反応理論 (Multidimensional Item**
1078 **Response Theory)** といいます。しかし数学的・技術的には同じであり、カテゴリカル因子分析をした結
1079 果から IRCCC を描くこともできますし、実際に分析するソフト上では使用する変数がカテゴリカル (順序尺
1080 度水準) であることを指定するだけです。われわれユーザはもはや悩む必要はなく、ただただデータに適した
1081 分析をするだけで良いのです。

1082 5.3.1 系譜の違いはどこに関係するか

1083 因子分析モデルと項目反応理論が、カテゴリカル因子分析として概念的に統合されることを話してきました
1084 た。本質的にはこのように違いがないのですが、系譜の違い、出自の違いはそれぞれの利用される文脈で、
1085 何を強調するかに影響してくることがあります。

1086 たとえば因子分析の文脈では、共通性が低い項目は削除し、綺麗な因子構造を目指そうという考え方があ
1087 ります。尺度作成の中で1つの項目は1つの因子に対応しているべきであるという考え方があり (**単純構造**
1088 **の原理 (Principle of Simple Structure)** といいます)、もし1つの項目が複数の因子の影響を受けて
1089 いるようであれば、「美しいので」削除されることがあります。因子分析は知能、性格、態度の研究で展開
1090 されてきたため、美しい「構造」を見つけ出すことに狙いがありますから、この美しさを損なうもの (項目) は取
1091 り除く、という方向に行きがちです。

1092 一方で、項目反応理論はテスト理論の生まれです。もちろん回答者の能力や特性を測定するのに優れた項
1093 目とそうでない項目、という峻別はしますが、中でも「この項目は測定能力の偏差値 30 程度の回答者を測定
1094 するのに適している」とか、「偏差値 70 程度の回答者を測定するのに適している」と判断できます。偏差値 30
1095 や 70 を測定するのに適した項目とは、非常に簡単 (ほとんどの人が正答する) だったり、非常に難易度が高
1096 い (ほとんどの人が誤答する) 項目です。心理尺度でいうと床効果、天井効果がみられる項目とされる、ど
1097 らかに偏った分布をもつ項目です。しかし、それを捨てるということにはならず、どのような項目でも回答者の
1098 能力を推定するための情報量がゼロではない、という考えから、さまざまな項目をどんどんためていく方向に
1099 いきがちです。項目反応理論の文脈では、あらゆる人に対する測定を準備しておく必要があり、むしろ回答者
1100 の特性にあわせて設問の方を選んだり、事前の項目特性から、前もってテストの項目構成をデザインする、と
1101 いうことを目的とするのです。

1102 このように使われるシーンによって、「構造」か「機能 (=得点)」のどちらに注目するかが変わり、結果的に実
1103 践の方針がちがってくることもあるのです。図 5.6 にこの心理学的系譜 (左ルート) とテスト理論的系譜 (右
1104 ルート) の流れを描いてみました。

1105 ポイントは「最終的には同じところに辿り着く」という点ですから、歴史的流れや個々のモデルの細かい数式
1106 を完全に理解していなくてもいいかもしれません。それよりは、心理学者として、あるいはテストをする側とし
1107 て、回答者に無理のない、それでいて必要な程度に精緻な情報が得られる適切な分析方法を選択できるよう
1108 になることが重要です。

1109 5.4 課題

1110 ■**信頼性についての考え方** 古典的テスト理論、因子分析モデル、項目反応理論、それぞれの信頼性の考
1111 え方を数式とともに自分のことばで説明できるようになろう。

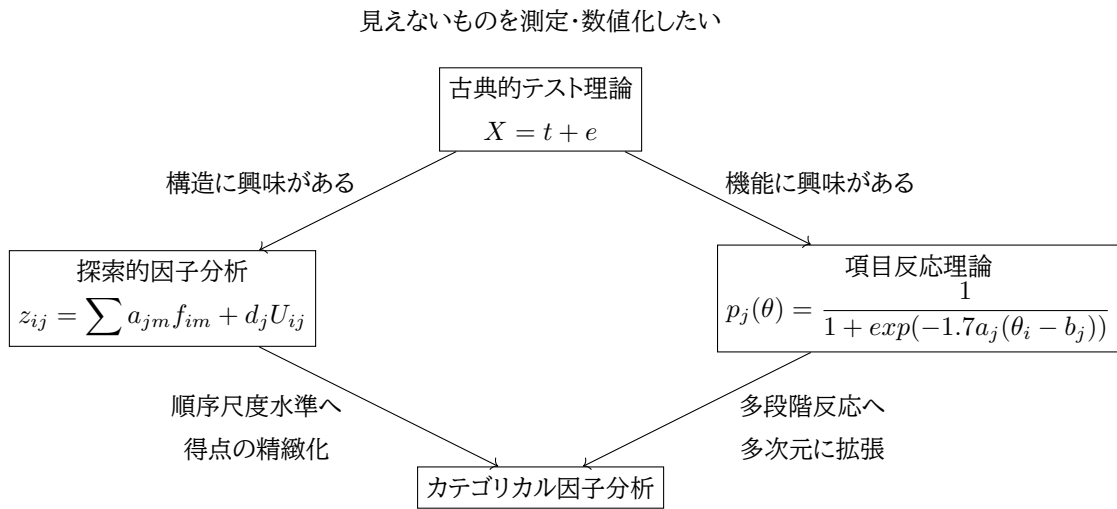


図 5.6 理論・モデルの流れ。左側のルートが心理学的系譜，右側のルートがテスト理論的系譜

- 1112 ■段階反応モデルの IRCCC 多段階モデルの IRCCC を，関数を描画するソフトなどを使って自分で描い
 1113 てみよう。

1114 第6章

1115 行列計算の基礎

1116 これまで古典的テスト理論、因子分析論、現代テスト理論を通じて、目に見えない潜在変数を数値化する
 1117 方法について学んできました。潜在変数という心のモデルは、心理学の中心的関心事であり、実際多くの調
 1118 査研究で潜在変数をモデルに組み込んで検証されています。その割には、どういったメカニズムで潜在変数
 1119 が見出されているのかについての理解は十分行き渡っていないようです。たとえば因子分析モデルは、統計
 1120 パッケージを使うと瞬時に「3 因子構造で因子負荷量はこれこれ、因子得点はこのようになっています」と答え
 1121 を出してくれます。しかし、なぜそのような数字になったのか、どのようにそれが算出されたのかを知らなけれ
 1122 ば、何もわかっていないのと同じではないでしょうか？ 因子は「機械がやってくれるもの」と思考停止しま
 1123 うと、結局のところ私たちの知りたいことには辿り着けませんし、誤用の元になってしまいます。

1124 なぜその肝心の箇所が放置されているかという点、数学的には線形代数 (linear algebra) と呼ばれる
 1125 計算が必要であり、そこについての文系数学的解説がないからです。線形代数はベクトルや行列の計算、文
 1126 字と式の便利な表現形式です。これを知ることの利点は、多くの数字のセットを簡単な記号で一般的に表現
 1127 できるようになることです。変数や回答者数が数十、数百、時には数万のサイズで得られた時、1つ1つの
 1128 データにアルファベットを割り振っていたのでは間に合いませんので、線形代数はデータ解析には必須の知
 1129 識です。

1130 本講義では、線形代数の基礎を導入した上で、最終的には潜在変数、共通因子や因子負荷量と呼ばれる
 1131 ものがどのように算出されるのかを理解することを目的としています。事前の知識は必要なく、また目的に必
 1132 要な最小限の知識だけで進めていきますので、一歩ずつ確実にフォローしてください*1。

1133 6.1 行列とベクトル

1134 行列やベクトルは、複数の数字をひとまとめにして扱うためのものです。まずはその基本的な形からみてい
 1135 きます。

1136 ■ベクトル 複数の数字を一行、あるいは一列にまとめて表現したものを、行ベクトル (row vector)、列
 1137 ベクトル (column vector) といいます。

1138 行ベクトルは次のように表します。

$$\mathbf{a} = (a_1 \quad a_2 \quad \cdots \quad a_m)$$

*1 ここからの話は小杉 (2018) の pp.148–179 に同内容のものがあります。もちろん線形代数のテキストとしては他にもいろいろあり、数学的な入門としては、基礎的には村上・佐藤・野澤・稲葉 (2016) が、発展的などころでは永田 (2005) が参考になるでしょう。より文系のデータ解析的解説が多いのは、絶版になってしまいましたが岡太 (2008) が最高です。

1139 列ベクトルは次のように表します。

$$\mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

1140 具体的には、 a_1 とか b_2 のところには数字が入っています。つまり次のような形です。

$$\mathbf{a} = (1 \quad 3 \quad 5), \mathbf{b} = \begin{pmatrix} 2 \\ 4 \\ 6 \\ 12 \\ 8 \end{pmatrix}$$

1141 ここで今回の \mathbf{a} は 3 つの要素が入っていますので、サイズは 3、同じく \mathbf{b} はサイズが 5 のベクトルです。行
1142 列の言い方に合わせて 1×3 の (行) ベクトル、 5×1 の (列) ベクトル、という言い方をすることもあります。
1143 このベクトルの中の数字は、とくに関係があるわけではありません。前に入っている数字がえらいとか、横にあ
1144 る方が重要だ、といったことはなく、ただただ数字をまとめて扱っているだけです。数字のセットを記号ひとつ
1145 で表せるので、ずいぶん楽ですよね。

1146 さて、行数も列数も 1 であるものつまり行列でない数字は、とくに **スカラー (scalar)** と呼びます。今まで
1147 は $1 + 2 = 3$ といった計算をしていましたが、この 1, 2, 3 はすべてスカラーだといえるわけです。

1148 ■**行列 行列 (matrix)** とは数を長方形に並べたものです。行列として並べられた数を成分といい、成分
1149 の横の並びを行、縦の並びを列と呼びます。

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix}$$

1150 この例は、 n 行 m 列の行列を表しています。お気づきかと思いますが、行列やベクトルを表す場合は、ア
1151 ルファベットを太字にするのが慣例です。たとえば、 A とか x は 1 つの数字を表していますが、 \mathbf{A} や \mathbf{x} であ
1152 れば行列やベクトルを表していることになります。成分を表す文字は、一般に a_{ij} のように、はじめの添え字で
1153 行番号、次の添え字で列番号を表します。行列の大きさは行数と列数とによって、 $n \times m$ のように表現しま
1154 す。 n と m が同じ、つまり行数と列数が同じであれば、これをとくに正方形行列といいます。正方形行列の例を次
1155 にあげておきます。

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix}$$

1156 正方形行列の中でも、 i 行 j 列目の値が j 行 i 列目の値と同じである行列 ($a_{ij} = a_{ji}$) のことを **対称行列**
1157 (**Symmetric Matrix**) といいます。

$$\mathbf{A} = \begin{pmatrix} 1 & 2 & 3 \\ 2 & 5 & 6 \\ 3 & 6 & 9 \end{pmatrix}$$

1158 この (正方) 対称行列の形は、データ解析の中ではよくでてきます。たとえば 3 つの変数 x_1, x_2, x_3 につい
1159 て、その相関係数を考えたいとしましょう。相関係数は 2 つの数字の組み合わせですから、 x_1 と x_2 、 x_1 と
1160 x_3 、 x_2 と x_3 について計算でき、それぞれ r_{12}, r_{13}, r_{23} と表したとします。 i と j の相関係数 r_{ij} は、 j と i

1161 の相関係数と同じ ($r_{ij} = r_{ji}$) であり, また $r_{jj} = 1.0$ なのは定義から明らかです。これを行列で表すと次の
1162 ようになります。

$$\mathbf{R} = \begin{pmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{pmatrix} = \begin{pmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{pmatrix}$$

1163 このように対称行列になっています。この行列をとくに**相関行列 (Correlation Matrix)** といいま
1164 す。また相関係数は標準化された共分散でもありました。標準化するまえの相関行列は、**分散共分散行列**
1165 **(Covariance Matrix)** と言います。その名前の通り, 自分自身との共分散が分散になるわけですから,
1166 右上から右下にかけての対角線上にある要素 (これをとくに**対角 (diagonal) 要素** と言います) が分散であ
1167 り, それ以外が共分散になっている行列です。

$$\mathbf{V} = \begin{pmatrix} s_1^2 & s_{12} & s_{13} \\ s_{21} & s_2^2 & s_{23} \\ s_{31} & s_{32} & s_3^2 \end{pmatrix}$$

1168 また, 正方行列の中でもとくに対角要素にのみ値があつて, それ以外はすべて 0 になっている行列のこ
1169 とを**対角行列 (diagonal matrix)**, 対角行列の中でもとくに, 対角項が 1 になっているものは**単位行列**
1170 **(identity matrix)** と呼びます。単位行列は \mathbf{I} とか \mathbf{E} で表されます。

$$\mathbf{I} = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

1171 これは後ほど, 掛け算をするときに「かけても変わらない状態」を表すために用いられます。

1172 6.2 行列の四則演算と操作

1173 行列の四則演算は, 通常のスカラーのそれとは異なります。改めて, 行列としての加減乗除を定義するのだ
1174 と思ってください。

1175 ■**加法・減法** まずは行列の足し算 (加法), 引き算 (減法) から説明します*2。これはそれぞれ対応する位
1176 置にある成分を加え合わせる (減じる) ことで表されます。

$$\mathbf{A} + \mathbf{B} = \begin{pmatrix} a_{11} + b_{11} & a_{12} + b_{12} & \cdots & a_{1m} + b_{1m} \\ a_{21} + b_{21} & a_{22} + b_{22} & \cdots & a_{2m} + b_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} + b_{n1} & a_{n2} + b_{n2} & \cdots & a_{nm} + b_{nm} \end{pmatrix}$$

1177 数値例をみておきましょう。

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} + \begin{pmatrix} 5 & 6 \\ 7 & 8 \end{pmatrix} = \begin{pmatrix} 1+5 & 2+6 \\ 3+7 & 4+8 \end{pmatrix} = \begin{pmatrix} 6 & 8 \\ 10 & 12 \end{pmatrix}$$

1178 これからわかるように, 行列の加法, 減法は大きさの等しい行列でないとなり立ちません。サイズが違うも
1179 のを足そうとすると, 演算できない箇所が出てしまうのです。このように行列では, 「計算できない」という状
1180 態になることが少なからずあります。行列のサイズに注意が必要, ということがお分かりいただけるかと思
1181 います。

*2 ベクトルは行列の中でも, 行数あるいは列数が 1 のものですので, これで一般的に表現します。

1182 ■乗法 続いて掛け算です。まずスカラーと行列の積を見てみましょう*3。

$$\lambda \mathbf{A} = \mathbf{A} \lambda = \begin{pmatrix} \lambda a_{11} & \lambda a_{12} & \cdots & \lambda a_{1m} \\ \lambda a_{21} & \lambda a_{22} & \cdots & \lambda a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \lambda a_{n1} & \lambda a_{n2} & \cdots & \lambda a_{nm} \end{pmatrix}$$

1183 実際の計算は、各成分をスカラー倍すればよいだけですので、比較的簡単ですね。

$$2 \times \begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} = \begin{pmatrix} 2 \times 1 & 2 \times 2 \\ 2 \times 3 & 2 \times 4 \end{pmatrix} = \begin{pmatrix} 2 & 4 \\ 6 & 8 \end{pmatrix}$$

1184 次はベクトルとベクトルの掛け算です。これは形が変わってしまうので、注意が必要です。まずは行ベクトル
1185 に列ベクトルをかける例からみていきましょう。

$$\mathbf{a} \mathbf{b} = (a_1 \ a_2 \ \cdots \ a_n) \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix} = \sum_{j=1}^n a_j b_j$$

1186 掛け算なのですが、足し合わせるという計算プロセスが入り込んでいるので、結果はスカラーになります。
1187 掛け算なのにどうして足し算の要素が入るんだ、というクレームは、今はなしです。このように計算することに
1188 決めたことで、あとあと便利なことが出て来ますから、作法にまず慣れてからにしましょう。数値例も確認して
1189 おきます。

$$(1 \ 2 \ 1) \begin{pmatrix} 3 \\ 4 \\ 2 \end{pmatrix} = 1 \times 3 + 2 \times 4 + 1 \times 2 = 13$$

1190 ここで注意して欲しいのは、両方のベクトルのサイズが同じ ($1 \times n$ ベクトルと、 $n \times 1$ ベクトル、いずれも
1191 サイズは n) ということです。サイズが違ると、演算が対応しない要素が出てくるので、計算できない、が答え
1192 になります。

1193 今度は向きを変えて、列ベクトルに右から行ベクトルをかけてみましょう。

$$\mathbf{a} \mathbf{b} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix} (b_1 \ b_2 \ \cdots \ b_n) = \begin{pmatrix} a_1 b_1 & a_1 b_2 & \cdots & a_1 b_n \\ a_2 b_1 & a_2 b_2 & \cdots & a_2 b_n \\ \vdots & \vdots & \ddots & \vdots \\ a_n b_1 & a_n b_2 & \cdots & a_n b_n \end{pmatrix}$$

1194 今度は行列になりました。かける順番が変わるとサイズが変わる (ここでは、上の例では 1×1 のサイズ、
1195 下の例では $n \times n$ のサイズ) ことに注意してください。スカラーの計算では順番を入れ替えても、たとえば
1196 $2 \times 3 = 3 \times 2$ のように同じ答えになりましたが、行列の場合は必ずしもそうはならない、ということです。

$$\begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} (3 \ 4 \ 2) = \begin{pmatrix} 1 \times 3 & 1 \times 4 & 1 \times 2 \\ 2 \times 3 & 2 \times 4 & 2 \times 2 \\ 1 \times 3 & 1 \times 4 & 1 \times 2 \end{pmatrix} = \begin{pmatrix} 3 & 4 & 2 \\ 6 & 8 & 4 \\ 3 & 4 & 2 \end{pmatrix}$$

1197 行列とベクトルの積や、行列と行列の積はこの応用になってきます。まず行列に列ベクトルを右からかける
1198 例を見てみましょう。結果は列ベクトルになります。

*3 式中にてでくる λ はギリシア文字でラムダといいます。小文字が λ 、大文字では Λ と書きます。

$$\mathbf{Ab} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} = \begin{pmatrix} \sum_{j=1}^m a_{1j}b_j \\ \sum_{j=1}^m a_{2j}b_j \\ \vdots \\ \sum_{j=1}^m a_{nj}b_j \end{pmatrix}$$

1199 ここでも掛け算なのに足し算のプロセスが入ってきています。注意深く記号を読んでみてください。数値例
1200 でも確認しておきます。

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \end{pmatrix} \begin{pmatrix} 2 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \times 2 + 2 \times 1 \\ 3 \times 2 + 4 \times 1 \end{pmatrix} = \begin{pmatrix} 4 \\ 10 \end{pmatrix}$$

1201 今度は行列に行ベクトルを左からかけましょう。結果は行ベクトルになります。

$$\mathbf{cA} = (c_1 \quad c_2 \quad \cdots \quad c_n) \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} = \left(\sum_{j=1}^n a_{j1}c_j \quad \sum_{j=1}^n a_{j2}c_j \quad \cdots \quad \sum_{j=1}^n a_{jm}c_j \right)$$

$$(1 \quad 3) \begin{pmatrix} 1 & 0 \\ 2 & 3 \end{pmatrix} = (1 \times 1 + 3 \times 2 \quad 1 \times 0 + 3 \times 3) = (7 \quad 9)$$

1202 さて、最後に行列と行列の積を考えます。行列 \mathbf{A} と \mathbf{B} の積が成立するのは、前者の列数と後者の行数と
1203 が等しいときに限られます。行列 \mathbf{A} のサイズが $n \times m$ 、行列 \mathbf{B} のサイズが $m \times l$ とすると、その積は $n \times l$
1204 の行列になります。計算手続きは、次のようになります。

$$\mathbf{AB} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} & \cdots & b_{1l} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \cdots & b_{ml} \end{pmatrix} = \begin{pmatrix} \sum_{j=1}^m a_{1j}b_{j1} & \cdots & \sum_{j=1}^m a_{1j}b_{jl} \\ \vdots & \ddots & \vdots \\ \sum_{j=1}^m a_{nj}b_{j1} & \cdots & \sum_{j=1}^m a_{nj}b_{jl} \end{pmatrix}$$

1205 どうにもこれはややこしいかもしれません。足し算や掛け算が入り乱れるし、計算途中でどの要素を計算し
1206 ているかわからなくなるからです。ベクトルと行列の積の時のように、前の行列の要素は左に進み、後ろの行
1207 列の要素は縦に進みますから、左手と右手で違う図形を描く認知課題のように、そもそも混乱しやすい作業
1208 なのです。

1209 しかし 2 つほど注意をしておく、間違いにくくなります。1 つは積によって得られる**結果の行列サイズを**
1210 **意識すること**です。先ほど、前の行列の列数と、後ろの行列の行数が同じでないと計算できないといいま
1211 した。つまり、 $n \times m$ 行列と $m \times l$ 行列でないと計算できない (m が同じ) ということです。また、結果は
1212 $n \times l$ 行列になります。前の行列の行数、後ろの行列の列数が結果のサイズです。ここに注目しておく、計算
1213 を始める前に、計算が可能かどうかと結果の行列サイズは想像がつかののです (図 6.1)。

1214 また、実際に計算する際は、**前の行列に横の、後ろの行列に縦の補助線を入れる**とわかりやすいかもしれ
1215 ません。こうすることで、間違えて計算を進めることがないようになるからです。

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix} \left(\begin{array}{c|c|c} 0 & 1 & 1 \\ 1 & 0 & 1 \end{array} \right) =$$

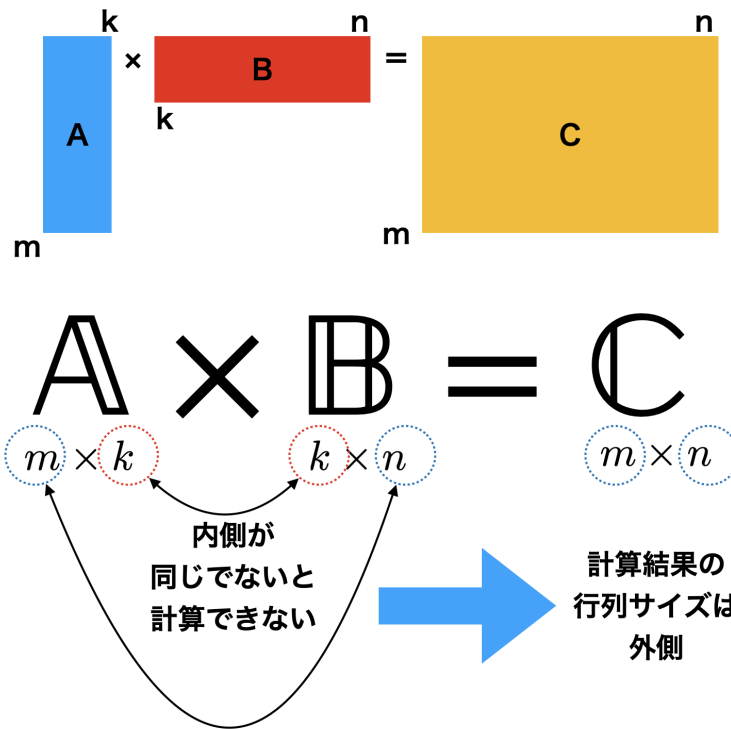


図 6.1 行列のサイズを把握する

$$\begin{pmatrix} 1 \times 0 + 2 \times 1 & 1 \times 1 + 2 \times 0 & 1 \times 1 + 2 \times 1 \\ 3 \times 0 + 4 \times 1 & 3 \times 1 + 4 \times 0 & 3 \times 1 + 4 \times 1 \\ 5 \times 0 + 6 \times 1 & 5 \times 1 + 6 \times 0 & 5 \times 1 + 6 \times 1 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 3 \\ 4 & 3 & 7 \\ 6 & 5 & 11 \end{pmatrix}$$

1216 ■転置 次に転置 (transpose) と呼ばれる操作を説明します。これは計算の便宜上、よく使われる行列操
1217 作のひとつです。

1218 大きさ $n \times m$ の行列 A における i 行 j 列成分を j 行 i 列成分とする $m \times n$ 行列のことを、元の行列 A
1219 の転置とよび、 A' や A^T と表します。行列を転ばせたようなイメージです。

$$A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nm} \end{pmatrix} \text{ のとき, } A' = \begin{pmatrix} a_{11} & a_{21} & \cdots & a_{n1} \\ a_{12} & a_{22} & \cdots & a_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ a_{1m} & a_{2m} & \cdots & a_{nm} \end{pmatrix}$$

1220 ベクトルも転置でき、行ベクトルを転置すると列ベクトルに、列ベクトルを転置すると行ベクトルになります。

$$\mathbf{a} = (a_1 \ a_2 \ \cdots \ a_n) \text{ のとき, } \mathbf{a}' = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

1221 また、転置には以下のような性質があります。これは知識として知っておくだけでよいでしょう。

- 1222 1. $(A')' = A$
- 1223 2. $(A + B)' = A' + B'$

1224 3. $(AB)' = B'A'$

1225 4. $(cA)' = cA'$

1226 ■**逆行列** 最後に**逆行列**のお話をします。逆行列は割り算のイメージです。ある行列にその逆行列をかける
1227 と単位行列になる、つまり割ると1になるような行列のことです。

1228 正確に表現すると、ある正方行列 A に対し、 $AX = I$ となるような行列 X が存在するとき、これを A
1229 の**逆行列 (inverse)** と呼び、 A^{-1} で表します。正方行列でない場合に逆行列はありませんし、正方行列で

1230 あっても逆行列が存在しない場合もあります。逆行列の例をみてみましょう。 $A = \begin{pmatrix} 2 & 1 \\ 5 & 3 \end{pmatrix}$ とすると、次の
1231 計算が成り立ちます。

$$AB = \begin{pmatrix} 2 & 1 \\ 5 & 3 \end{pmatrix} \begin{pmatrix} 3 & -1 \\ -5 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

1232 このとき B は A の逆行列、すなわち $A^{-1} = B$ といえます。

1233 とくに対角行列の逆行列は、対角成分の逆数をそれぞれ対角成分とする行列になります。

$$D = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{pmatrix} \text{ のとき, } D^{-1} = \begin{pmatrix} \frac{1}{d_1} & 0 & \cdots & 0 \\ 0 & \frac{1}{d_2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{d_n} \end{pmatrix}$$

1234 逆行列は、行列の世界の割り算のようなものです。これで一通り四則演算の定義ができました。

1235 6.3 行列を使うと便利なこと

1236 さて、ここまでで行列の計算の話をしてきましたが、どこが良いのかいまいちピンとこない、という人もい
1237 るかもしれません。そこで最後にどうしてこのような計算をするのか、何が良いのかを説明してみたいと思
1238 います。

1239 6.3.1 行列と方程式

1240 線形代数は「便利な書き方」の学問です。便利な書き方をするためにルールが作られていますから、ルール
1241 から学ぶと「なんでそんな変な操作をするんだ」という気持ちになるのもわかります。

1242 では何が便利になるのでしょうか。これは方程式を解くことと関係があります。たとえば、以下のような連立
1243 方程式があったとしましょう。

$$\begin{cases} x - 2y - 5z = 3 \\ 5x + 4y + 3z = 1 \\ 3x + y - 3z = 6 \end{cases}$$

1244 これは行列で表現すると、次のようになります。

$$\begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 6 \end{pmatrix}$$

1245 この左辺を行列とベクトルの式の計算ルールにのっって展開してみてください。ちゃんと最初の連立方程

1246 式の左辺になることがわかると思います。かけて足して、という面倒な計算ルールは、連立方程式を簡単に表
1247 記するためのものだったのですね。

1248 最終的にはこの方程式を解いて、次のように答えを求めます。

$$\begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix} \begin{pmatrix} x = -1 \\ y = 3 \\ z = -2 \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 6 \end{pmatrix}$$

1249 皆さんも学校で習ったように、このような連立方程式を解く方法として、加減法や代入法というのがあります。
1250 ですがここはひとつ、行列を使った解法を考えてみましょう。

1251 そのような解法のひとつ、消去法は、ひとつの方程式を何倍かして、他の方程式に加えることにより、方程
1252 式をどんどん簡単にしていくというものです。まず、第一の式を5倍、あるいは3倍して、第二、第三の式から
1253 x の項を消去します。

$$\begin{cases} x - 2y - 5z = 3 \\ -14y - 28z = 14 \\ -7y - 12z = 3 \end{cases}$$

1254 第二の式の係数を簡単におきましょう。

$$\begin{cases} x - 2y - 5z = 3 \\ y + 2z = -1 \\ -7y - 12z = 3 \end{cases}$$

1255 第二の式を7倍して、第三の式から y を消去します。

$$\begin{cases} x - 2y - 5z = 3 \\ y + 2z = -1 \\ 2z = -4 \end{cases}$$

1256 あとはこれの3行目から $z = -2$ が得られ、芋づる式に $x = -1$ 、 $y = 3$ が得られました。

1257 この操作は、式を一本ずつ、あるいは2つの式を組み合わせ文字を消していく消去法を係数全体に行う
1258 操作になっています。実際、ここで操作される係数だけ見ていくと、次のようになります。

■第一段階

$$\begin{pmatrix} 1 & -2 & -5 \\ 0 & 1 & 2 \\ 0 & -7 & -12 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 3 \\ -1 \\ 3 \end{pmatrix}$$

■第二段階

$$\begin{pmatrix} 1 & -2 & -5 \\ 0 & 1 & 2 \\ 0 & 0 & 2 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 3 \\ -1 \\ -4 \end{pmatrix}$$

1259 さらにこの方法を改良した、ガウス-ジョルダンの消去法というものがあります。この手法による係数の変
1260 化を、行列表記で見ていくことにします。

1261 まず第一段目は同じです。

$$\begin{pmatrix} 1 & -2 & -5 \\ 0 & 1 & 2 \\ 0 & -7 & -12 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 3 \\ -1 \\ 3 \end{pmatrix}$$

1262 次に、第二の方程式を用いて第一と第三の式から y の係数を消してしまいます。

$$\begin{pmatrix} 1 & 0 & -1 \\ 0 & 1 & 2 \\ 0 & 0 & -2 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 \\ -1 \\ 4 \end{pmatrix}$$

1263 最後に、第三の式の z の係数を 1 にして、第一、第二式の z の係数を消してしまいましょう。

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} -1 \\ 3 \\ -2 \end{pmatrix}$$

1264 最後の形を見ると、左辺は単位行列になっていますから、

$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} -1 \\ 3 \\ -2 \end{pmatrix}$$

1265 と解が求められたことがわかります。ここで注目すべきは、連立方程式の解を求めるプロセスは係数行列を
1266 単位行列に変えていくプロセスだった、ということです。係数行列が単位行列になれば、それはもう答えを出
1267 したことになるのです。

1268 さて、係数行列を A とすると、その逆行列 A^{-1} があれば $A^{-1}A = I$ となるのです。であれば、連立方
1269 程式の右辺にあったベクトルに A^{-1} をかけてやれば、一気に答えが求まるではないですか。

1270 実際に見て見ましょう。

$$\begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 3 \\ 1 \\ 6 \end{pmatrix}$$

1271 この連立方程式に対して、次のような操作をします。

$$\begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix}^{-1} \begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix}^{-1} \begin{pmatrix} 3 \\ 1 \\ 6 \end{pmatrix}$$

1272 とします*4。すると左辺は単位行列になりますから、次のように計算すれば一気に答えが求まることにな
1273 るのです。

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x \\ y \\ z \end{pmatrix} = \begin{pmatrix} -1 \\ 3 \\ -2 \end{pmatrix}$$

1274 つまり、連立方程式を解くという問題が、係数行列の逆行列を求める問題になります。また、逆行列は存在し
1275 ないこともある、ということでしたが、その場合その連立方程式は解けない、ということになります。

1276 6.4 課題

1277 ■線形代数の練習問題 行列 $A = \begin{pmatrix} 1 & 2 & -1 \\ 3 & 1 & 0 \end{pmatrix}$ のとき、次の計算をしなさい。なお、 I_n とは $n \times n$ の単
1278 位行列、 O とはすべての要素が 0 の適当なサイズの正方行列であることを表します。

1279 1. $A'A$

1280 2. AA'

1281 3. AI_3

1282 4. A

*4 数値的には $\begin{pmatrix} 1 & -2 & -5 \\ 5 & 4 & 3 \\ 3 & 1 & -3 \end{pmatrix}^{-1} = \begin{pmatrix} 15/28 & 11/28 & -1/2 \\ -6/7 & -3/7 & 1 \\ 1/4 & 1/4 & -1/2 \end{pmatrix}$ という行列です

1283 ■線形代数の練習問題その 2 行列 $A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$, $B = \begin{pmatrix} 3 & 1 & 5 \\ 2 & 4 & 6 \end{pmatrix}$, $C = \begin{pmatrix} 4 & 1 \\ 2 & 3 \end{pmatrix}$, 列ベクトル

1284 $x = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$, 行ベクトル $y = (2 \ 8)$ とするとき, 次の計算をしなさい。なお, 計算が定義されていないものに

1285 ついては「計算できない」と回答しなさい。

1286 1. $A + B$

1287 2. $A - C$

1288 3. AB

1289 4. AC

1290 5. $B'A$

1291 6. Ay'

1292 7. xy

1293 8. xB

1294 9. $x'B'$

1295 10. yx

1296 ■連立方程式を解く 次の連立方程式を解きなさい。

$$\begin{cases} x - 2y + 3z = 1 \\ 3x + y - 5z = -4 \\ -2x + 6y - 9z = -2 \end{cases}$$

第 7 章

行列による関係の表現

前回から線形代数の話をしています。線形代数は数字をセットで扱うための表現方法、計算方法ですから、多変量データを分析しようという時には必須の技術になります。前は線形代数の導入ですから、計算方法を解説してきましたが、今回はこの計算方法を使って具体的にデータをどのように表現し、どのように計算するのかを見ていくことになります。

7.1 データの行列表現

ここまで行列の形ばかり見て来ましたが、狙いはあくまでも調査研究など、多変量データを扱う場面での利用です。なぜ多変量データ分析をする際にこのような知識が必要なのか、思うかもしれません。ですが、得られるデータは行列として扱うと表現が大変便利なのです。たとえば質問項目が m 個あって、調査対象者 n 人から回答を得たとすると、データは次のように表現できます。

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1m} \\ x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{nm} \end{pmatrix}$$

データ全体をこうして、ひとつの記号で表現できたら便利ですね。これらを使ったデータの表記に慣れておきましょう。

各反応の平均点は以下のように表現されます。まず、要素がすべて 1 からなるベクトルを次のように表します。

$$\mathbf{1} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

わかりにくいかもしれませんが、この $\mathbf{1}$ は太字でベクトルを表しており、スカラーの 1 とは違うことに注意してください。

1314 さて、各項目の和はベクトルの掛け算の定義によって次のように表現できます。

$$\mathbf{X}'\mathbf{1} = \begin{pmatrix} \sum_{i=1}^n x_{i1} \\ \sum_{i=1}^n x_{i2} \\ \vdots \\ \sum_{i=1}^n x_{im} \end{pmatrix}$$

1315 これを使って平均値 (列) ベクトル \mathbf{m} を次のように表すことができます。

$$\mathbf{m} = \frac{1}{n}\mathbf{X}'\mathbf{1} = \begin{pmatrix} 1/n \sum_{i=1}^n x_{i1} \\ 1/n \sum_{i=1}^n x_{i2} \\ \vdots \\ 1/n \sum_{i=1}^n x_{im} \end{pmatrix} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_m \end{pmatrix}$$

1316 ここで \bar{x}_1 は 1 つめの添字 i を足し合わせて割ることでなくして書いていますから、 \bar{x}_1 のように省略して書くこと
1317 があります。このときの 1 は「第一番目の変数」という意味であり、個人の情報がなくなっている変数を意味し
1318 ていることに注意してください。

1319 さて、平均からの偏差を要素に持つ行列 \mathbf{V} を考えたとします。

$$\begin{aligned} \mathbf{V} &= \mathbf{X} - \mathbf{1}\mathbf{m}' \\ &= \mathbf{X} - \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} (\bar{x}_1 \quad \bar{x}_2 \quad \cdots \quad \bar{x}_m) \\ &= \mathbf{X} - \begin{pmatrix} \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_m \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_m \\ \vdots & \vdots & \ddots & \vdots \\ \bar{x}_1 & \bar{x}_2 & \cdots & \bar{x}_m \end{pmatrix} \end{aligned}$$

1320 この行列 \mathbf{V} のサイズは $n \times m$ であることに注意してください。これはまた、次のように表すこともできます。

$$\mathbf{V} = \left(\mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}'\right)\mathbf{X}$$

1321 ここで \mathbf{I} は適切なサイズの単位行列です*1。

1322 これを使うと、たとえば分散共分散行列 \mathbf{S} は次のようになります。

$$\mathbf{S} = \frac{1}{n}\mathbf{V}'\mathbf{V} = \begin{pmatrix} s_{11}^2 & s_{12} & \cdots & s_{1m} \\ s_{21} & s_{22}^2 & \cdots & s_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ s_{m1} & s_{m2} & \cdots & s_m^2 \end{pmatrix}$$

*1 適切なサイズってなんだよ、と思いますよね。これは計算に合うようなサイズ、という意味です。具体的に考えてみますと、 \mathbf{I} の後ろは $\mathbf{1}\mathbf{1}'$ です。 $\mathbf{1}$ は $n \times 1$ の列ベクトルで、転置したものと掛け合わせますから、 $\mathbf{1}\mathbf{1}'$ のサイズは $n \times n$ です。行列の引き算は同じサイズでないと成立しませんから、ここでの \mathbf{I} も $n \times n$ でなければなりません。カッコの中身が $n \times n$ で、それにサイズ $n \times m$ である \mathbf{X} をかけますから、計算結果や右辺のサイズは $n \times m$ になります。

1323 ここで s_j とあるのは第 j 変数の標準偏差を, s_{jk} とあるのは第 j 変数と第 k 変数の共分散です。添え字
1324 は変数番号になっています。また, ここでもサイズに注目してください。 $V'V$ は, サイズで言うと $m \times n$ と
1325 $n \times m$ の積ですから, $m \times m$ になります。この行列は正方対称行列です。

1326 また, 対角項に各変数の標準偏差 s_j が入った行列 Q を以下のように定めるとしましょう。次のような行列
1327 です。

$$Q = \begin{pmatrix} s_1 & 0 & \cdots & 0 \\ 0 & s_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & s_m \end{pmatrix}$$

1328 そうすると, この逆行列をつかって標準得点行列 Z を次のように表すことができます。

$$Z = VQ^{-1}$$

1329 さらに, これを用いて相関行列 R を次のように表すことができます。

$$R = \frac{1}{n} Z'Z = \begin{pmatrix} 1 & r_{12} & \cdots & r_{1m} \\ r_{21} & 1 & \cdots & r_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ r_{m1} & r_{m2} & \cdots & 1 \end{pmatrix}$$

1330 データのサイズにかかわらず, 一般的にこのように表現できるのはとてもわかりやすいですね。

1331 7.2 線形モデルの行列表現

1332 ベクトルや行列の記法をつかうと, 回帰分析や重回帰分析の式がとても単純な形で表現できます。

1333 回帰分析は, $Y = aX + b + e$ という式で表現できる, ということでしたが, 式中の X や Y は観測され
1334 たデータですので, $X = \{x_1, x_2, \dots, x_n\}$, $Y = \{y_1, y_2, \dots, y_n\}$ というベクトルだと考えることができま
1335 す。ですから, 正確に書けば, ベクトルを使って次のように書くべきです。

$$Y = aX + b + e$$

1336 これは, 要素を表現しながら書くと*2次ようになります。

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = a \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} + \begin{pmatrix} b \\ b \\ \vdots \\ b \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix}$$

1337 列ベクトル X に定数 a をかけて得られるのは同じサイズの列ベクトル, 列ベクトル同士は足しても同じサ
1338 イズの列ベクトルですから, 左辺と右辺はどちらも列ベクトルで, 対応関係が取れていることになります。

1339 ここで少し表現の工夫をします。説明変数 X のベクトルの左に数字の 1 だけが入った列を作ります。また,
1340 係数もまとめてベクトル β を次のように用意します。

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \beta = \begin{pmatrix} b \\ a \end{pmatrix}$$

*2 エレメントワイズ element-wise の表現, と言ったりします。

1341 このようにすると、回帰分析の式は

$$Y = X\beta + e$$

1342 と表すことができます。とても簡単な表現になりました (試しに各行を行列の計算式に則って計算してみてください。
1343 ださい。うまく表現できていることがわかると思います)。

1344 さらにこの表現はありがたいことに、複数の説明変数がある重回帰分析の時でも同じ形で表すことがで
1345 きます。重回帰分析は、これまでの書き方ですと $Y = a_1X_1 + a_2X_2 + \dots + a_nX_n + b + e$ というよう
1346 にしていました。ここで、係数と切片を 1 つの行列で表現する時わかりやすくするために、少し書き換えて
1347 $Y = \beta_0 + \beta_1X_1 + \beta_2X_2 + \dots + e$ としましょう。記号が変わっただけで中身は同じ、意味も同じです*3。た
1348 だ、切片 b を β_0 として右辺の一番前に持って来ました。というのも、そうするとベクトルで書く時にわかりや
1349 すいからです。

1350 説明変数行列の左端に 1 を入れたベクトルを追加し、回帰係数 β もセットにして、次のように表現します。

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1m} \\ 1 & x_{21} & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nm} \end{pmatrix}, \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{pmatrix}$$

1351 とすると、重回帰分析の式は次のように簡単な表現に変わります。

$$Y = X\beta + e$$

1352 なんと、説明変数が m 個に増えたのに、式の形は (単) 回帰分析のそれと同じではありませんか！

1353 このように、行列表現をすると多くの変数を一気に扱い、表現できるのです。このため、多変量解析ではベ
1354 クトルの表記が基本になります。サイズを気にせず一般的に表現できるからです。

1355 実際にこれらの式を読む時は、行列のサイズをイメージしながら読むと良いでしょう。たとえば左辺の Y は
1356 サイズ n のベクトルなので、右辺の $X\beta$ もサイズ n の縦ベクトルになるはずなのです。実際、 X は
1357 $n \times (m + 1)$ の行列で、 β は $(m + 1) \times 1$ のベクトルですから、計算結果は $n \times 1$ 、つまりサイズ n の縦ベ
1358 クトルです。

1359 7.3 デザイン行列

1360 ところで、心理統計と言えば平均値の差を見ることだ、という話はこれまで散々聞いてきたところかと思
1361 います。心理学は要因計画を立て、標本の平均値差から母集団に話を一般化するために、推測統計の知見を
1362 使って、帰無仮説検定やベイズ推定法を駆使するというやつです。この要因計画は実は線形モデルの一環で
1363 あり、**一般線形モデル (General Linear Model)** と呼ばれています。これを行列で表現することを、ここ
1364 では少し考えてみたいと思います。

1365 まずは回帰分析と要因計画は何が同じで何が違うのかを、はっきりさせましょう。同じところは線形モデル
1366 であるというところ、違うところは、回帰分析は説明変数も従属変数も連続変数であるのに対し、要因計画で
1367 は一般に説明変数が離散変数であること、でした。離散変数であるとは、言い方を変えると名義尺度水準の
1368 数字だということです。すなわち「統制群」か「実験群」か、という違いを表すのに、0, 1 と言った数字を割り

*3 厳密に言えば記述統計学として誤差を最小にするように推定した係数はアルファベット b_0, b_1, \dots で、推測統計学として母数の推定値として算出した係数はギリシア文字 β_0, β_1, \dots で表現する、というルールです。すでに習ったように、最小二乗法での推定値と最尤法での推定値は、誤差が正規分布する場合一致しますので、ただ書き変わっただけだと思っていただいて問題ありません。

1369 振ったものです。これは別に 3 と 12523, という数字を割り振ったと言ってもいいのです。だって名義尺度水
1370 準は、数字と対象が一対一対応していれば良いのですから。

1371 ここでは数学的に話をしやすくするために、統制群を 0, 実験群を 1 とするとしましょう。線形モデルとい
1372 う枠組みは一緒なので、従属変数 y_i が説明変数 x_i によって変わるわけですが、ここではこの x が
1373 $\{0, 1\}$, というわけです。線形モデルを要素ごとに表現すると次のようになります。

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

1374 この時、 i さんは統制群に割り振られていたとすると、 $x_i = 0$ ですから、この式は次のようになります。

$$y_i = \beta_0 + e_i$$

1375 逆に、 i さんが実験群に割り振られていたとしますと、 $x_i = 1$ ですから、この式は次のようになります。

$$y_i = \beta_0 + \beta_1 + e_i$$

1376 これを見るとわかるように、両群のベースライン β_0 は同じで、そこに効果 β_1 が乗っかるかどうか、が興味
1377 的になります。この式の右辺に i は誤差成分しかなく、誤差を抜きにすると従属変数は β_1 だけ変化するは
1378 ずだ、というところから、平均値の差を検証しましょうと言ってることになるのです。

1379 これも x_i が個人ごとに変わるベクトルだと考えると、行列表現では次のようになります。

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

1380 同じ形ですね！ただし、ここでベクトル \mathbf{x} は、その中身が $\mathbf{x} = (0, 0, 0, 1, 1, 0, 1, 0, \dots)$ のように実験群か
1381 統制群かを分けるフラグが入っているだけになります。

1382 以上は実験群と統制群という 2 群の話でしたが、3 群以上になっても基本的なアイデアは同じです。たと
えば表 7.1 のようなデータセットがあったとしましょう。

表 7.1 群間要因 (3 水準) のデータセット例

参加者番号	群わけ	従属変数
1	A	3
2	A	3
3	A	4
4	A	4
5	B	6
6	B	7
7	B	8
8	B	9
9	C	7
10	C	6
11	C	5
12	C	4

1383

1384 ここで群ごとの効果を表現したいとすると、次のように書くことになります。

$$\begin{aligned}
 y_1 &= \beta_0 + \beta_1 + e_1 \\
 y_2 &= \beta_0 + \beta_1 + e_2 \\
 y_3 &= \beta_0 + \beta_1 + e_3 \\
 y_4 &= \beta_0 + \beta_1 + e_4 \\
 y_5 &= \beta_0 + \beta_2 + e_5 \\
 y_6 &= \beta_0 + \beta_2 + e_6 \\
 y_7 &= \beta_0 + \beta_2 + e_7 \\
 y_8 &= \beta_0 + \beta_2 + e_8 \\
 y_9 &= \beta_0 + \beta_3 + e_9 \\
 y_{10} &= \beta_0 + \beta_3 + e_{10} \\
 y_{11} &= \beta_0 + \beta_3 + e_{11} \\
 y_{12} &= \beta_0 + \beta_3 + e_{12}
 \end{aligned}$$

1385 添字の対応に注意しながらみてくださいね。群 A の効果は β_1 、群 B の効果は β_2 、群 C の効果は β_3 になり
 1386 ます。

1387 この β それぞれを該当するところ (割り当てられた群) だけに対応させつつ、統一的表現形である
 1388 $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ にするには、次のように書く必要があります。

$$\begin{pmatrix} 3 \\ 3 \\ 4 \\ 4 \\ 6 \\ 7 \\ 8 \\ 9 \\ 7 \\ 6 \\ 5 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ e_7 \\ e_8 \\ e_9 \\ e_{10} \\ e_{11} \\ e_{12} \end{pmatrix}$$

1389 このような表記になった時の \mathbf{X} のことをとくに、実験のデザインを表している行列ということで、**デザイン行**
 1390 **列 (design matrix)** といいます。デザイン行列は自分で書くと面倒な感じがしますが、とにかくこのような
 1391 書き方で $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ の統一表現は可能です。

1392 ところでこのデザイン行列、 \mathbf{X} のサイズは今回 $n \times (m + 1)$ になっていますね (m は水準数)。2 水準の
 1393 ときは 2 列で済んだものが、3 水準になると 4 列になるのはおかしくないですか？そうです、1 つ大事なポ
 1394 イントを忘れていました。各群の値はベースライン β_0 からの相対的な違いです。相対的な違いというのは、
 1395 言い換えると $\sum \beta = 0$ 、すなわち全部の水準の和が 0 である必要があるのです。この式は今回の例だと
 1396 $\beta_1 + \beta_2 + \beta_3 = 0$ であり、これを移項すると明らかなように $\beta_3 = 0 - \beta_1 - \beta_2$ です。つまり総和が決まっ
 1397 ているので、自由に大きさを推定できるのは水準数 -1 になります。

1398 これを踏まえてデザイン行列を次のように書き換えることができます。

$$\begin{pmatrix} 3 \\ 3 \\ 4 \\ 4 \\ 6 \\ 7 \\ 8 \\ 9 \\ 7 \\ 6 \\ 5 \\ 4 \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \\ 1 & -1 & -1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ e_7 \\ e_8 \\ e_9 \\ e_{10} \\ e_{11} \\ e_{12} \end{pmatrix}$$

1399 この式は先ほどの式と内容的には同じで、表現の仕方が違うだけですが、 $\sum \beta = 0$ の条件がなければ計
 1400 算結果は算出されません。計算するための制約が少ないと、答えが出なくなるのです。 $\sum \beta = 0$ の制約条
 1401 件を別途書き加えてもいいですが、制約条件も含めた行列の書き方ができるといわけですね。

1402 線形モデルの統一的表现、あるいは行列の計算方法に少しは慣れてきたでしょうか。これがさらに多くの変
 1403 数、未知数を扱うことになると、その利点はよりはっきりしてきます*4。

1404 7.4 因子分析モデルの行列表現

1405 ということ、因子分析モデルの代数的表現ですが、これも行列を使って表現すると非常にシンプルに表現
 1406 できるということを説明していきましょう。

1407 内容はまったく同じですが、確認しておきましょう。標準得点行列 Z を因子負荷行列 A と因子得点行列
 1408 F をつかって、次のように表します。

$$Z = FA' + UD \quad (7.1)$$

1409 ここで、各行列の要素のサイズ感をつかんでおきましょう。まず R というのは相関行列ですから、 m 個項
 1410 目があるのでサイズは $m \times m$ の正方行列になります。次に F ですが、これは因子得点の行列です。得点は
 1411 人数分ありますから行は n 、因子の数が列になるのでこれを p とすると $n \times p$ です。 A は因子負荷行列。因
 1412 子負荷行列は因子の数と項目の数の組み合わせだけあるわけですから、 $m \times p$ になりますね。 U は独自因
 1413 子得点です。得点ですから人数分、独自性は各項目にありますから、サイズとしては $n \times m$ になります。最
 1414 後に D ですが、これは独自因子の負荷量です。項目の数だけあるのですが、列ベクトルや行ベクトルで表現
 1415 すると計算の時にサイズが変わって不便なことになります。ですから、対角項に d_j をもつ正方行列 $m \times m$
 1416 として表現しています。

1417 サイズを確認したところで、実際に行列計算をしてみましょう。

*4 ここでは触れませんが、被験者内計画・反復測定になっても行列表現はできます。個人差を表すデザイン行列を別途加えること
 になります。混合計画になると非常に複雑になりますが、それでも一般的な表現は可能です。

$$\begin{aligned}
R &= \frac{1}{N} Z' Z && Z \text{ を因子分析のモデル式にして} \\
&= \frac{1}{N} (FA' + UD)' (FA' + UD) && \text{前の項の転置を中に入れます} \\
&= \frac{1}{N} \{ (FA')' + (UD)' \} (FA' + UD) && \text{転置のカッコを外すときは順番を入れ替えて転置} \\
&= \frac{1}{N} (AF' + D'U') (FA' + UD) \\
&= \frac{1}{N} AF'FA' + \frac{1}{N} AF'UD + \frac{1}{N} D'U'FA' + \frac{1}{N} D'U'UD
\end{aligned} \tag{7.2}$$

1418 と、このように展開できました。記号を見ているとわかりにくいので、サイズ感を確認しましょう。最終的には、
1419 次のようになっています。

$$R = \frac{1}{N} \underset{m \times m}{A} \underset{m \times pp}{F'} \underset{pp \times nn}{F} \underset{nn \times pp}{A'} + \frac{1}{N} \underset{m \times pp}{A} \underset{pp \times nn}{F'} \underset{nn \times mm}{U} \underset{mm \times m}{D} + \frac{1}{N} \underset{m \times mm}{D'} \underset{mm \times nn}{U'} \underset{nn \times pp}{F} \underset{pp \times m}{A'} + \frac{1}{N} \underset{m \times mm}{D'} \underset{mm \times nn}{U'} \underset{nn \times mm}{U} \underset{mm \times m}{D}$$

1420 ここで、要素ごとに計算していた時のことを思い出してください。第二項 $\frac{1}{N} AF'UD'$ と第三項
1421 $D'U'FA'$ の中にある、 $F'U$ と $U'F$ のところは、共通因子得点と独自因子得点の積ですし、いずれ
1422 も標準化されていますから、 $\frac{1}{N}$ と合わせて考えると、これは相関係数を表していることになります。また、共
1423 通因子と独自因子は相関しませんので、これはイコール 0 となり、この 2 つの項が消えてしまうのです。
1424 また、第一項の $\frac{1}{N} F'F$ は、共通因子同士の相関を表しています。 $F'F = C$ とすると、これは因子得点
1425 間相関 C を表すことになります。これが直交であると仮定する、つまり他の因子と相関しないと考えると、
1426 $C = I$ 、つまり単位行列です。単位行列は計算に影響を与えませんから、 $AF'FA' = ACA' = AIA' =$
1427 AA' となり、この式は簡単に次のように変形できます。

$$R = AA' + D^2 \tag{7.3}$$

1428 先ほどの代数的展開を、そのまま行列で表現しただけですが、この方がシンプルに表現できていますね。この
1429 表現は、因子分析の第一定理と第二定理の両方を含んで一度に表せているのです。

1430 いかがでしょうか。行列表現の便利さがわかっていただければ、と思います。しかしこれでもまだ謎は残り
1431 ますね。我々が追っている謎は、行列からどのようにして因子負荷量を計算するのか、です。それを知るため
1432 には、もう 1 つ線形代数から明らかになる特徴を知らなければなりません。次回をどうぞお楽しみに。

1433 7.5 課題

1434 ■行列計算を確認しておこう $V = (I - \frac{1}{n}11')X$ の要素を書き下してみよう。平均偏差行列ができて
1435 いるでしょうか。

1436 ■行列計算を確認しておこう 2 S, Z, R も、面倒でも要素レベルまで書き下してみよう。

1437 ■行列計算を確認しておこう 3 因子分析モデルの行列計算の結果出てくる、 $R = AA' + D^2$ の要素
1438 を確認し、エレメントワイズで表現していた式との対応を確認しよう。

1439 第 8 章

1440 固有値と固有ベクトルと因子分析モデル 1441 の関係

1442 8.1 固有値と固有ベクトル

1443 今回は正方行列にみられるおもしろい特徴である、固有値 (eigenvalue) と固有ベクトル
1444 (eigenvector) についての話から始めます。ある正方行列 A , 列ベクトル x , スカラー λ が次のような関係
1445 にあった時, λ を固有値, x を固有ベクトルと言います。

$$Ax = \lambda x$$

1446 一見すると, x が両辺に入っていますから, A が λ に置き換わった等式に見えます。しかし一方は行列
1447 で, 他方はスカラーです。こんな奇妙なことが本当にあるのでしょうか? 具体的な数値例をみてみましょう。

$$A = \begin{pmatrix} 1 & 6 \\ 2 & 5 \end{pmatrix}, x = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

1448 を例にします。この時次の関係が成り立ちます。

$$Ax = \begin{pmatrix} 1 & 6 \\ 2 & 5 \end{pmatrix} \begin{pmatrix} 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 7 \\ 7 \end{pmatrix} = 7 \begin{pmatrix} 1 \\ 1 \end{pmatrix} = 7x$$

1449 確かに成立する組み合わせがありますね。この行列 A に対して, 7 が固有値, $(1, 1)$ が固有ベクトルに
1450 なっています。また, 実はこの行列 A については, -1 も固有値であり, そのときの固有ベクトルは $(-3, 1)$ も
1451 固有ベクトルです。

1452 この固有値分解こそ, 因子分析を元とする多変量解析の中心的な数学原理なのです。多変量解析の世界
1453 においては, 分散共分散行列やそれを標準化した相関行列など, 変数同士の関係を分析のスタートにおく
1454 でした。これらの行列は正方行列ですから, その固有値や固有ベクトルを計算することで正方行列の特徴を
1455 別の視点から分解して考えられるようになります。

1456 8.1.1 固有値の特徴

1457 この固有値の数学的特徴は色々あるのですが, データ分析をする上で重要な点を押さえておきましょう。

1458 固有値の特徴として, **固有値の総和が正方行列の対角要素の総和に合致する**, というのがあります。数式
1459 で表現すると, 次のようになります。

$$\sum_{i=1}^N \lambda_i = \text{trace}(\mathbf{A}) = \sum_{i=1}^N a_{ii}$$

1460 ここで a_{ij} は行列 \mathbf{A} の要素であり、 a_{ii} は i 行 i 列目、つまり対角要素です。この正方形行列 \mathbf{A} のサイズ
1461 は N で、対角要素の総和をとくに**トレース (trace)** といい $\text{trace}(\mathbf{A})$ と表します。それが固有値の総和と
1462 イコールになる、ということを表しています。サイズ N の正方形行列からは固有値が N 個算出できることがわ
1463 かっており、それをすべて足し合わせたものがトレースと同じになっているのですね。先ほどの例で言えば、
1464 \mathbf{A} のトレースは $1 + 5 = 6$ で、固有値の総和は $7 - 1 = 6$ であり、確かにこの関係が成立していることがわ
1465 かります。

1466 分散共分散行列のトレースは、分散の総和を意味します。項目同士の関係を表した行列であれば、分散は
1467 その項目から得られる情報の大きさであり、それを総和するということは、その調査研究・項目群から得られ
1468 る情報の総和であると言ってもいいでしょう。相関行列のトレースは、対角項に入っているのが $r_{ii} = 1.0$ で
1469 すから、項目の数と一致します。1つの項目の情報量を 1.0 に基準化して N 項目分の情報がある、というこ
1470 とを表しています。

1471 固有値と行列の関係は冒頭で示したように、 $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ であり、正方形行列の特徴をスカラーにしてしま
1472 うというものです。得られる N 個の固有値は、元の正方形行列のエッセンスをスカラーにして表現しているわけ
1473 です。

1474 ところで \mathbf{A} を n 次正方対称行列、つまり $n \times n$ サイズの対称行列だとすると、 n 個の固有値が求められ
1475 ます。これを $\lambda_1, \lambda_2, \dots, \lambda_n$ として、対応する固有ベクトルを $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ とします。ここで各固有ベクト
1476 ルのノルムが 1 であるとしましょう。行列と固有値・固有ベクトルの関係から、

$$\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i$$

1477 となりますが、このベクトルを並べた行列 $\mathbf{X} = (\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_n)$ を考えると、

$$\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{\Lambda}$$

1478 と書くことができます。ここで $\mathbf{\Lambda}$ は

$$\mathbf{\Lambda} = \begin{pmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_n \end{pmatrix}$$

1479 のような行列です。

1480 この両辺に \mathbf{X}' をかけると

$$\mathbf{A}\mathbf{X}\mathbf{X}' = \mathbf{X}\mathbf{\Lambda}\mathbf{X}'$$

1481 となりますが、固有ベクトルの性質とノルムを整えていることから $\mathbf{X}\mathbf{X}' = \mathbf{I}$ であり、そこから

$$\mathbf{A} = \mathbf{X}\mathbf{\Lambda}\mathbf{X}'$$

1482 と書くことができます。

1483 ここであらためて要素に注目すると、行列 \mathbf{A} が次のように分解されていることがわかります。

$$\mathbf{A} = \lambda_1\mathbf{x}_1\mathbf{x}_1' + \lambda_2\mathbf{x}_2\mathbf{x}_2' + \dots + \lambda_m\mathbf{x}_m\mathbf{x}_m' = \sum \lambda_i\mathbf{x}_i\mathbf{x}_i'$$

1484 となります。

1485 この分解例は 2×2 の簡単な例で確認しておきましょう。たとえば $A = \left(\begin{pmatrix} a \\ b \end{pmatrix}, \begin{pmatrix} c \\ d \end{pmatrix} \right) = \begin{pmatrix} a & c \\ b & d \end{pmatrix}$ とい

1486 う行列と、対角行列 $\Psi = \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix}$ があったとして、 $A\Psi A'$ の計算をしてみたいと思います。

$$\begin{aligned} A\Psi A' &= \begin{pmatrix} a & c \\ b & d \end{pmatrix} \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} \\ &= \begin{pmatrix} \alpha a & \beta c \\ \alpha b & \beta d \end{pmatrix} \begin{pmatrix} a & b \\ c & d \end{pmatrix} \\ &= \begin{pmatrix} \alpha a a + \beta c c & \alpha a b + \beta c d \\ \alpha a b + \beta c d & \alpha b b + \beta d d \end{pmatrix} \\ &= \alpha \begin{pmatrix} a a & a b \\ a b & b b \end{pmatrix} + \beta \begin{pmatrix} c c & c d \\ c d & d d \end{pmatrix} \\ &= \alpha \begin{pmatrix} a \\ b \end{pmatrix} (a \ b) + \beta \begin{pmatrix} c \\ d \end{pmatrix} (c \ d) \end{aligned}$$

1487 と、このようにスカラーとベクトルの積和の形に書き換えられるのですね*1。

1488 さて、これらをまとめて考えると、

- 1489 1. 固有値分解は行列を列ベクトルとその転置ベクトルの積の形に分解する。
- 1490 2. 固有値の総和は元の行列の対角要素の総和である。
- 1491 3. 元の行列の対角要素は各項目の分散を表している。

1492 ということです。固有ベクトルは全体の情報量をそのままに重要度の大きさに並べ替えたもの、固有値分
1493 解は行列をその要素の重要度ごとに分解していくことである、といえます。これこそ**因子分析**で取り出そうと
1494 している因子であり、固有値の大きさはその因子の重要度として、共通次元の判別（どこまで共通次元とみな
1495 すか）に使われるのです。

1496 8.2 固有値と固有ベクトルを求める

1497 ここで少し数学の方に話を戻して、固有値と固有ベクトルの計算方法を考えましょう。元の式を書き換えて
1498 次のような方程式を考えます。

$$(A - \lambda I)x = 0$$

1499 固有ベクトルは $x = \mathbf{0}$ すなわち全部ゼロであれば当然成り立ちますから（自明な解といいます）、これは除
1500 外することになります（ $x \neq \mathbf{0}$ ）。行列の表現は連立方程式の解を求めることと同じなものでした。 $A - \lambda I$ を連立
1501 方程式の係数行列だと考えれば、それが**逆行列**を持つと左辺にそれをかけてしまえば全部ゼロの答えになっ
1502 てしまいますから、そうでない答えを求めるには、この係数行列が逆行列を持たないことが重要です。

*1 ただし、この計算が可能なのは分解する元の行列が実対称行列だからです。実対称行列は固有値と固有ベクトルで対角化可能であることが証明できます。実対称行列の固有値は全て実数です。固有値が全て実数であれば適当な直交行列をつかって対角化でき、実対称行列の固有ベクトルは互いに独立するのでこれらを使って直交行列を作ることができるからです。これらの性質については線形代数のテキストなどの証明を参照してください。また計算プロセスからもわかるように、同じ要素を持つ列ベクトルと行ベクトルの積ですから、結果は対称になってしまうからです。

1503 さて、この授業の中では説明してきませんでしたが、方程式が解を持つかどうかを決定する計算方法があ
1504 ります。これを**行列式 (determinant)** といい*2、この値がゼロでなければその方程式は解を持つ、というこ
1505 とがわかっています。説明しなかったのは、この値を求める計算がとて面倒だからで、詳しくは線形代数の
1506 テキストにお任せするとして*3、ここでは簡便のために 2×2 方程式の行列について紹介します。

1507 2×2 の係数行列、 $\mathbf{P} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ の逆行列は次の式で求められることがわかっています。

$$\mathbf{P}^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix}$$

1508 この式から考えると、 $ad - bc$ のところが 0 になるとこの計算はできませんから、逆行列が存在しないこと
1509 になります。この $ad - bc$ にあたるところが行列式であり、 $|\mathbf{P}|$ とか $\det(\mathbf{P})$ のように表します。 $ad - bc$ が
1510 0 でなければ方程式は解けるのですが、今回の場合は解けると自明になってしまうので困ります。今回は
1511 $ad - bc = 0$ でなければならないのです。つまり一般的に書くと次のようになります。

$$|\mathbf{A} - \lambda \mathbf{I}| = 0$$

1512 $\mathbf{A} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$ とすると、この式は次のようになります。

$$\begin{vmatrix} a - \lambda & b \\ c & d - \lambda \end{vmatrix} = 0$$

1513

$$(a - \lambda)(d - \lambda) - bc = 0$$

1514 この方程式をとくに**固有方程式**といいます。これを解いてやれば良いことになりますね。具体的に
1515 $\begin{pmatrix} 1 & 6 \\ 2 & 5 \end{pmatrix}$ の例で計算してみましょう。

$$\begin{vmatrix} 1 - \lambda & 6 \\ 2 & 5 - \lambda \end{vmatrix} = 0$$

1516

$$(1 - \lambda)(5 - \lambda) - 12 = 0$$

1517

$$\lambda^2 - 6\lambda - 7 = 0$$

1518

$$(\lambda - 7)(\lambda + 1) = 0$$

1519 ここから $\lambda = 7, -1$ が得られますね。

1520 では固有ベクトルはどうなるでしょうか。固有値 7 の例で計算してみます。

$$\begin{pmatrix} 1 & 6 \\ 2 & 5 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix} = 7 \begin{pmatrix} x \\ y \end{pmatrix}$$

1521

$$x + 6y = 7x \rightarrow 6x = 6y \rightarrow x = y$$

1522

$$2x + 5y = 7y \rightarrow 2x = 2y \rightarrow x = y$$

1523 あれっ？ なんじゃこりゃ？ と思った人もいるかもしれませんが。でもこれ、間違いではありません。実は固有
1524 ベクトルは大きさが定まっておらず、今回の例で言えば $x = y$ つまり $x : y = 1 : 1$ の関係であればあらゆる

*2 行列式は数値であり、解が求まるかどうか決定 determinant する、という意味なのに、日本語訳はなぜか「式」といいます。変なの。

*3 たとえば村上他 (2016) の第 3 章をみてください。

1525 値が成立してしまうのです。固有ベクトルが (1, 1) でも (2, 2) でも (100, 100) でも、 $Ax = \lambda x$ の関係に影
 1526 響きませんから、普通はベクトルの長さ (ノルム (norm)) を 1.0 に規格化するという方策が取られます*4。
 1527 先ほど「固有ベクトルを適当な大きさに選んでやれば」というような表現をしましたが、それはベクトルの大き
 1528 きはいくらでもいいからできることなのです。

1529 さてここでみたように、 2×2 の方程式であれば固有方程式を解くことはできるのですが、行列のサイズが
 1530 どんどん大きくなると一般的に解けなくなっていくことは想像にかたくないと思います。実際我々は正方行列
 1531 として、項目の情報が詰まった分散共分散行列とか相関行列を使いますから、それが 2 項目しかないなんて
 1532 ことはなくて、もっともっと大きなサイズになります。そうすると計算機を使って近似的に答えを求めて行くこと
 1533 になります。

1534 8.3 固有値と固有ベクトルの幾何学的意味

1535 固有値、固有ベクトルについて、今度は違う側面から見直してみましよう。

1536 ある正方行列から固有値 λ と固有ベクトル \mathbf{a} が得られたとします。このベクトル \mathbf{a} のすべての要素を定数
 1537 c 倍したベクトル $\mathbf{b} = c\mathbf{a}$ を考えると、これもやはり同じ関係が成り立ちます。

$$A\mathbf{b} = \lambda c\mathbf{a} = \lambda \mathbf{b} \quad (8.1)$$

1538 先ほどの計算でも明らかになりましたが、固有ベクトルの値は絶対的なものではなく、要素間の相対的大きさ
 1539 を反映しているに過ぎないのでしたね。

1540 さてこれを幾何学的に、図形として考えてみましょう。要素が 2 つのベクトルは、2 次元座標に表現できま
 1541 す。ベクトル $\mathbf{x} = (x, y)$ という座標を表しているというわけです。固有ベクトルも要素が 2 つであれば、座標
 1542 で表現できます。先ほどの、要素を c 倍しても固有ベクトルとしての性質は変わらない、という話は、「固有ベ
 1543 クトルは大きさに意味はなく、方向を表したもの」ということになります。では何の方向を指し示しているので
 1544 しょうか。

1545 固有値と固有ベクトルの話の最初にあった、 $Ax = \lambda x$ というのを見直してみましよう。 x がなんらかの
 1546 座標を表していると考え、それに正方行列をかけるとはどういう意味でしょうか。次の計算式を見てくだ
 1547 さい。

$$A\mathbf{x} = \begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix} \begin{pmatrix} 1 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 6 \end{pmatrix} \quad (8.2)$$

1548 これをみると、座標 $\mathbf{x} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$ に A をかけたことで、座標が $\begin{pmatrix} 2 \\ 6 \end{pmatrix}$ に変わった、と見ることもできますね。こ
 1549 のように、ある座標が別の座標に移ることをとくに「変換」と呼びます*5。つまり正方行列はなんらかの変換を
 1550 施すものだ、と考えることができます。

1551 今回の例では、行列 A には次のような性質があります。

$$\begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix} = 2 \begin{pmatrix} 1 \\ 0 \end{pmatrix} \quad (8.3)$$

$$\begin{pmatrix} 2 & 0 \\ 0 & 3 \end{pmatrix} \begin{pmatrix} 0 \\ 1 \end{pmatrix} = 3 \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad (8.4)$$

*4 ノルムとは、要素の二乗和の平方根、 $\sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$ のことです。

*5 より一般的にいうと、以下ようになります。: 集合 X の各元 x に集合 Y の元 $f(x)$ を対応させる対応 f のことを、集合 X から集合 Y への写像 (mapping), 関数 (function), あるいは変換 (transformation) という。

1552 そう、お気づきのように、これは固有値・固有ベクトルです。この行列の固有値・固有ベクトルはそれぞれ
 1553 $\lambda_1 = 2, \boldsymbol{x}_1 = (1, 0), \lambda_2 = 3, \boldsymbol{x}_2 = (0, 1)$ であることがわかります。この固有値、固有ベクトルの組み合わ
 1554 せをじっくりとみていると、おもしろい特徴がわかって来ます。

1555 今回の行列から得られた 2 つの固有ベクトル、 $(1, 0)$ と $(0, 1)$ は、2 次元平面の単位ベクトルと呼ばれ
 1556 るものです。2 次元座標の任意の点は、これら 2 つのベクトルの任意の線型結合で表現できます。座標
 1557 (a, b) は $a \times (1, 0)$ と $b \times (0, 1)$ からなるベクトルですから、 $\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} a \\ b \end{pmatrix}$ というように、です。

1558 このように、単位ベクトルは 2 次元世界の基礎となる単位ともいべきもので、ここで $(1, 0)$ は x 座標の、
 1559 $(0, 1)$ は y 座標の基盤となるベクトルであるということが出来ます (これをとくに**基底**といいます)。つまり、
 1560 正方行列 A から得られる固有ベクトルは、その正方行列が作る空間の基盤を明らかにするものであったの
 1561 です。

1562 では固有値はどうでしょうか？ 今回は x 座標を 2 倍、y 座標を 3 倍に引き延ばす変換をしたわけですが、
 1563 この座標の歪み (重み) が固有値に対応していますね。つまり、固有値と固有ベクトルは新しい座標に変
 1564 換する、その変換先の空間的性質を表していることとなります。元の座標空間は $(1, 0), (0, 1)$ で作られる空
 1565 間ですが、変換先の空間は $(2, 0), (0, 3)$ で作られている空間、ということになります。

1566 つまり、正方行列は空間を変換するもの、あるいは正方行列の中に固有ベクトルを基底とした空間がある
 1567 もの、ということです。

1568 すべての正方行列に、こういった「変換」という解釈ができるのであれば、相関行列にも同様のことがいえ
 1569 るでしょう。相関行列は正方行列ですので、固有値分解できるのです。相関行列を固有値分解することは、相
 1570 関行列の中に潜む次元 (dimension) を抽出して行くことです。固有ベクトル (因子負荷行列) は、正方行列
 1571 によって変換される、変換先の単位ベクトルのことだったのです。そして、固有値はその次元のゆがみ (重み、
 1572 重要性) という意味があったのです。

1573 8.4 因子分析の数学的理解

1574 さあ因子分析に戻って考えてみましょう。因子分析のモデルは次のようなものでした。

$$R = AA' + D^2 \quad (8.5)$$

1575 ここで、左辺の正方行列を相関行列 R とし、 $\lambda \boldsymbol{x}\boldsymbol{x}' = \boldsymbol{a}\boldsymbol{a}'$ となるようにベクトルの大きさを整えてみましょ
 1576 う。これは λ を分解してベクトルの中に溶け込ませるようなものですから、 $\boldsymbol{a} = \sqrt{\lambda}\boldsymbol{x}$ とすればよいでしょう。
 1577 すると相関行列は次のように分解できます。

$$R = \boldsymbol{a}_1\boldsymbol{a}_1' + \boldsymbol{a}_2\boldsymbol{a}_2' + \cdots + \boldsymbol{a}_m\boldsymbol{a}_m' + \boldsymbol{d}\boldsymbol{d}' \quad (8.6)$$

$$R = \boldsymbol{a}_1\boldsymbol{a}_1' + \boldsymbol{a}_2\boldsymbol{a}_2' + \cdots + \boldsymbol{a}_m\boldsymbol{a}_m' + \boldsymbol{d}\boldsymbol{d}' \quad (8.7)$$

1579 ここでは共通因子の数が m 個だとわかっている体で分解していますが、基本的にはサイズ N の行列から
 1580 は N 個の固有値がずらーっと並ぶわけです。それをどこかで「共通しているのはここまで」と判断し、残りは
 1581 誤差であるとしてまとめて $\boldsymbol{d}\boldsymbol{d}'$ にしているだけです。このように数学的にはここからが共通因子、ここから
 1582 が独自因子といった区別をすることなく、最後のひとかけらまで固有値分解を行なっているのですが、その次
 1583 元の重要度でもって共通因子と誤差因子に (研究者が恣意的に) 分割しているのが因子分析のやっている
 1584 ことなのです。

1585 一般に、 N よりも m のほうがグッと少なくなります。たとえば YG 性格検査では $N = 120$ であり、 m は
 1586 せいぜい 5 から 10 数個です。120 項目のつくる 120 次元空間の中で、そこに働きかけても方向の変わらな
 1587 い基礎的な少数の次元にのみ注目すれば、効率よく情報圧縮ができるというものです。

1588 相関係数を固有値分解すると、その固有値はすべて足し合わせるとサイズ N になるのです。元のデータ
 1589 から計算される相関行列は、1つの項目が一単位分の情報を持っていると考えますが、固有値分解はそれ
 1590 を次元の重要性順に並べ替えます。固有値は項目いくつ分の重要度があるかということを表す指標だと考え
 1591 ることができます。どこから共通因子でどこからが誤差か、ということを考えるときに、たとえば固有値が 1.0
 1592 よりも小さくなるようであれば、項目 1 つ分の情報もないのだからということ誤差因子だと判断することが
 1593 あります。

1594 ところで因子分析モデルの A 、因子負荷行列ですが、これは共通因子の固有ベクトルをセットで扱ったも
 1595 のです。つまり、 $A = a_1, a_2, \dots, a_m$ と縦ベクトルを並べたものになっています。エレメントワイズで表現す
 1596 ると次のようになります。

$$A = \left(\left(\begin{array}{c} a_{11} \\ a_{12} \\ \vdots \\ a_{1N} \end{array} \right), \left(\begin{array}{c} a_{21} \\ a_{22} \\ \vdots \\ a_{2N} \end{array} \right), \dots, \left(\begin{array}{c} a_{m1} \\ a_{m2} \\ \vdots \\ a_{mN} \end{array} \right) \right)$$

1597 ここで AA' の間に単位行列 I を挟んでも、別に結果は変わりませんよね。単位行列はかけても変わらな
 1598 いのが特徴ですから、 $AA' = AIA$ です。ここでの I のサイズは $m \times m$ であることに注意しつつ聞いてく
 1599 ださい。

1600 $I = TT' = T'T$ になるような m 次の行列を使うと、 $AA' = ATT'A'$ の関係は保たれたままです。
 1601 この T はこれまた行列の座標を変えてしまう変換行列であり、これを挟むことができるということは**因子負
 1602 荷量の値はなんでもありだ**ということになってしまいます。この変換行列 T はとくに**回転行列 (rotation
 1603 matrix)**とも呼ばれますが、因子分析はこのようにどんな回転でもできる、**回転に関する不定性**があるの
 1604 です。あらゆる因子負荷量の組み合わせがあり得る、というのは困りものなので、なんらかの形で因子負荷行
 1605 列 A に制約をかける必要があります。言い方を変えれば、色々な制約の中ではあっても好きな値を取ること
 1606 ができますから、ユーザにとって便利な基準を考えてやれば良いでしょう。因子分析では一般に、**因子軸の回
 1607 転**を行います、それはこうした理由からです。

1608 ちなみに TT' に $\text{diag}(T\Phi T') = I_m$ という制約のある行列 Φ を挟んでやっても、元の計算モデルに影
 1609 響はありません。この Φ は**因子間相関 (factor correlations)**とよばれ、相関を持った因子軸の回転、す
 1610 なわち**斜交回転 (oblique rotation)**も色々なものが考えられています*6。

1611 ずいぶんややこしい話になってきました。この後は実践形式で因子分析を理解していければと思います。

1612 8.5 課題

1613 ■固有値問題 行列 $\begin{pmatrix} 8 & 1 \\ 4 & 5 \end{pmatrix}$ の固有値・固有ベクトルを求めよ。

*6 これに対して $TT' = I$ のような因子間相関がない (単位行列) な回転を**直交回転 (orthogonal rotation)** といいます。

1614 第9章

1615 Rをつかっての行列計算

1616 9.1 Rによる行列計算

1617 今回は統計環境 R をつかって、演習によってこれまで習ったことを整理していきましょう。

1618 9.1.1 環境の準備 (確認)

1619 まずは環境の準備です。すでに R や RStudio のインストールは終わっているものとして話を進めます*1。
1620 次の3つのステップをたどって、実行の準備をしてください。

1621 ■RStudio の起動 RStudio を起動してください。パッケージのインストールをするときなど、管理人権限が必要になるケースがあります。

1623 ■プロジェクトを開く RStudio を使う時の基本は、プロジェクトによる管理です。分析するデータ、テーマ、内容によってプロジェクトを切り替えながら使います。メニューバーから File > New Project と進んでください。すると「New Directory」「Existing Directory」「Version Control」の選択肢が出てきます。New Directory は新しいディレクトリ (=フォルダ) を作ってそこで関連ファイルをまとめますよ、という意味です。Existing Directory はすでに存在するディレクトリ (=フォルダ) をまとめる場所にしますよ、という意味です。みなさんの環境に応じて使い分けてください。すでにプロジェクトが存在する場合は、プロジェクトを開く (File→Open Project) から選択します。RStudio の右上などで、プロジェクトが開いていることを確認しましょう。

1631 ■R スクリプトを開く 今日のコードを書くための R スクリプトファイルを準備します。File > NewFile > R Script と進み、何も書いてない R スクリプトの画面を表示させてください。真っ白いファイルが開いたら問題ありません。

1634 これでスクリプトのところにコードを書いていく準備ができました。

1635 9.1.2 Rにおけるデータの型

1636 統計環境 R は Excel や Numbers や Calcs など*2の表計算ソフトとは違って、データを行列・ベクトルとして扱うことができます。たとえば次のコードを実行してみてください。

*1 まだだよー！という人は RStudio インストールといったキーワードで検索すると、親切にも案内してくれているサイトに色々出会えるでしょう。

*2 Excel は Microsoft Office に含まれる表計算ソフト。Numbers は Apple の Mac や iPad, iPhone で使える表計算ソフト、Calcs は Libre Office というフリーソフトウェアの表計算ソフトです。

code : 9.1 ベクトルデータの保持

```

1638 1 A <- 1:9
1639
1640 2 B <- c(3, 4, 5)
1641
1642 3 C <- matrix(A, ncol = 3)
1643
1644 4 D <- matrix(A, ncol = 3, byrow = T)

```

1644 ■コード解説

1645 1行目 1:9 をオブジェクト A に代入しています。ここでコロン: は連続した数字を表しており、
1646 c(1,2,3,4,5,6,7,8,9) と同じ意味です。

1647 2行目 要素 3, 4, 5 からなるベクトルを作って B に代入しています。c() は結合する (combine) という関
1648 数です。

1649 3行目 9つの要素があるベクトルを持ったオブジェクト A を matrix 型に変更しています。その時の列数
1650 は 3(ncol=3) です。

1651 4行目 同じくオブジェクト A を matrix 型に変換、列数は 3 ですが byrow = T で「行ごとに並べる」ス
1652 イッチをオン (TRUE) にしています。

1653 それぞれのオブジェクトに格納されているものを確認してみましょう。何がどうなっているかわかるでしょ
1654 うか (出力 9.1)。

R の出力 9.1: R におけるベクトルと行列

```

> A
[1] 1 2 3 4 5 6 7 8 9
> B
[1] 3 4 5
> C
      [,1] [,2] [,3]
[1,]    1    4    7
[2,]    2    5    8
[3,]    3    6    9
> D
      [,1] [,2] [,3]
[1,]    1    2    3
[2,]    4    5    6
[3,]    7    8    9

```

1655
1656 R のオブジェクトには色々な型がありますが、基本はベクトルです。数字のセットとして扱うのです。このと
1657 きの R のベクトルに行・列の区別はありません。行列であることを明示するためには、matrix という関数を
1658 当てはめることになります。ベクトルに行、列の意味を持たせれば、matrix 型で一行あるいは一列の行
1659 列である、という指定をしてください。

1660 これまで、あるいは他の授業で学んだことのある R のオブジェクトの形としては、list 型が多かったので
1661 はないかと思います。list という形は、要素がベクトルでも数字でも文字列でもなんでもかまわない、とい
1662 うものでした。これを矩形 (長方形) に整え、変数名としての列名や、行番号としての行名をもっているものが
1663 data.frame 型です。data.frame 型は list 型の特殊系なわけです。matrix 型は矩形のオブジェクト
1664 ではありますが、data.frame 型とは違います。matrix 型にも列名、行名をつけることはできますが、行列

1665 であると明示してあるように、計算するときには行列演算の対象になるのです。data.frame 型は行列の四則
1666 演算はできません。

1667 また一点注意が必要なこととして、R はベクトルの再利用をするということです。次の一行 (code:9.2) を
1668 実行して中身を見てみましょう。

code : 9.2 ベクトルデータの保持

```
1669 1      E <- matrix(A, ncol = 3, nrow = 6)
1670 2      G <- matrix(A, ncol = 3, nrow = 4)
1671
1672
```

1673 ■コード解説

1674 1 行目 1:9 をオブジェクト E に代入している。ここで列数は 3, 行数は 6 を指定している。

1675 2 行目 1:9 をオブジェクト G に代入している。ここで列数は 3, 行数は 4 を指定している。

R の出力 9.2: ベクトルの再利用

```
> E <- matrix(A, ncol = 3, nrow = 6)
> G <- matrix(A, ncol = 3, nrow = 4)
警告メッセージ:
matrix(A, ncol = 3, nrow = 4) で:
  データ長 [9] が行数 [4] を整数で割った、もしくは掛けた値ではありません
> E
      [,1] [,2] [,3]
[1,]    1    7    4
[2,]    2    8    5
[3,]    3    9    6
[4,]    4    1    7
[5,]    5    2    8
[6,]    6    3    9
> G
      [,1] [,2] [,3]
[1,]    1    5    9
[2,]    2    6    1
[3,]    3    7    2
[4,]    4    8    3
```

1676
1677 オブジェクト A の要素数は 9 です。しかし行列 E を作る時に、3 行 6 列、すなわち 18 の要素が必要
1678 であることを要求しました。そうすると R は (勝手に and/or 親切にも?) A を再利用して足りないところ
1679 を埋めてしまいます。今回はちょうど 2 回使えば埋まりましたので、エラーや警告もなくそのまま通っていま
1680 す。つくられた E の要素を確認し、再利用されていることをみてください。ちなみに G は 3 行 4 列、すな
1681 わち 12 の要素が必要であることを要求しました。そうすると R は警告を出して、再利用が途中で途切れ
1682 ますよと言ってくれます。とはいえ警告に過ぎないのでそのまま計算を続けることができ、作られた G には
1683 1, 2, 3, 4, 5, 6, 7, 8, 9, 1, 2, 3 と 2 周目の最初の 3 要素だけ使われています*3。

1684 このように R は、可能な時には良かれと思って勝手に再利用してしまいますから、注意してくださいね。

*3 そんなことよりオブジェクトが A,B,C ときてるんだから E のつぎは F だろ、と思う人もいるかもしれませんが。F でもいいのですが、実は大文字の F と T はそれぞれ FALSE/TRUE の略語であり予約後語です。上書きして使うこともできますが、なるべく避けたいほうがいいでしょう。ちなみに RStudio などこれら一文字を使うと予約語になっているので表示の時にハイライトされます。

1685 9.1.3 行列方と行列の演算

1686 それでは R を使ったの行列計算を進めましょう。

1687 加法減法はサイズが同じでないといけない、という制約はありますが、計算記号などに違いはありません。

1688 乗法は、スカラーとベクトル、スカラーと行列であれば普通に記述していただいて結構です。

code : 9.3 ベクトルの和・差, スカラーをかける

```
1689 1 x <- 1:3
1690 2 y <- 8:10
1691 3 x + y
1692 4 x - y
1693 5 2 * x
1694 6 y / 3
1695 7 A <- matrix(c(1, 2, 3, 4), ncol = 2)
1696 8 B <- matrix(c(5, 6, 7, 8), ncol = 2)
1697 9 A + B
1698 10 A * 3 + B * 2
1699
1700
```

1701 ■コード解説

1702 1行目 1:3 をオブジェクト x に代入している。

1703 2行目 1:9 をオブジェクト y に代入している。

1704 3行目 ベクトルの和 ($x + y$)

1705 4行目 ベクトルの差 ($x - y$)

1706 5行目 ベクトルにスカラーをかける ($2x$)

1707 6行目 ベクトルをスカラーで割る。スカラーの逆数をかけることと同じ。 ($y/3$)

1708 7行目 行列 A をつくる。サイズは 2×2

1709 8行目 行列 B をつくる。サイズは 2×2

1710 9行目 行列の和 ($A + B$)

1711 10行目 行列にスカラーをかけ、それを足し合わせる。スカラー倍は各要素にかかる ($3A + 2B$)

1712 出力結果はここでは提示しませんが、皆さんそれぞれ計算できているか確認しておいてください。

1713 続いて行列の乗法ですが、サイズが変わるような行列としての計算の場合ほとくに `%*%` という記号を使う
 1714 ことになります。ただのアスタリスク*ではなく、それを % で囲むことで行列の積になっていることを表現してい
 1715 るのです。また行列の行列を反転させるのは、**転置**という操作になりますが、これは `t()` という関数を使いま
 1716 す。行列の積を計算する例を `code:9.4` に示しました。計算結果が出力されますので、どういう計算をしてい
 1717 るか一行ずつ確認しながら進めてください。

code : 9.4 行列の積

```
1718 1 a <- c(1, 2, 1)
1719 2 b <- c(3, 4, 2)
1720 3 a * b
1721 4 a %*% b
1722 5 a %*% t(b)
1723 6 A <- matrix(1:9, ncol = 3)
1724 7 A * a
1725
```

```

1726 8 A %*% a
1727 9 a %*% A
1728 10 B <- matrix(1:6, nrow = 3, byrow = T)
1729 11 C <- matrix(c(1, 0, 0, 1, 1, 1), ncol = 3)
1730 12 B %*% C
1731 13 B %*% t(C)

```

1733 ■コード解説

1734 1行目 ベクトル $a = (1, 2, 1)$ を作る

1735 2行目 ベクトル $b = (3, 4, 2)$ を作る

1736 3行目 掛け算記号*を使っているが、これはベクトルの掛け算ではなく要素ごとの掛け算を意味する。

1737 4行目 ベクトルの掛け算記号 %*% を使っており、 ab の計算をしている。この時、ベクトル a は行ベクトル、
1738 b は列ベクトルと解釈されます。行列計算に適した形であり、デフォルトでは行方向が優位です。

1739 5行目 ベクトルの掛け算だが $t()$ で転置、すなわち b' を行っており、行列に適した形に変換されて
1740
$$\begin{pmatrix} 1 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 3 & 4 & 2 \end{pmatrix}$$
 と解釈されている。 3×1 と 1×3 の計算なのでできあがる行列は 3×3 のサイズ
1741 になる。

1742 6行目 行列 A を作る

1743 7行目 行列 A に対してベクトルをかけているように見えるが、ベクトルの掛け算でなく要素の掛け算記号*
1744 なので、各行を1倍、2倍、1倍する結果になっている。

1745 8行目 行列とベクトルの積 Aa であり、行列計算可能な 3×3 と 3×1 と解釈され、結果のサイズは 3×1
1746 になっている。

1747 9行目 ベクトルと行列の積 aA であり、行列計算可能な 1×3 と 3×3 と解釈され、結果のサイズは 1×3
1748 になっている。

1749 10-11行目 行列 B, C を作る

1750 12行目 行列の積 BC を計算している。 3×2 と 2×3 の行列の積なので、結果のサイズは 3×3 になっ
1751 ている。

1752 13行目 行列の積 BC' を計算しようとしているが、 3×2 と 3×2 の積は計算できないので、エラーが
1753 返ってくる。

1754 続いて**逆行列**の例を見てみましょう。逆行列の計算は、手計算でやると大変なのですが、R では関数
1755 `solve()` を使うと計算できます。

code : 9.5 逆行列の計算

```

1756 1 A <- matrix(c(2,1,5,3), ncol=2)
1757 2 solve(A)
1758 3 A %*% solve(A)
1759 1760

```

1761 ■コード解説

1762 1行目 行列 A をつくる

1763 2行目 逆行列 A^{-1} の計算

1764 3行目 $AA^{-1} = I$ より、逆行列になっていたことの確認。

1765 逆行列はかけると単位行列になるので、スカラーでいうところの逆数を意味します。逆数をかけると1にな
 1766 るように、行列の場合は逆行列をかけると単位行列 I になるのです。これの何が嬉しいかというと、行列が
 1767 連立方程式の係数を表していた場合、逆行列をかけることで連立方程式が解けるのです (セクション 6.3,
 1768 Pp.67 参照)。次の連立方程式を使って、実際に計算してみましょう。

$$\begin{cases} x - 2y - 5z = 3 \\ 5x + 4y + 3z = 1 \\ 3x + y - 3z = 6 \end{cases}$$

code : 9.6 連立方程式を解く

```
1769 1 A <- matrix(c(1,5,3,-2,4,1,-5,3,-3),ncol=3)
1770 2 b <- c(3,1,6)
1771 3 solve(A) %*% b
1772
1773
```

1774 1 行目 係数行列 A をつくる

1775 2 行目 右辺のベクトルを作る

1776 3 行目 $Ax = b$ から $A^{-1}Ax = A^{-1}b$ より、 $x = A^{-1}b$ として方程式の解を求める

1777 元の方程式との対応を確認しながら、「行列ってすごい、便利!」と思っていただければ幸いです。

1778 9.2 データの行列表現

1779 さて連立方程式が解けるのは嬉しいのですが、データ解析に直結するかと言われたら少し違いますね。
 1780 データの記述統計量など、数的処理をする時に行列を使う便利さを体感してみましょう。

表 9.1 投手のデータ

Name	team	height	weight	salary	Win	Save
菅野 智之	Giants	186	92	65000	14	0
西 勇輝	Tigers	181	82	20000	11	0
秋山 拓巳	Tigers	188	101	3200	11	0
大野 雄大	Dragons	183	83	13000	11	0
石川 柊太	Softbank	185	86	4800	11	0
千賀 滉大	Softbank	187	90	30000	11	0
涌井 秀章	Eagles	185	85	12500	11	0
森下 暢仁	Carp	180	76	1600	10	0
大貫 晋一	DeNA	181	73	2500	10	0
小川 泰弘	Swallows	171	80	9000	10	0
⋮	⋮	⋮	⋮	⋮	⋮	⋮

1781 表 9.1 に示したのは、データ解析基礎でも扱った 2000 年代の野球選手のデータです*4。それをプロジェ
 1782 クトフォルダに置き、code:9.6 のコードを実行してください。

*4 伴走サイト https://kosugitti.github.io/psychometrics_syllabus/ より、サンプルデータにある「野球選手のデータ 10 年度分」を選んでください。ファイルは UTF-8 形式で保存されています。

code : 9.7 データファイルの読み込み

```

1783 1 library(tidyverse)
1784 2 dataset <- read_csv("baseballDecade.csv") %>%
1785 3   dplyr::filter(Year=="2020年度") %>%
1786 4   dplyr::filter(position == "投手") %>%
1787 5   dplyr::select(Name, team, height, weight, salary, Win, Save) %>%
1788 6   na.omit() %>%
1789 7   arrange(-Win) %>%
1790 8   select(height, weight, salary) %>%
1791 9   as.matrix()
1792
1793

```

- 1794 1 行目 読み込み他, ファイル操作を便利にするパッケージ tidyverse を読み込みます*5。
- 1795 2-9 行目 データファイルを変形し, dataset オブジェクトに代入しています。ここで %>% はパイプ演算子
- 1796 と呼ばれるもので, 左の操作を右の操作へつなげる役割をします。
- 1797 3 行目 データを 2020 年度のものだけに絞るため, dplyr パッケージの filter 関数を適用して
- 1798 います。
- 1799 4 行目 データを投手のものだけに絞るため, dplyr パッケージの filter 関数を適用しています。
- 1800 5 行目 変数を必要なものだけに絞るため, dplyr パッケージの select 関数を適用しています。
- 1801 6 行目 データセットから欠損値を除いています。
- 1802 7 行目 表示用に, データを勝利数 (変数 Win) の多い順に並べ直しています。
- 1803 8 行目 以下の計算に使う変数だけに絞り込んでいます。
- 1804 9 行目 データセットを行列方に変換しています。

1805 ともかくこれで dataset オブジェクトは行列型の, 335 行 3 列の行列になっています。今から行う計算

1806 は, 第 7 講でやったように, データ行列 X から, 平均ベクトル m , 平均偏差行列 V , 分散共分散行列 S , 標

1807 準偏差を対角に持つ行列 Q , 標準化スコア行列 Z , 相関行列 R を作る, という手順です。数式とコードを示

1808 しますので, 確認しながら進めてください。

$$m = \frac{1}{n} X' \mathbf{1} \quad (9.1)$$

$$V = X - \mathbf{1} m' \quad (9.2)$$

$$S = \frac{1}{n} V' V \quad (9.3)$$

$$Z = V Q^{-1} \quad (9.4)$$

$$R = \frac{1}{n} Z' Z \quad (9.5)$$

code : 9.8 データの行列演算

```

1813 1 n <- nrow(dataset)
1814 2 one <- rep(1, n)
1815 3 m <- t(dataset) %*% one / n
1816 4 V <- dataset - one %*% t(m)
1817 5 S <- t(V) %*% V / n
1818

```

*5 tidyverse パッケージがない人はインストールしてください。関連パッケージをいろいろまとめたパッケージ群パッケージなので, 少し時間がかかります。

```

1819 6 SD <- diag(S) %>% sqrt()
1820 7 Q <- diag(SD)
1821 8 Z <- V %*% solve(Q)
1822 9 R <- t(Z) %*% Z / n
1823 10 cor(dataset)
1824

```

1825 1 行目 行数をオブジェクト n に入れる。行数を数える関数が `nrow` です。

1826 2 行目 データセットの列数と同じ長さの、1 だけを繰り返し入れたベクトルを作ります。繰り返しは `rep` 関数を使います。

1828 3 行目 平均ベクトル m を作る操作です。数式 9.1 に対応する操作です。

1829 4 行目 平均ベクトルを使って平均偏差行列 V を作る操作です。数式 9.2 に対応しています。

1830 5 行目 分散共分散行列 S を作る操作です。数式 9.3 に対応しています。

1831 6 行目 `diag` 関数を使って、行列 S の対角要素を抜き出し、その平方根を取り出しています。SD は標準偏差が入ったベクトルになります。

1833 7 行目 ベクトルを `diga` 関数にいれると、ベクトルの要素を対角に持つ正方行列ができます。それを Q として表現しています。

1835 8 行目 標準スコア行列 Z を作る操作です。平均偏差行列 V に Q^{-1} をかけています。数式 9.4 に対応しています。

1837 9 行目 相関行列 R を作る操作です。数式 9.5 に対応しています。

1838 10 行目 できあがった行列 R を検算するために、 R の持っている相関行列を作る関数 `cor` を実行しています。作った R と同じになっていることを確認してください。

1840 この一連のスク립トは、行列のサイズが変わっても適用できます。1 つの式表現で、どんなデータサイズでも一般的に表現できているところが行列表記のすごいところ。最終的に `cor` 関数があるのなら、こんな手間をかけなくてもいいじゃないかと思うかもしれません。しかし理屈を知っていて使うことと、理屈を知らずに使うことは違いますからね。

1844 9.3 R による固有値計算

1845 つづいて行列計算のとくにおもしろいところ、固有値計算をしたいと思います。

1846 固有値の計算は、小さなサイズであれば手計算でもできますが、大きなサイズになると n 次連立方程式を解くことになるのでとても大変です。解の公式で解くということもできませんので、近似計算手続きが必要です。その方法としてパワー法、ヤコビ法、ハウスホルダー法などいろいろ考えられていますが、数値計算の細かい理論に入っていくのは少し寄り道が過ぎますので割愛して、 R の関数を使って計算しましょう。

1850 先ほど計算した相関行列 R を使ってこれを確認します。

code : 9.9 相関行列の固有値分解

```

1851 1 eig <- eigen(R)
1852 2 eig$values
1853 3 sum(eig$values)
1854 4 sum(diag(R))
1855 5 eig$vector
1856 6 eig$vector[,1]^2 %>% sum
1857
1858

```

1859 固有値分解をする関数は `eigen` です。計算結果を `eig` オブジェクトに入れ、固有値 `values` と固有ベク

1860 トル vector を表示させてみました (出力 9.3)。固有値が大きい方から 1.7043160, 0.9486112, 0.3470728
 1861 となっています。ここでは変数が 3 つあって、それぞれの情報の大きさは相関行列で標準化されていて 1.0
 1862 です。この行列がもっている 3.0 の分散を、1.7 : 0.9 : 0.3 に再分配したことになります。ちなみに固有値の
 1863 総和 `sum(eig$values)` は、行列のトレース `sum(diag(R))` と等しくなっていることが確認できます。
 1864 また、固有ベクトルはそれぞれの固有値に対して計算されます。各列が各固有値に対応しており、第一固有
 1865 値 1.7 に対応する固有ベクトルは (0.68, 0.68, 0.26) です。固有ベクトルは向きだけあって大きさがありませ
 1866 んから、この数字はノルムすなわち二乗和したものが 1 になるようにして算出されているのがわかります。

R の出力 9.3: 固有値と固有ベクトル

```
> eig <- eigen(R)
> eig$values
[1] 1.7043160 0.9486112 0.3470728
> sum(eig$values)
[1] 3
> sum(diag(R))
[1] 3
> eig$vector
      [,1]      [,2]      [,3]
[1,] 0.6839162 -0.1734073  0.70865257
[2,] 0.6812174 -0.1959156 -0.70537929
[3,] 0.2611541  0.9651668 -0.01586178
> eig$vector[,1]^2 %>% sum
[1] 1
```

1867

1868 このような数値計算を駆使して、因子分析や回帰分析が実際に計算されているのですね。当然ですが、手
 1869 計算よりも機械で計算させた方が圧倒的に早いですね。計算機の発展スピードのおかげで、今は何万、何十
 1870 万というデータでも扱うことができる用意になりました。またデータのサイズが変わってもスクリプトを変える必
 1871 要がありません。

1872 こうした原理を知っておくと、ビッグデータなど恐るに足らず、というところではないでしょうか。

9.4 課題

1873

1874 ■線形代数の練習問題 行列 $A = \begin{pmatrix} 1 & 3 \\ 2 & 4 \end{pmatrix}$, $B = \begin{pmatrix} 3 & 1 & 5 \\ 2 & 4 & 6 \end{pmatrix}$, $C = \begin{pmatrix} 4 & 1 \\ 2 & 3 \end{pmatrix}$, 列ベクトル $x = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$,

1875 行ベクトル $y = (2 \ 8)$ とするとき、次の計算を R で実行するコードを書きなさい。なお、計算できないもの
 1876 については「計算できない」としてコードは書かないこと。

- 1877 1. $A + B$
- 1878 2. $A - C$
- 1879 3. AB
- 1880 4. AC
- 1881 5. $B'A$
- 1882 6. Ay'
- 1883 7. xy

1884 8. $\mathbf{x}B$

1885 9. $\mathbf{x}'B'$

1886 10. $\mathbf{y}x$

1887 ■連立方程式を解く 次の連立方程式を R で解きなさい。

$$\begin{cases} x - 2y + 3z = 1 \\ 3x + y - 5z = -4 \\ -2x + 6y - 9z = -2 \end{cases}$$

第 10 章

R を使った因子分析と尺度作成法

前回に続いて、R を使った演習を進めましょう。今回は R を使った因子分析および尺度作成です。

10.1 調査研究の手順

まずは心理尺度作成の手順を大まかに理解しておきましょう。

■構成概念の設定 心理尺度を作り始めるにあたって、まずは何を測りたいのかを明確にする必要があります。構成概念妥当性 (Construct Validity) についてのそもそも論ですね。ある概念を測定したいとして、そういう概念が本当に存在するのか、どこの誰がどのように持つ概念なのかを明確にする必要があります。たとえば文部科学省がスローガンのように掲げる生きる力というのを測ってみたいと思っても、それは何なのかわかりません。死んでなければ生きていますから、生きる力はあるように思えます。でもそういうことじゃないらしい。じゃあどうということなのか、ということを考えて行かなければなりません。あるいは、みなさんはいちびりという関西弁をご存知ですか。いたずらっ子、わざと変なことをして注目を引きたい子、のような意味なのですが、このいちびり感は関西の人間にしかわからない感覚かもしれません。であれば調査対象者が限定的な、一般的ではない概念ですから、因子分析の行う多くの人の反応から共通成分を抜き出すという考え方には適さない概念です。より一般的なものに考え直す必要があるでしょう。

ほかにも、心理学的な態度なのか、パーソナリティのような行動の傾向なのか、抑うつ傾向のように心身両方の問題なのか、そしてそれらが紙とペンで測定可能なものなのかどうか、改めて考えましょう。何を当たり前のことを、と思うかもしれませんが、意外とおかしな問題設定に入り込んでないとも限らないのです。また、心理尺度や心理学的な概念はすでにさまざまなものが存在します。これらを見比べて、自分の測定したい概念が他の概念とどのように同じで、どのように違うのかを論理的に考えられなければなりません。どこかにある尺度とほとんど同じであればオリジナリティが存在しないだけでなく、車輪の再発明というムダな努力をすることになります。完全にオリジナルなものを考えるのは難しいかもしれませんが、新しい概念を使うことによってこれまでの概念で説明できたことを含み、さらにこれまでの概念で説明できなかったことも説明できるようになる、という利点が必要です。これらは弁別的妥当性 (distinctive validity) にも関わる問題です。

■項目の選出 測定したい概念が明確になれば、これに関わる項目を作り出す必要があります。ある態度を強く持っている人はどのような振る舞いをし、どのような振る舞いをしないのか。どのような意見に賛成し、どのような意見には反対するのか、といったことをさまざまな角度から検証します。「目に見えないものを測定する」のが心理測定であり、質問紙調査というのはこの見えない的に向かって項目という矢を射掛けるようなものです。矢の数はなるべく多く、また人は嘘をついたり間違えたりする生き物ですから多角的に聞くことで、捉えるべき本質に迫っていく必要があります。言葉を使って聞くわけですから、古すぎる・新しすぎる表現は避

1918 け、専門用語は使わずなるべく平易な言葉で聞く必要があります。ワーディング (言い回し) についても細部
1919 まで注意しましょう*1。

1920 ■予備調査の実施 ある程度項目が集まれば、予備調査を行います。用意した項目の 5 倍から 10 倍の人
1921 を集めるのが良い、と一般的に言われています (Grimm and Yarnold, 1994)。ここでの予備調査項目は、
1922 十分考えられたものではあると思いますが、分析してみると意外と不適切な項目というのが出てくるものです
1923 ので、多めに用意されていることでしょう。必然的に、予備調査の対象サイズも百人以上の大きなものになる
1924 のが一般的です。

1925 ■探索的因子分析 この後詳しく説明しますが、ここが尺度作成に関わる分析のメインです。ここで因子数を
1926 決める (あるいは確認する) ことになり、因子的な妥当性みたり、不適切な項目は除外したりします。別の聞き
1927 方の方が良いような項目があれば、項目の入れ替えを行なって再調査することもあります。この探索的な因
1928 子分析と予備調査を繰り返して、何度行っても同じ因子構造になり、同じ項目が同じ因子に含まれるというの
1929 が確認できれば、項目セットは完成すると言えるでしょう。

1930 ■項目反応理論 ある因子に関連する項目とそのデータセットを抜き出して、項目反応理論 (段階反応モデ
1931 ル) を実行しましょう。なぜ一部を抜き出すかという、項目反応理論モデルは 1 次元性を仮定したモデルで
1932 あることが基本だからです。これを多次元に展開した多次元段階反応モデルもありますので、直接そちらを
1933 使っても構いません。ともかくこれを行うことで、反応段階数を確認でき、また項目情報曲線やテスト上表曲
1934 線を書くことで、この尺度がどのような領域を得意とする項目群からなるのかを記述できることとなります。こ
1935 れに関しては次回より詳しく説明します。

1936 ■本調査へ 最終的に項目群が確定したら、大規模な調査を行ってテストの標準化を目指します。因子構造
1937 などはほぼ確定しているはずですから、このテストを使うとどのようなデータの散らばりが得られるのかを確
1938 定するわけです。この最後のステップの目的は標準化、つまりこのテストで測定される人の平均や散らばり、
1939 分布の形状を確認することにあります。使うときに逐一分析しなくてもいいように、項目回答パターンがどれ
1940 ぐらいであれば、全体のどのあたりに位置するかといったプロフィールが作れるとなおいいでしょう。

1941 以上が大体の流れになります。ここにあるように、予備調査を何度も繰り返して因子構造を確定し、そこか
1942 ら本調査を行うことで、ほぼこの尺度はどれぐらいの点数で分布するか、といったことがわかるようになります。
1943 かつては因子分析を 1 回実行するだけでも多大な時間がかかりましたから、心理尺度を作るというのは
1944 それだけでライフワークになるような作業でした。幸い最近では分析のスピードが飛躍的に向上しましたので、
1945 簡単に分析してやり直すということが出来ます。しかしだからといって、構成概念妥当性を疎かにしてはいけ
1946 ませんし、「やったらこうなった」というようなやり逃げ研究になるのではなく、使うための尺度を作るのだとい
1947 うことは忘れてはいけません。

1948 10.2 共通性の推定

1949 それでは実際にデータを使ってのやり方を説明します。改めて、因子分析モデルについて確認しておきま
1950 しょう。正方行列を相関行列 R とし、 $\lambda_1 \mathbf{x}\mathbf{x}' = \mathbf{a}\mathbf{a}'$ となるようにベクトルの大きさが整えられているとすれ
1951 ば、次のように表現できるのでした。

$$R = \mathbf{a}_1\mathbf{a}_1' + \mathbf{a}_2\mathbf{a}_2' + \cdots + \mathbf{a}_m\mathbf{a}_m' + \mathbf{d}\mathbf{d}' \quad (10.1)$$

*1 たとえば「テレビやラジオをよく見聞きますか」という質問は悪い質問です。テレビしか見ない人、ラジオしか聞かない人がどう答えていいかわからず。これはダブルバーレルと呼ばれる悪手ですが、ほかにも色々ありますので調査法の専門書を参考にしてみてください。

このことから、相関行列を固有値分解すれば共通因子負荷量が得られることがわかります。しかしこれには1つ問題があります。この式の最後の項目、 dd' がわからないのです。行列の形で書くと因子分析モデルは次のようになります。

$$R = AA' + D^2$$

ここで D は対角項に d_j^2 をもつ対角行列です。共通因子を得るために分解するべき行列は、次のようになります。

$$R\ddagger = R - D^2 = AA'$$

R は相関行列であり、その対角は 1.0 ですが、分解すべき行列 $R\ddagger$ は対角項に $1 - d_j^2 = h_j^2$ が入った行列でなければなりません。それを固有値分解してはじめて、共通因子成分が得られるわけです。そしてこの d_j^2 の大きさは、事前にわかっていないので、計算の最初にはなんとかしてこれを埋める必要があります。これを **共通性の推定問題** といいます。

実際の計算においては色々な方法が考えられてきています。推定せずに 1.0 で分析しちゃう方法、その行における相関係数の最大値を入れる方法、その項目を他の項目から回帰分析した時の R^2 値を入れる方法 (SMC 法) などです。最初の「共通性を推定しない」という方法は、**主成分法 (principle component method)** と呼ばれます。**主成分分析** という分析法と同じやり方だからです。第二、第三の方法は、ひとまずその値で計算しますが、その後固有値分解をした行列から相関行列を再構成する、つまり $R\ddagger \rightarrow AA' \rightarrow R\ddagger_2 \rightarrow AA' \rightarrow \dots$ という繰り返しをおこない、数字が変化しなくなるまで反復するという方法で、**主因子法** と呼ばれます*2。この主因子法と同じ結果になるとされているのが、**最小二乗法** です。回帰分析の時に聞いたことのある方法ですね。それと同じ考え方で、因子数が決まっていれば、因子構造の形とデータとの誤差が最も小さくなるように共通性を推定するという方法です。原理的に主因子法と最小二乗法は同じ結果になるとされています。最小二乗法が出てきたらひょっとして、と思うかもしれませんが、そうその通りで、**最尤法** もあります。得られたデータが多次元正規分布からのサンプルだと考え、多次元正規分布の形にぴったりフィットするように推定パラメータを定めていく方法です。

このように、探索的に因子分析をする場合は、共通性をどのように推定するかということを指定しなければなりません。

10.3 因子数の決定

先ほどの最小二乗法、最尤法の説明の時に「因子数が決まっていれば」とありました。そうです、探索的な方法の場合は、因子分析をするときに因子数を決めて、何因子の答えを出すかを定めてやらなければなりません。それには色々な方法が考えられています。ここでは古典的な方法 4 つと、洗練された方法 1 つを紹介します。

■スクリープロットを見る 古典的な方法その 1 です。「見る」と書いてあるように、目で見て判断します。**スクリープロット (scree plot)** とは、固有値を大きい順にならべ、横軸に大きさの順番、立て軸に固有値の値をとった折れ線グラフのことです。これを見て、大きな変化が見られたところを基準とし、因子の数を決めます。

■固有値 1.0 という基準 固有値を算出し、1.0 以上であれば共通因子、それ未満であれば誤差因子とまとめる基準のことを言います。1.0 というのは項目の分散の大きさであり、固有値はこの分散の大きさを再配

*2 反復せずに最初に推定した初期値でいく方法もあります。これは反復しない主因子法ですが、計算機能力が十分高い今では反復するのが基本であり、とくに断りなく主因子法とあれば反復して求めていると考えていいでしょう。

1986 分する手続きなのでした。再配分した結果、1 項目分も情報が無いような次元は誤差みたいなものでしょ、と
1987 いう判断をすることになります。

1988 ■**累積寄与率をみる** 固有値の大きさを項目数で割ることで、1 つの因子が全体分散の何 % を説明す
1989 るかを算出できます。これをとくに**寄与率 (contribution ratio)** といいます。これをどんどん積み重ねて
1990 いて (累積), 累積寄与率が 50% を超えるところまでを共通成分とみなす, という方法です。因子分析は項
1991 目情報の圧縮をしていることでもあり, 共通因子に入れないものはゴミとして捨てるようなものですが, 全体
1992 の半分以上をゴミとして捨てたらいかんでしょ, という基準です。

1993 ■**解釈可能性を考える** 因子数に迷った時は, とりあえず何パターンかで分析してみて, 最も解釈しやすい
1994 数で OK にしましょう, というやり方です。なんと主観的な方法でしょうか。数学的な基準ではなく, 研究者の
1995 カンに頼る方法であり, 客観性や再現性の問題を考えると容認し難い方法ではありますが, 実際にはよく使
1996 われているやり方です。

1997 ■**並行分析による方法** ここまでの方法は幾分古く, 見た目や感覚にたよるものでしたが, **並行分析**
1998 (**parallel analysis**) は数学的な基準で分析を行うものです。これは分析するデータセットと同じサイズの
1999 乱数を発生させ, その固有値構造とデータの固有値構造を比較するというものです。乱数で作られる構造
2000 は, 心理的な反応パターンに依らないのですから, 当然無意味な因子ばかりが出てきます。実際のデータに
2001 基づく相関係数が持つ構造は, 心理的なパターンを反映しているのですから, それとは違う有意なパター
2002 ンになっているはず。この有意性, 無意味性を比較して, 乱数で作る因子以上の意味ある次元がでて
2003 いれば, 共通因子として採用するという方法です。

2004 以上, 5 つほど説明してきましたが, これについては実際に見てもらったほうが早いと思いますので, 今か
2005 ら演習を始めたいと思います。

2006 10.4 探索的因子分析の実際

2007 10.4.1 環境の準備 (確認)

2008 まずは環境の準備です。すでに R や RStudio のインストールは終わっているものとして話を進めます。い
2009 つもの 3 つのステップをたどって, 実行の準備をしてください。

2010 ■**RStudio の起動** RStudio を起動してください。パッケージのインストールをするときなど, 管理人権限が
2011 必要になるケースがあります。

2012 ■**プロジェクトを開く** プロジェクトは前回作成した物で良いと思います。同じフォルダの中で, スクリプト
2013 ファイルだけ変えれば良いでしょう。メニューからプロジェクトを開き (File→Open Project), RStudio の
2014 右上などで, プロジェクトが開いていることを確認しましょう。

2015 ■**R スクリプトを開く** 今日のコードを書くための R スクリプトファイルを準備します。File > NewFile
2016 > R Script と進み, 何も書いてない R スクリプトの画面を表示させてください。真っ白いファイルが開いた
2017 ら問題ありません。

2018 また今日は psych パッケージ (Revelle, 2021) を用います。準備がまだの人は Package タブからインス
2019 トールするか, コンソールで `install.packages("psych")` と入力してパッケージを導入しておいてくだ
2020 さい。

2021 データセットも psych パッケージが持っているものを使います。

code : 10.1 サンプルデータを使う

```
2022 1 rm(list=ls())
2023 2 library(tidyverse)
2024 3 library(psych)
2025 4 dat <- psych::bfi %>% dplyr::select(-gender,-education,-age)
2026 5 dat
2027 6 help(bfi)
2028
2029
```

■コード解説

1 行目 環境の初期化

2-3 行目 パッケージの読み込み

4-5 行目 サンプルデータ bfi を使う。性別、教育、年齢の情報は因子分析には使わないので除外する。

6 行目 データセットのヘルプを表示

ヘルプ画面にあるように、これは 25 のパーソナリティについての調査項目からなる、2800 人分のデータです。A(Agreeableness, 同調性), C(Conscientiousness, 几帳面さ), E(Extraversion, 外向性), N(Neuroticism, 神経質さ), O(Openness, 開放性) という Big 5 の各要素に対応した項目が 5 つずつあります。加えて性別 (gender), 最終学歴 (education), 年齢 (age) といった 3 つの変数も入っています。今回使うのは、5 つの因子それぞれに対応すると考えられる各 5 つの項目、計 25 項目分です。

10.4.2 並行分析による因子数の決定

code : 10.2 データの構造を見る

```
2041 1 corMat <- cor(dat,use="pairwise")
2042 2 corMat
2043 3 eigen(corMat)
2044 4 fa.parallel(dat)
2045
2046
```

■コード解説

1-2 行目 データの相関行列を計算し、表示させる

4 行目 相関行列の固有値分解を実行

5 行目 並行分析の実行

因子分析は相関行列から分析を始めますので、まずその計算をしました。オプション use="pairwise" は欠損値が含まれるデータセットに対し、ペアで残っている箇所は使って分析するという指定です*3。4 行目で固有値分解を試みました。結果として出力 10.1 のようなものが示されています。

*3 これがなかったら、欠損値があるので計算できません、で終わります。

R の出力 10.1: 固有値の推移

```
eigen() decomposition
$values
[1] 5.0369025 2.7440855 2.1076322 1.8318415 1.5356864 1.1131589 0.8462367 0.8114075
[9] 0.7349482 0.6956449 0.6810036 0.6571432 0.6281130 0.5964129 0.5626284 0.5405207
[17] 0.5236707 0.4984540 0.4899130 0.4549419 0.4328324 0.4092023 0.4071932 0.3852111
[25] 0.2752152
```

2054

2055 この下に固有ベクトル (\$values) も続いているのですが、長くなるので省略しました。ここで固有値が 25 個算
 2056 出されていることがわかります。また、これらを総和すると 25 になります ($tr(\mathbf{R}) = \sum_i = 1^{25} r_{ii}$)。このよう
 2057 うに標準化されたデータの分散を再構築し、第 1 固有値は 5.03, すなわち項目 5 個分の情報を持っているよ
 2058 うになったのです。第 2 固有値が 2.744, 第 3 固有値が 2.107, 第 4 固有値が 1.183... と続きます。また、こ
 2059 れを相対比率に変えると、 $5.03/25 = 20.12$ ですから、第 1 固有値だけで全体の 20% を説明できることにな
 2060 ります。第 2 固有値は 2.744 ですから、 $2.744/25=0.1096$ で全体の 10% ほど、第 1 固有値と合わせて
 2061 31.12% の累積寄与率があることになります。累積寄与率は第 5 固有値までで 53.02% になりますから、25
 2062 項目すべてを使わずとも 5 因子だけで情報の半分は説明できることになります！

2063 また、第 7 固有値の大きさは 0.846 であり、それ以降の固有値はすべて 1.0 未満、すなわち項目 1 つ分も
 2064 情報を持っていない因子ということになります。であれば共通因子として掬い上げる必要はなく、少なくとも
 2065 7 番目以降 (基準 2), あるいは 6 番目以降 (基準 3) は誤差と考えてまとめてしまってもいいかもしれませ
 2066 んね。

2067 続く `fa.parallel` 関数は並行分析を行う関数で、同時にスクリープロットも表示してくれます (図
 2068 10.1)。

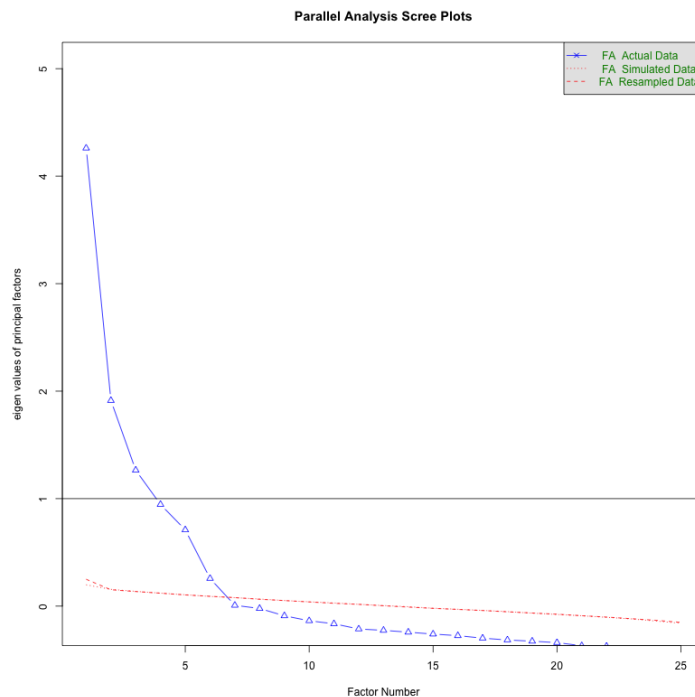


図 10.1 スクリープロットと並行分析

この図は先ほど算出されたような固有値^{*4}を大きい順に並べ、線で繋いだものが示されています。青い線が実データ、赤い点線が乱数によるにデータサイズが同じだけの偽データです。基準 2 に該当する 1.0 のところで線が引かれていて、これよりも大きい固有値を持つのが 3 つあることがわかります。また、第 1 固有値から第 2 固有値にかけて大きな傾斜があります (4.26063957 → 1.91279965)。最初のうちはまあ当然として、第 5 と第 6 の間も少し大きなギャップがあるようです (0.70938377 → 0.25749866)。それ以降はだらだらと数値が減衰していますね。このような大きなギャップのところは、大きな意味の違いがあると考えてもいいでしょう。基準 1 の観点からいくと、第 5 固有値までを共通因子とし、第 6 因子以降は誤差としてまとめてしまってもいいでしょう。また、並行分析の観点では、赤い線よりも上にある青いラインは有意義ということですから、6 因子構造がいい、ということになりますね。

このように、さまざまな基準で考えても、因子数は 5 か 6 のどちらかになって決定打に欠けます。こういうときは解釈しやすい 5 因子にしよう (基準 4)、ということになったりします。

ちなみにこの固有値構造、最終的には負の値が出ていますので、あまり良いデータとは言えません。25 項目を用意していたのに、8 因子以降は 0 以下、つまり情報がないようなものです。これは項目同士が似かよっていて、あまり違った情報を持ってこなかったということでもあります。実際のデータを分析する時は、こうした点にも注意が必要です。

10.4.3 因子分析の実行

それでは因子数を 5 と定めて因子分析を実行してみましょう。

code : 10.3 因子分析の実行

```
1 result <- fa(dat, nfactors = 5, fm = "ml", rotate = "geominQ")
2 print(result, sort = T)
```

これは 1 行目で実行し、2 行目で表示しているだけですが、内容を少し説明します。psych パッケージの持つ fa 関数は、データの他に因子数 nfactors、共通性の推定方法 fm、回転方法 rotate などをオプションで定めることができます。共通性の推定方法は、ここでは ML(最尤法) を選択しました。回転方法とは因子軸の回転 (rotation of factor axis) について定めるところです。これについては少し説明が必要ですね。

■因子軸の回転 因子分析は、相関行列という項目間関係が作り出す空間に、その基礎となる軸を見出すものでした。これを図 10.2 のように表現してみます。項目同士はその相関関係から空間を作っていますが、ここでは 2 因子、F1 と F2 という軸をもとに座標のように書いています。項目同士が近くにあるのは、相関あっていることを意味します^{*5}。ここで項目 1 の座標を見ると、第一因子 (F1) に 0.3、第二因子 (F2) に 0.6 とありますが、これが項目 1 の因子に対する因子負荷量ということになります。

因子分析では因子負荷量をもとに、「この項目と関係の深い因子はどれか」「因子はどのような項目から構成されているか」ということを考えていくのですが、項目 1 は第一、第二因子それぞれにちよとずつ関係していることがわかります。そのほかの項目も、2 つの因子に負荷していますがなかなか決め手がない、というところですね。

そこで座標を回転させます。どの項目がどの因子に関係しているのかを、よりはっきりさせるように軸をグルリと回すのが因子軸回転の目的です。図 10.3 のように回転させますと、回転後の軸 ($F1'$, $F2'$) に下ろした垂線が座標、すなわち因子負荷量になりますから、項目 1 は第一因子 ($F1'$) に 0.9、第二因子 ($F2'$) に

*4 正確には R^+ の固有値です。

*5 原点からのベクトルと考えたときに、2 つの項目ベクトルが作る角度 θ のコサインを取った $\cos \theta$ が相関係数になります。

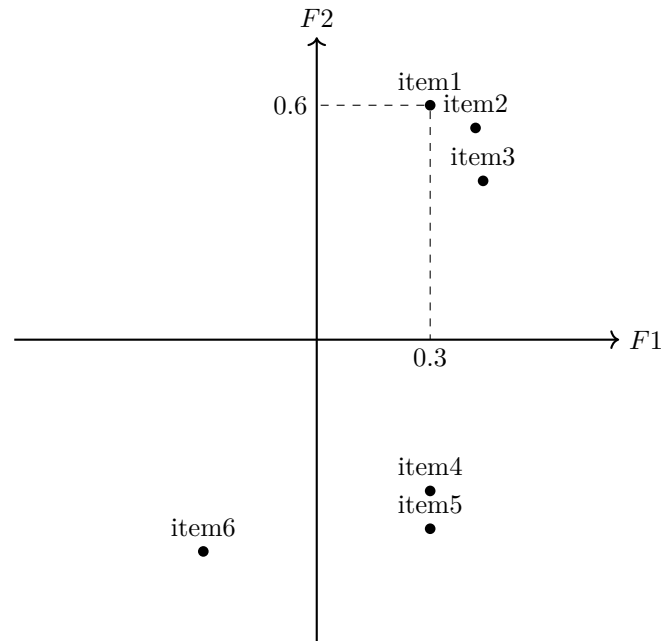


図 10.2 項目の空間にある基底

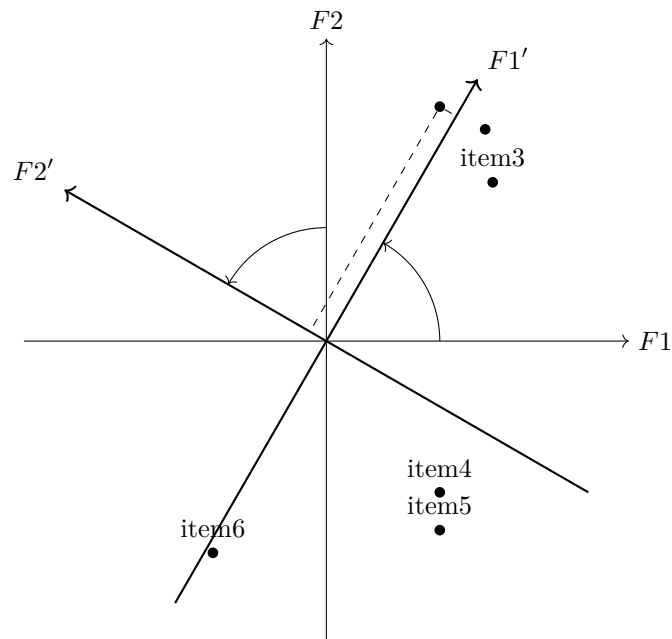


図 10.3 因子軸の回転。60 度回転させてみました。

2107 -0.07 になりました。こうしてみると、項目 1 は第一因子と非常に深く関係しているし、第二因子とはほとんど
 2108 関係がない、とみることができます。ほかの項目も一方の因子に大きく負荷し、他方の因子にほとんど負荷し
 2109 ていないようになりますから、どの因子がどういう項目と関係するかを評価するのにメリハリが効いてやりや
 2110 すくなります。このように、利用しやすく座標を変えるのが因子軸の回転、という技なのです。

2111 ここで最初の因子軸 $F1, F2$ は直交していましたね。座標ですから当然です。直交している、すなわち因子
 2112 間相関がないように回転することを直交回転 (orthogonal rotation) といいます。

2113 しかし、因子軸の回転の目的は、わかりやすくするためなのですから、いっそ図 10.4 のように、角度をつけ
2114 た方がさらにうまくメリハリをつけることができますね。この図では、項目 4,5 が明らかに第二因子に影響して
いる、ということが見て取れるようになりました。

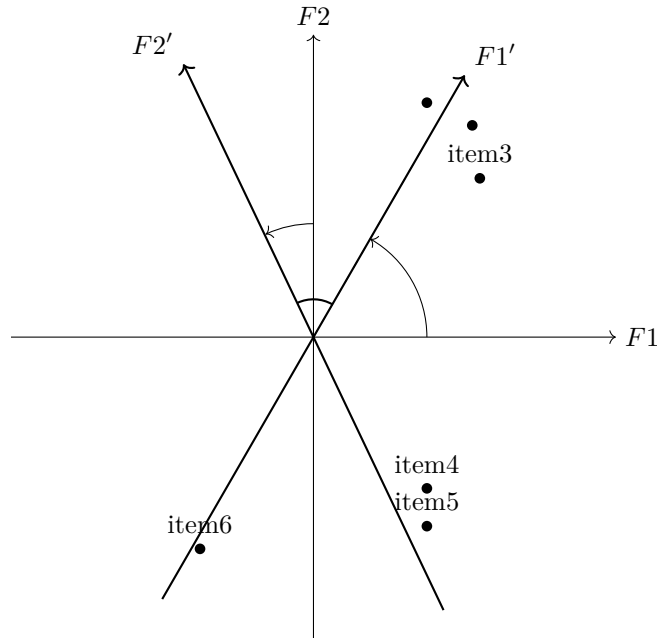


図 10.4 角度をつけて回転させる斜交回転の例。角度があるとは、因子同士が相関するという。

2115
2116 このように、因子軸が直角で交わっていない場合の回転方法を**斜交回転 (oblique rotation)** といいま
2117 す。因子軸の回転は一般に、斜交回転した方がメリハリが効いて因子を解釈しやすくなります。斜交回転した
2118 とき算出される因子同士の相関を考えて、これが十分小さい、すなわち直角とほぼ見做せる程度だとおもわ
2119 れたら直交回転を行う、とするのが良いでしょう。

2120 ■**行列による説明** 因子分析の行列モデルは次のようなものでした。

$$R = AA' + D^2$$

2121 ここで AA' に単位行列 I を挟んでも、 $AA' = AIA'$ で結果は変わりませんね。ここで $I = TT'$ とな
2122 るような行列 T があっても同じです。たとえば 2 因子構造の時に、次のような行列 T を考えたとする
2123 と、 $TT' = I$ であることがわかります。

$$T = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$$

2124

$$TT' = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix} \begin{pmatrix} \cos\theta & \sin\theta \\ -\sin\theta & \cos\theta \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I$$

2125 このような T は回転行列と言われます。原点を中心に、座標を角度 θ で回転させるからです。 θ の値は任意
2126 なので、因子負荷行列は、 θ に応じて無数にあることになります。であれば、使い勝手の良い角度で回転して
2127 やればいいとも言えるわけです。

2128 **直交回転** の場合は角度が直角、すなわち相関がないわけですから、数学的には回転行列が $TT' = I$ の
2129 ようになった状態を言います。一方で**斜交回転** の場合は相関があるわけですが、

$$I = T'\Phi T$$

2130 となるように置いておけば、 $R = AT'\Phi TA' + D^2 = BB' + D^2$ の形が保たれているので因子分析
 2131 モデルとしては違いがないわけで、このような Φ を任意に定めることができます。この時 Φ は因子間相関
 2132 (factor correlations) と呼ばれます。

2133 直交回転も斜交回転も、目的は解釈をしやすくすることであり、そのための基準はいくつか考えら
 2134 れます。それらの基準に応じて、たとえば直交回転ではバリマックス回転 (varimax rotation)、斜交回転
 2135 ではプロマックス回転 (promax rotation) などがあり、ここではジオミン回転 (geomin rotation) を
 2136 選びました*6。

2137 10.4.4 因子分析の結果をみる

2138 出力 10.2 に因子負荷行列を表示しました。あの A が計算されていますよ！

R の出力 10.2: 因子負荷行列

```
Factor Analysis using method = ml
Call: fa(r = dat, nfactors = 5, rotate = "geominQ", fm = "ml")
Standardized loadings (pattern matrix) based upon correlation matrix
  item ML2 ML1 ML5 ML3 ML4 h2 u2 com
N1  16  0.83 -0.07 -0.24  0.02 -0.04 0.71 0.29 1.2
N2  17  0.80 -0.02 -0.23  0.03  0.02 0.66 0.34 1.2
N3  18  0.68  0.15 -0.01 -0.03  0.03 0.53 0.47 1.1
N5  20  0.48  0.26  0.15  0.01 -0.13 0.34 0.66 1.9
N4  19  0.46  0.43  0.03 -0.12  0.09 0.48 0.52 2.2
E2  12  0.10  0.69 -0.08  0.00 -0.04 0.55 0.45 1.1
E1  11 -0.07  0.58 -0.07  0.12 -0.09 0.37 0.63 1.2
E4  14  0.00 -0.55  0.33  0.00 -0.07 0.52 0.48 1.7
E5  15  0.13 -0.43  0.03  0.25  0.22 0.40 0.60 2.4
E3  13  0.05 -0.38  0.27 -0.02  0.30 0.44 0.56 2.8
A3   3  0.02 -0.11  0.65  0.05  0.03 0.51 0.49 1.1
A2   2  0.03 -0.02  0.59  0.11  0.03 0.40 0.60 1.1
A5   5 -0.09 -0.21  0.58  0.01  0.04 0.48 0.52 1.3
A4   4 -0.03 -0.05  0.45  0.21 -0.15 0.29 0.71 1.7
A1   1  0.16 -0.12 -0.39  0.02 -0.03 0.15 0.85 1.6
C2   7  0.09  0.10  0.07  0.63  0.09 0.43 0.57 1.2
C4   9  0.18  0.06  0.05 -0.62 -0.05 0.47 0.53 1.2
C3   8  0.01  0.06  0.09  0.56 -0.04 0.32 0.68 1.1
C5  10  0.19  0.18  0.00 -0.55  0.08 0.43 0.57 1.5
C1   6  0.01  0.03 -0.02  0.52  0.18 0.32 0.68 1.2
O3  23 -0.01 -0.15  0.07 -0.01  0.62 0.47 0.53 1.1
O1  21 -0.03 -0.10  0.01  0.04  0.53 0.32 0.68 1.1
O5  25  0.14 -0.04  0.08 -0.03 -0.52 0.27 0.73 1.2
O2  22  0.20  0.00  0.18 -0.07 -0.44 0.24 0.76 1.8
O4  24  0.11  0.33  0.13 -0.02  0.39 0.26 0.74 2.4
```

2139
 2140 出力 10.2 にあるように、各行に N,E,A,C,O の項目が並んでいます。この順番が入れ替わっているの
 2141 は、表示オプションで sort=T を入れたからで、因子負荷量の大きい順に並べ替えて表示させたからです。
 2142 現に左上から右下にかけて、因子負荷量の大きさの順に並んでいますね。列として ML1,2,3,4,5 とあるの

*6 ジオミン回転には直交モデルと斜交モデルがあり、それぞれ geominT, geominQ という名前で実装されています。

2143 は最尤法 (ML) で抽出した因子の 1,2,3,4,5 を表しています。h2 は h^2 , u2 は $1 - h^2 = d^2$ をあらわ
 2144 しています (独自性 (uniqueness) の u)。com は複雑性指標 (index of complexity) と呼ばれ、
 2145 $\sum (a_j^2)^2 / \sum a_j^4$ で計算される指標です。共通性は高く、独自性、複雑性の低い項目がいい項目だといえる
 2146 でしょう。ちなみに、出力オプションとして cut=0.3 を追加すると、因子負荷量が 0.3 に満たないところの表
 2147 示はなくなります。どの項目がどの因子と関係あるのかがわかりやすくなるのでおすすめです。

2148 出力 10.3 にはその他の情報として、各因子の因子負荷行列の二乗和 (SSloadings)、分散説
 2149 明率 (Proportion Var)、累積分散説明率 (Cumulative Var)、全体を 1 とした時の説明比率
 2150 (Proportion Explained)、その累積 (Cumulative Proportion) が表示されています。回転するま
 2151 えは固有値と因子負荷量の二乗和が合致するのですが、回転した後の場合はそうはならないので、これら
 2152 をどれくらい説明できているかを考えます。とくに類遺跡分散説明率が、今回は 5 因子で 41% しかありませ
 2153 ん。59% をゴミとして捨てたようなものですから、これでよかったのかな、と思案するところです。

R の出力 10.3: 因子分析結果 2

	ML2	ML1	ML5	ML3	ML4
SS loadings		2.50	2.24	2.05	1.95 1.61
Proportion Var		0.10	0.09	0.08	0.08 0.06
Cumulative Var		0.10	0.19	0.27	0.35 0.41
Proportion Explained		0.24	0.22	0.20	0.19 0.16
Cumulative Proportion		0.24	0.46	0.66	0.84 1.00

2154

2154 その下出力 10.4 には、因子間相関が表示されています。これをみると、第 2 因子 (ML1) と第 3 因子
 2155 (ML5) との間に -0.34 という中程度の負の相関がみられます。直交回転はこれを 0 として解釈してしま
 2156 いますから、スッキリはするのですがデータには適合してないところかもしれません。このように、分析する時は
 2157 まず斜交回転を行い、因子間相関がなべて低い場合は直交で分析をやり直す、というステップを経ます。
 2158

R の出力 10.4: 因子分析結果 3

With factor correlations of

	ML2	ML1	ML5	ML3	ML4
ML2	1.00	0.11	0.06	-0.12	0.06
ML1	0.11	1.00	-0.34	-0.22	-0.16
ML5	0.06	-0.34	1.00	0.18	0.23
ML3	-0.12	-0.22	0.18	1.00	0.17
ML4	0.06	-0.16	0.23	0.17	1.00

2159

10.5 因子分析の後で

2160 因子数を定め、因子負荷量が明らかになると、ここから逆算的に因子得点を計算できます。fa 関数には
 2161 scores オプションがあり、これをつけて実行すると因子得点行列を表示させることができます。因子得点は
 2162 因子数 × 人数分 (ここでは 2800 人分) ありますので、大変大きなサイズのデータですから注意してくださ
 2163 い。ともあれ、そこに含まれているのは、各因子についての個々人の相対的な関与度ということになります。
 2164

code : 10.4 因子得点を算出する

2165

2166

2167

2168

```
1 result <- fa(dat, nfactors = 5, fm = "ml", rotate = "geominQ", scores = T)
2 result$scores
```

2169 モデルの仮定にあった通り、因子負荷量は標準化されたスコアです。すなわち平均情報が取り払われている
 2170 ことに注意してください。たとえば 7 件法で「4. どちらでもない」を挟んで、ポジティブ・ネガティブな意味を
 2171 持っていたとしても、ここでは個々人間の相対的な大小関係しか意味しないことになります。

2172 そういう意味で、絶対的なスケールの情報を残したまま得点を考えてみたいということもあるでしょう。そう
 2173 いう時は**簡便的因子得点**といって、因子に関する項目の素点から平均点を算出し、個々人の得点とする方
 2174 法もあります。最もこのやり方は、平均点の情報は保持されますが、因子分析で除去できた d_j^2 、すなわち独
 2175 自性成分を再び取り込むことになっている点に注意が必要です。次のコード code:10.4 は、簡便法による
 2176 因子得点と因子分析の結果から計算された因子得点との相関係数を算出する例です。高い相関を示すところ
 2177 (ML2 と N で 0.9475742, ML5 と A で 0.83845417) もありますが、低いところもありますので、どちら
 2178 の方法で計算するのか、その場合の長所短所をよく理解して使うようにしましょう。

code : 10.5 二種類の因子得点計算法比較

```

2179 1 dat %>%
2180 2   dplyr::mutate(
2181 3     A = (A1 + A2 + A3 + A4 + A5)/5,
2182 4     C = (C1 + C2 + C3 + C4 + C5)/5,
2183 5     E = (E1 + E2 + E3 + E4 + E5)/5,
2184 6     N = (N1 + N2 + N3 + N4 + N5)/5,
2185 7     O = (O1 + O2 + O3 + O4 + O5)/5
2186 8   ) %>%
2187 9   dplyr::bind_cols(., result$scores %>% as.data.frame) %>%
2188 10  dplyr::select(ML1, ML2, ML3, ML4, ML5, A, C, E, N, O) %>%
2189 11  cor(use="pairwise")
2190
2191

```

2192 10.6 さいごに

2193 今回は探索的因子分析の実行例をみてみました。いかがだったでしょうか。計算自体にはほとんど時間が
 2194 かかりませんので、色々な因子数や回転方法、抽出法を試して見るのも良いかもしれません。

2195 このように比較的簡単に計算できるようになってきましたが、大事なはその目的と意味です。関数の名
 2196 前を覚えて、機械的に XX 法 XXX 回転で答えが出たから正しいのだ、などと思わないことです。因子数の
 2197 決め方の時にも薄々感じていただけたかと思いますが、数学的なモデルとは裏腹に、実際にやる時は職人的
 2198 ノウハウも結構必要です。スクリーンプロットを「みて」とか、回転法を「選んで」とか、因子得点の計算の仕方を
 2199 「比較して」、といったところは人間的要素であり、やり方が変われば結果が変わる可能性があります。論文な
 2200 どに示されているこれらの方法、指標も、変えようと思えば変えられるところだということです。科学的真実の
 2201 探究が目的であれば、良い見栄えや自分に都合の良い結果のための手法選択は無意味であることは理解し
 2202 てもらえるところだと思います。皆さんも決して目的を見失わず、批判的に検証できるようになってください。

2203 探索的因子分析は因子を見つけ出すことができます。が、それは決して人間の心の何かを取り出したわけ
 2204 ではないのです。人間がそもそも持っている心の状態を取り出したのではなく^{*7}、人間に「項目群」「紙とペ
 2205 ン」という刺激を与え、反応を強いたのです。その結果、「項目群にみられる相関行列に潜む次元」が見いだ
 2206 されただけであり、それが人間の持つ特性だと考えるのは、あくまでの研究者の主観的主張にすぎません。実
 2207 際にあり得る話ですが、「学校生活は嫌ですか」「学校は楽しくないと思うことがありますか」「学校の友達とは
 2208 遊びたくないと思うことがありますか」といったネガティブな項目を山のように投げかけ、因子分析することで

^{*7} 因子分析をすると「因子を抽出した」と論文などに記載される場合がありますが、このようにあたかも先にあるエッセンスを抜き出したのだ、という表現は適切ではないと批判されることもあります。

2209 「学校を嫌う項目群=因子」を抽出し、子供は学校が嫌いなのだという結論を出すのがおかしいことなのはわ
2210 かりますよね。項目群のなかにネガティブなものしかなければ、因子もそれに応じたものしか出ません。そして
2211 因子得点が相対的なものでしかないので、素点を見ずに「学校嫌い得点が高い、低い」として議論を進めて
2212 も仕方がないのです。そうした場合は簡便的因子得点を計算して、平均点がちゃんと中点「4. どちらでもな
2213 い」を挟んでいるかどうかで嫌っているかどうかを判断しなければなりません。このように、自分で項目に埋め
2214 込んだ呪いを自分で取り出して見つけ出したと騒ぐ、**てっちゃんの手品**にならないようにすることは、常に心
2215 がけなければなりません*8。

2216 こうした批判的観点を持つためには、数理的にはどういうことをやっていて、どこからが人為的な問題なの
2217 かをしっかり把握しておくことが役に立つでしょう。

2218 10.7 課題

2219 今日の授業でおこなったすべての次の計算をする R コードを提出してください。ファイル形式は R スクリ
2220 プトか Rmd とします。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱い
2221 になります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を
2222 受けてください。

*8 てっちゃんの手品については、[小杉 \(2018\)](#) を参照してください。

第 11 章

R を使った項目反応理論

今回も R を使った演習です。R を使った行列計算、R を使った因子分析と進めてきましたが、今回は R を使った項目反応理論を実践的に理解していきましょう。

11.1 項目反応理論の実際

11.1.1 データの素描

項目反応理論は、テスト理論に関するモデルですから、サンプルデータとして二値データを使いましょう。表 11.1 に使うデータの一部を示しています^{*1}。

表 11.1 IRT で使うサンプルデータ (一部)

V1	V2	V3	V4	V5	V6	V7	V8	V9	V10
0	0	1	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	1
0	1	0	1	0	0	0	0	0	1
0	1	1	1	1	1	0	0	0	1
0	1	1	1	1	0	0	0	0	1
1	1	1	1	0	0	0	0	1	1

このように、テスト理論で使うデータは 0/1 の二値データで、1 が正答、0 が誤答を表しています。これらのデータを使って R で分析をします。IRT を実行するパッケージとして、`ltm`(Rizopoulos, 2006) や `irtos`(Partchev, Partchev and Suggests, 2017) などがあります。ここでは `ltm` パッケージの方を使っていきましょう。

分析に先立って、データの記述統計量を確認しておきます。code: 11.1 を実行してください。

code : 11.1 データの読み込みから記述統計量の計算まで

```
1 rm(list=ls())
2 library(tidyverse)
3 library(ltm)
4 dat <- read_csv("IRTsample.csv")
5 ltm::descript(dat)
```

*1 このデータ全体については、`IRTsample.csv` というファイル形式で配布していますので、それを読み込んで使ってください

2243 ■コード解説

2244 1 行目 環境の初期化

2245 2-3 行目 パッケージの読み込み

2246 4-5 行目 サンプルデータ IRTsample.csv を読み込む。文字化けする人は、read_csv 関数のオプション
2247 として locale=locale(encoding="UTF-8") を追加してください。

2248 6 行目 ltm パッケージに含まれる descript 関数を実行

2249 いろいろな情報が一気に出力されるので、順にみていきます。まず R の出力 11.1 のところです。ここには
2250 各変数の 0 と 1 の比率が入っています。この例では、V1 は 0 が 71.0%、1 が 29.0% 含まれているというこ
2251 とです。最後の logit は、0 から 1 の間に入る数字 p にたいして、 $\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$ で計算される
2252 量です。ここでは p は正答率であり、この logit の値が負であれば 0 の比率が、正であれば 1 の比率が多
2253 かった、ということの意味します。0/1 がちょうど半々であれば、logit の値は 0 になる、そういう数字です。
2254 これを見るだけでも、V8 はずいぶん正答率が低かったんだな、V3 は逆に正答率が高かったんだな、とい
2255 うことがわかります。

R の出力 11.1: 記述統計の出力 1

```

Proportions for each level of response:
  0    1  logit
V1 0.710 0.290 -0.8954
V2 0.114 0.886  2.0505
V3 0.186 0.814  1.4762
V4 0.206 0.794  1.3492
V5 0.670 0.330 -0.7082
V6 0.832 0.168 -1.5999
V7 0.666 0.334 -0.6901
V8 0.884 0.116 -2.0309
V9 0.652 0.348 -0.6278
V10 0.254 0.746  1.0774

```

2256

2257 次のセクション R の出力 11.2 には、正答数の分布が示されています。10 問全部正解したのは 3 人だけ、
2258 6 問正解した人が 104 人だった、といったことがわかります。

R の出力 11.2: 記述統計の出力 2

```

Frequencies of total scores:
  0  1  2  3  4  5  6  7  8  9 10
Freq 4 19 41 62 97 104 66 54 36 14  3

```

2259

2260 次のセクション R の出力 11.3 には、点双列相関係数 (Point Biserial correlation) が計算されてい
2261 ます。この相関係数は 2 値と連続値の相関係数であり、ここでは各変数 (二値) と合計得点 (0 から 10) の相
2262 関係数を計算しています。この合計得点の計算方法には、その変数自身を含める場合 (Included) と含めな
2263 い場合 (Excluded) の二種類があります。具体的に説明したほうがわかりやすいでしょう。たとえば変数 V1
2264 についていうと、「その変数自身を含める場合」の合計得点とは $V1 + V2 + V3 + V4 + V5 + V6 + V7 +$
2265 $V8 + V9 + V10$ で計算しています。「含めない場合」は $V2 + V3 + \dots + V10$ とするわけです。含める場

2266 合の相関係数が 0.5207, 含めない場合の相関係数が 0.2480 です。当然自身が含まれているほうが合計得
 2267 点に深く関わっているので相関係数が高くなります。裏を返せば, 含める場合と含めない場合の相関係数に
 2268 大きな差があるとき, その変数はテスト全体に大きく貢献する重要な変数だといえるでしょう。

R の出力 11.3: 記述統計の出力 3

Point Biserial correlation with Total Score:

	Included	Excluded
V1	0.4228	0.2118
V2	0.3959	0.2502
V3	0.5202	0.3567
V4	0.5209	0.3501
V5	0.5449	0.3470
V6	0.4635	0.2983
V7	0.5440	0.3452
V8	0.2901	0.1350
V9	0.4992	0.2890
V10	0.5890	0.4180

2269

2270 次のセクション R の出力 11.4 には, α 係数 (alpha coefficient) が算出されています。これも全体の
 2271 場合と, 特定の変数を抜いた場合とを計算しています。心理尺度における信頼性係数としての α 係数は,
 2272 0.8 以上であれば信頼性がある, といった目安がありますが, 二値データの場合はここに示したようにそれよ
 2273 りもずいぶんと低くなります。これは二値データの分散が小さくなることから仕方のないことですし, 項目反応
 2274 理論で分析する上での信頼性は情報量関数として表されますから, 問題ではありません。

R の出力 11.4: 記述統計の出力 4

Cronbach's alpha:
value

All Items	0.6341
Excluding V1	0.6303
Excluding V2	0.6195
Excluding V3	0.5978
Excluding V4	0.5986
Excluding V5	0.5982
Excluding V6	0.6100
Excluding V7	0.5987
Excluding V8	0.6380
Excluding V9	0.6130
Excluding V10	0.5817

2275

2276 最後にセクション R の出力 11.5 では, 変数同士の関連がとくに弱いものが示されています。変数同士の
 2277 ペアについて度数分布表を作成し, χ^2 検定によって p.value を算出しています。この値が小さければ, 2
 2278 つの変数の関係が強いことを示していると解釈できるのですが, ここにあげているように大きな数字になっ
 2279 ているということは, 変数間の関係が弱いと読み取ることができます。

R の出力 11.5: 記述統計の出力 4

```

Pairwise Associations:
  Item i Item j p.value
1       2     8  0.625
2       1     8  0.409
3       1     2  0.347
4       6     8  0.303
5       4     8  0.234
6       7     8  0.222
7       5     8  0.195
8       1     4  0.191
9       1     3  0.166
10      3     8  0.124

```

2280

11.1.2 IRT モデルの実際

2281 データの感じがわかったところで、IRT モデルを実行していきましょう。

2282 今回は 3 つのモデルを試してみます。確認のために、少しモデル式に言及しておきます。困難度母数だけに
 2283 注目するのが 1 パラメータ・ロジスティックモデル (1 parameter logistic model, 1PL) で、モデル
 2284 式は次の通りです。
 2285

$$p_j(\theta) = \frac{1}{1 + \exp(-1.7(\theta - b_j))}$$

2286 ここで $p_j(\theta)$ は能力が θ のひとが項目 j に正答する確率であり、 b_j が**困難度母数** difficulty parameter
 2287 と呼ばれます。これは開発者の名前を使って、別名**ラッシュモデル** (rasch model) とも呼ばれています。

2288 項目の識別力にも注目するのが 2 パラメータ・ロジスティックモデル (2 parameter logistic model,
 2289 2PL) で、モデル式は次の通りです。

$$p_j(\theta) = \frac{1}{1 + \exp(-1.7a_j(\theta - b_j))}$$

2290 ここで a_j は**識別力母数** (discriminant paramter) と呼ばれます。

2291 最後に 3 パラメータ・ロジスティックモデル (3 parameter logistic model, 3PL) では、困難度、
 2292 識別力に加えて**当て推量母数** (guessing parameter) というのを推定します。モデル式は次の通りです。

$$p_j(\theta) = c + \frac{1 - c}{1 + \exp(-1.7a_j(\theta - b_j))}$$

2293 もちろん式ではピンと来ないという人もいると思いますので、実際にどのように表されるかみていきま
 2294 しょう。

1PL モデルの場合

2295 ltm パッケージを使って 1PL モデルを実行します。code: 11.2 を実行してください。

code : 11.2 1PL モデルの実行

```

2297 1 result.1pl <- rasch(dat)
2298

```

```

2299 2 print(result.1pl)
2300 3 plot(result.1pl,type="ICC")
2301 4 plot(result.1pl,type="IIC")
2302 5 plot(result.1pl,type="IIC",items = 0)
2303

```

2304 ■コード解説

2305 1 行目 1PL モデル (ラッシュモデル) の実行
 2306 2 行目 結果の出力
 2307 3 行目 項目特性曲線 (Item Characteristic Curve) のプロット。plot 関数のオプションで表示する
 2308 情報を ICC と指定。
 2309 4 行目 項目情報曲線 (Item Information Curve) のプロット。plot 関数のオプションで表示する情
 2310 報を IIC と指定。
 2311 5 行目 テスト情報曲線 (Test Information Curve) のプロット。plot 関数のオプションで表示する
 2312 項目番号をゼロにすると TIC になる。
 2313 2 行目で結果の出力を行っていますが、出力 11.6 のように表されます。

R の出力 11.6: 2PL モデルの結果出力

```

Call:
rasch(data = dat)

Coefficients:
  Dffc1t.V1  Dffc1t.V2  Dffc1t.V3  Dffc1t.V4  Dffc1t.V5  Dffc1t.V6
    0.992    -2.208    -1.615    -1.480     0.787     1.744
  Dffc1t.V7  Dffc1t.V8  Dffc1t.V9  Dffc1t.V10  Dscrmn
    0.768     2.187     0.699    -1.189     1.126
Log.Lik: -2468.743

```

2314
 2315 1PL モデルで表現する項目の違いは**困難度母数**だけですから、この数字 (Dffc1t) が小さければ簡
 2316 単な項目、大きければ難しい項目だったことがわかります。たとえば V2 の困難度母数は $b_2 = -2.208$ です
 2317 が、記述統計のところで見たと同じ正答率 (88.6%) を考えると、簡単な問題だったことがわかりますね。**識別力母**
 2318 **数 (は) Dscrmn**として 1.126 と推定されています。これは項目を通じて一定です。このことは図 11.1 の左に
 2319 ある ICC をみて、すべての曲線の傾きが同じことで確認できます。

2320 図 11.1 の左が数値的出力結果を図示したものです。これを見た方がわかりやすいかもしれません。x 軸
 2321 は θ 、つまり測定する潜在変数の値になっており、テストで言えばある人の学力 θ_i があつたときに、各問いに
 2322 どれぐらいの確率で正答するかを表していることになります。これをみると V2 がもっとも簡単で、V8 が最も
 2323 難しかったテストであることがわかります。

2324 また図 11.1 の中央にあるのが IIC で、それぞれの項目がどのあたりの能力値 θ を測定するときにもっと
 2325 も敏感になるかを表しています。これは項目反応理論という**信頼性**の表現であることを思い出してください。
 2326 最後に図の左側にあるのは、この IIC を 10 問分足し合わせたテスト全体の情報曲線、TIC になります。こ
 2327 れをみると、 $\theta = 0$ よりやや右側のあたりにピークがきているので、このテスト全体としては平均より高い学
 2328 力の人を測定するとき、最も鋭敏に働くことがわかるでしょう。

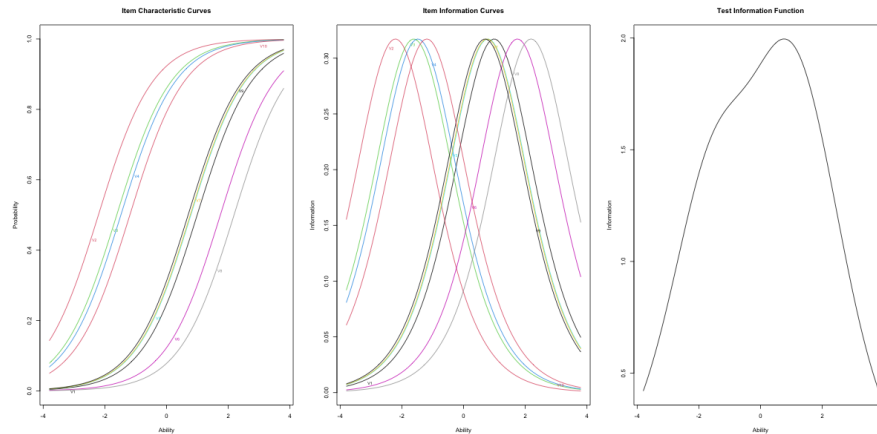


図 11.1 1PL のさまざまなプロット。左から ICC,IIC,TIC

2329 2PL モデルの場合

2330 ひきつづき ltm パッケージを使って、2PL モデルを実行していきましょう。code:11.3 を実行してくだ
2331 さい。

code : 11.3 2PL モデルの実行

```

2332 1 result.2pl <- ltm(dat~z1)
2333 2 print(result.2pl)
2334 3 plot(result.2pl,type="ICC")
2335 4 plot(result.2pl,type="IIC")
2336 5 plot(result.2pl,type="IIC",items = 0)
2337
2338

```

2339 ■コード解説

2340 1 行目 1PL モデル (ラッシュモデル) の実行

2341 2 行目 結果の出力

2342 3 行目 項目特性曲線 (Item Characteristic Curve) のプロット

2343 4 行目 項目情報曲線 (Item Information Curve) のプロット

2344 5 行目 テスト情報曲線 (Test Information Curve) のプロット

2345 2 行目で結果の出力を行っていますが、出力 11.7 のように表されます。

R の出力 11.7: 2PL モデルの結果出力

```
Call:
ltm(formula = dat ~ z1)

Coefficients:
      Dffc1t  Dscrmn
V1      1.602   0.603
V2     -2.078   1.228
V3     -1.213   1.892
V4     -1.199   1.620
V5      0.759   1.176
V6      1.698   1.175
V7      0.740   1.174
V8      4.133   0.516
V9      0.866   0.828
V10     -0.880   2.020

Log.Lik: -2446.104
```

2346

2347 2PL モデルで表現する項目の違いは、**困難度母数**と**識別力母数**の 2 つです。識別力母数を入れること
 2348 で、IIC 曲線の傾きもデータに合わせて調整できるようになりました。識別力母数のデフォルトは (推定しな
 2349 いのであれば) 1.0 ですから、ここから大きく逸脱するような項目には注意です。今回は V8 をみると、困難度
 2350 が非常に高いこともわかりますが、識別力が非常に低くなっているため、この項目に誤答したからと言ってす
 2351 ぐさま成績が悪い、と判断できるほどではないことがわかります。

2352 最後の一行に Log.Lik とありますが、これは**対数尤度 (log likelihood)** を表しています。これが大き
 2353 いほど、データとモデルの適合度が高いことを意味します。1PL のときの Log.Lik が -2468.743 でしたか
 2354 ら、1PL モデルより 2PL モデルの方が当てはまりがよかったと言えるでしょう。

2355 プロットによる出力も見てください (図 11.2)。識別力母数が入ったことで、ICC の傾きが項目によっ
 てさまざまになりますし、IIC や TIC の見た目もすっかり変わってしまいましたが、読み取り方は同じです。

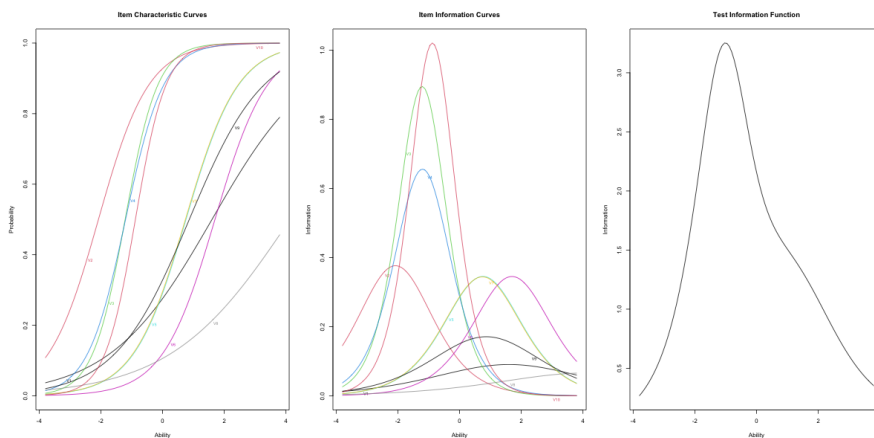


図 11.2 2PL のさまざまなプロット。左から ICC, IIC, TIC

2356

2357 3PL モデルの場合

2358 最後に ltm パッケージを使って、3PL モデルを実行していきましょう。code:11.4 を実行してください。

code : 11.4 3PL モデルの実行

```
2359
2360 1 result.3pl <- tpm(dat)
2361 2 print(result.3pl)
2362 3 plot(result.3pl,type="ICC")
2363 4 plot(result.3pl,type="IIC")
2364 5 plot(result.3pl,type="IIC",items = 0)
2365
```

2366 コードは一行目の関数以外ほとんど同じですので、逐一解説は致しませんが、出力のところを見ておきま
2367 しょう (出力 11.8)。

R の出力 11.8: 3PL モデルの結果出力

```
Call:
tpm(data = dat)

Coefficients:
      Gussng Dffclt Dscrmn
V1      0.230   1.407  15.123
V2      0.638  -0.610   2.232
V3      0.225  -0.825   2.907
V4      0.324  -0.637   2.324
V5      0.149   0.931   2.992
V6      0.043   1.613   1.706
V7      0.074   0.879   1.509
V8      0.000   4.490   0.471
V9      0.133   1.144   1.255
V10     0.278  -0.438   3.085

Log.Lik: -2431.371
```

2368
2369 3PL モデルは困難度、識別力に加えて**当て推量母数**が計算されます。これは ICC の形 (図 11.3 の左) を
2370 見ればわかりますが、曲線全体が底上げされたようになっているものがあります。この最低ラインの高さが当
2371 て推量母数の大きさです。図や数値から、V2 のそれが大きい数字であることが読み取れます。ICC は θ の
2372 関数である正答率ですが、これが θ の値にかかわらず一定の値を取るということは、どれだけ能力が低く
2373 もその確率で正答してしまうことを意味しています。なので「適当に推測して答えを当てられる大きさ」という
2374 意味で当て推量母数と呼ばれているのです。

2375 今回当て推量母数を含めて考えたことで、IIC や TIC がまた大きく変化しました。この数字が高すぎる
2376 V2 は、テストとしてはあまり良い項目ではなかったのかもしれませんが、また、モデルの適合度としては 3PL モ
2377 デルが最もよかったのですが、2PL モデルとそれほど大きく違うわけでもありません。どのモデルでどのよう
2378 に推定すべきか、推定されたモデルはどのような仮定を持っているのかについて、しっかり理解した上で利用
2379 するようにしましょう。

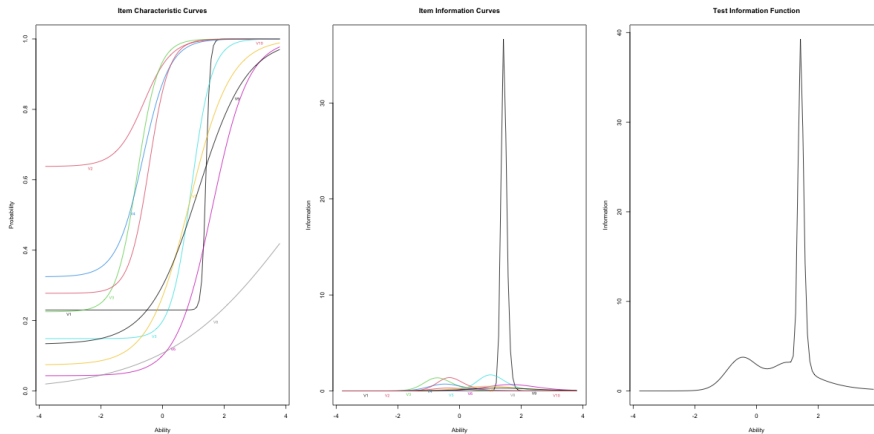


図 11.3 3PL のさまざまなプロット。左から ICC,IIC,TIC

11.2 段階反応モデルの実際

2380

2381 続いて**段階反応モデル**の練習に進みましょう。項目反応理論は正答/誤答の二値データに対して用いるも
 2382 のでしたが、これを多段階反応に拡張したモデルの1つが**段階反応モデル**です。多段階反応とは、「まったく
 2383 当てはまらない」を1、「非常に当てはまる」を7とする、といったような**リッカート尺度**のようなスタイルの調
 2384 査項目です。反応させるカテゴリの段階数に応じて、5件法とか7件法と呼ばれたりします。

2385 リッカート尺度の場合は、シグマ法をつかって累積度数からカテゴリに該当する値を算出し、それを元に態
 2386 度得点を作るという方法でした。とはいえこのやり方でやっても、いい尺度であればそのまま素点を尺度値と
 2387 してもほとんど問題がないため、シグマ法が回りくどい方法だと思われたのか、使われなくなってしまいま
 2388 した。しかし実際には「いい尺度であれば」という条件を満たしているかどうか、しっかり確認しなければなりま
 2389 せん。歪んだ分布や偏ったカテゴリ反応なのに、素点をそのまま尺度値とするのは不適切なことがあります。
 2390 また、リッカート尺度で撮られたデータは因子分析をすることによって、解釈度に分類されて、因子ごとのスコ
 2391 アを算出して利用されることが一般的です。しかし、尺度値が適切でなければ、ピアソンの積率相関係数も適
 2392 切ではなく、そこから計算が始まる因子分析のモデルにも疑義が生じます。

2393 カテゴリ反応はあくまでも並びの問題であって、**順序尺度水準**の情報しか持たないと考えるのであれば、
 2394 その背後にある態度のような連続体を仮定し、この態度 θ が大きくなることで徐々に反応カテゴリが変わる、
 2395 という**段階反応モデル**のほうが適切な分析方法になるわけです。

2396 項目反応理論(の段階反応モデル)を使って確認することは決して難しくありません。今回は因子分析の
 2397 時に使った bfi データの一部を使って段階反応モデルを実践してみましょう。一部を使うとしたのは、段階
 2398 反応モデルは項目反応理論の中間ですから、潜在変数が1次元(一因子)である、という仮定があるからで
 2399 す。ここではN因子(Neuroticism, 神経質さ)を使った例を示しています。

code : 11.5 GRM の実行

2400

```

2401 1 library(psych)
2402 2 dat <- bfi %>%
2403 3   dplyr::select(-gender, -education, -age) %>%
2404 4   dplyr::select(starts_with("N"))
2405 5 result.grm <- grm(dat)
2406 6 result.grm

```

```

2407 7 plot(result.grm, type = "ICC", items = 1)
2408 8 plot(result.grm, type = "IIC", items = 1)
2409 9 plot(result.grm, type = "ICC", items = 4)
2410 10 plot(result.grm, type = "IIC", items = 4)
2411 11 plot(result.grm, type = "IIC", items = 0)
2412

```

2413 ■コード解説

2414 1 行目 psych パッケージを読み込む。サンプルデータ bfi を用いるためです。

2415 2-4 行目 データの加工。bfi データから、性別、最終学歴、年齢などの情報を取り除き、また N で始まる変数だけ選出 (select) しています。

2416

2417 5 行目 段階反応モデルの実行。

2418 6 行目 結果の出力。

2419 7 行目 項目 1 についての ICC、厳密には項目反応カテゴリ特性曲線 (Item Response Category Characteristic Curve) と呼びます。

2420

2421 8 行目 項目 1 についての IIC を表示させています。

2422 9 行目 項目 4 についての ICC を表示させています。

2423 10 行目 項目 4 についての ICC を表示させています。

2424 11 行目 項目群についての TIC を表示させています。

2425 code:11.5 の 6 行目で、分析結果の数値的出力をさせていますので、それをみておきましょう (出力 11.9)。

2427 ここで、Dscrnm とあるのは識別力母数ですが、困難度母数に対応するのが Extrmt1 から Extrmt5 にあたります。この尺度は 6 件法ですので、反応カテゴリが変わる切れ目、すなわち閾値が 5 つあるのです。これらが 2PL モデルでいうところの困難度に対応すると考えてもいいでしょう。

R の出力 11.9: GRM モデルの結果出力

```

Call:
grm(data = dat)

Coefficients:
      Extrmt1 Extrmt2 Extrmt3 Extrmt4 Extrmt5 Dscrnm
N1    -0.810  -0.093   0.342   0.985   1.720   3.125
N2    -1.366  -0.555  -0.111   0.648   1.483   2.890
N3    -1.187  -0.298   0.123   0.876   1.767   2.026
N4    -1.564  -0.355   0.238   1.240   2.280   1.277
N5    -1.296  -0.125   0.494   1.479   2.530   1.112

Log.Lik: -21721.91

```

2430

2431 数値を見てもピンとこないかもしれませんので、図でこの各項目の特徴を確認しておきましょう。変数

2432 N1 と N4 の IRCCC を図 11.4 に示しました。この図に表されている複数の曲線が、 θ に対応した各カテゴリ

2433 に反応する確率を表しています。図の左、N1 は綺麗なカーブが描かれており、各カテゴリのピークが見て取

2434 れますが、N4 は反応カテゴリ 3 の山が埋もれてしまっていることがみてとれます。この項目については、6 段

2435 階ではなく 5 段階で反応を求めるのがよかったのかもしれない*2。

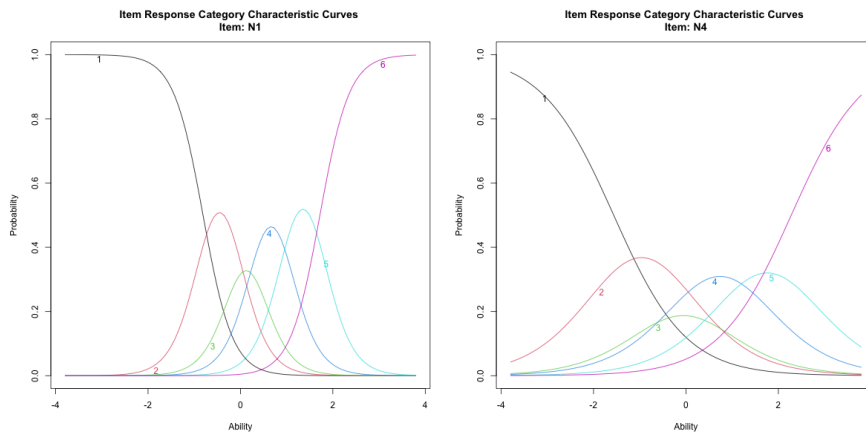


図 11.4 GRM の IRCCC

2436 原理的には難しいことのように感じますが、コードを書いて実践する分にはほとんど苦勞がなく、二値モデ
 2437 ルと同じように考えることができます。コードには IIC や TIC を描くものもありますが、1PL,2PL,3PL モデル
 2438 と変わらない感じで多段階モデルが実行できますね。みなさんも心理尺度を使う時は、盲目的にスコアリン
 2439 グし、因子分析するのではなく、どういう仮定の元でどういう計算をしているかを理解しながら、データや調査
 2440 協力者にとって最適な分析をするように心がけてください。

2441 11.3 カテゴリカル因子分析との対応

2442 さて多段階反応による分析も簡単にできることがわかりましたが、気になるのは段階反応モデルが一因子
 2443 構造しか仮定していないところです。bfi データはビッグファイブと呼ばれるように五因子 (5 次元) の仮定が
 2444 あるのですが、一因子ずつ grm を実行し、あとで統合するのは大変です。

2445 でも大丈夫。簡単に多次元に展開する方法があります。いわゆる一般的な因子分析は、ピアソンの積率相
 2446 関係数を元に (積率相関係数の相関行列 R の固有値分解から) 因子を求めていくのですが、この積率相関
 2447 係数が間隔尺度水準を仮定したもので、シグマ法で計算された尺度値を使っていないのであれば疑義が残
 2448 る、という話でした。この元になる相関係数が、潜在的連続体を仮定した順序尺度水準向けの相関係数であ
 2449 る、ポリコリック相関係数であれば、その問題は解決します。反応段階が順序尺度を仮定した因子分析のこ
 2450 とをとくにカテゴリカル因子分析と言います。カテゴリカル因子部な席は、前回紹介した psych パッケージに
 2451 少しオプションを足すだけで簡単に実行できるものです。

code : 11.6 GRM の実行

```
2452
2453 1 dat <- bfi %>% dplyr::select(-gender, -education, -age)
2454 2 result.fa <- fa.poly(dat, nfactors = 5, rotate = "geominQ")
2455 3 print(result.fa, sort = T, cut = 0.3)
2456
```

2457 ■コード解説

2458 1 行目 サンプルデータ bfi の加工をしています。

*2 ちなみに N1 は「すぐ怒る (Get angry easily)」, N4 は「割と凹む (Often feel blue)」です。日本語訳は小杉の意識で定訳ではありません。あしからず。

2459 2 行目 ポリ個リック相関係数に基づく因子分析の実行

2460 3 行目 結果の出力。

2461 コードの 2 行目を見ていただくとわかる通り、ただの `fa` 関数ではなく `fa.poly` としただけで、あとは因子
2462 分析と同じです。この `.poly` のところがポリコリック相関係数を使うよう指定しているところで、これで実質的
2463 に測定尺度水準が順序尺度水準であると想定した因子分析をしたことになっています。

2464 カテゴリカル因子分析は、段階反応モデルと数学的に等価であることが知られています。出自が違います
2465 ので、分析の数字や出力方法に違いがありますが*³、実質的には同じことをしていると理解してください。ま
2466 た、R ではたった数文字追加するだけで適切な分析方法になるのですから、使わない手はないと思います。

2467 カテゴリカル因子分析は、因子分析の文脈から順序尺度水準の項目を扱えるように、と発展してきたもの
2468 です、これに対して、項目反応理論の文脈から心理尺度を扱えるように、と発展してきたルートもあります。こ
2469 れは**多次元項目反応理論 (multi-dimensional item response theory)** と呼ばれるもので、これを
2470 提供する R パッケージもあります。それが `mirt` パッケージ (Chalmers, 2012) です。最後にこれを使った
2471 コードの例を示しておきます。この計算には時間がかかりますが、完全情報最尤推定によって項目反応理論
2472 で算出するような、精度の高い**因子得点**を得ることができるのは利点だと言えるでしょう。

code : 11.7 mirt の実行

```
2473
2474 1 library(mirt)
2475 2 result.mirt <- mirt(dat, 5)
2476 3 print(result.mirt)
2477 4 result.mirt %>% summary(rotate = "geominQ")
2478
```

2479 11.4 課題

2480 今日の授業でおこなったすべての次の計算をする R コードを提出してください。ファイル形式は R スクリ
2481 プトか Rmd とします。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱い
2482 になります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を
2483 受けてください。

*³ たとえば因子分析の文脈では、因子負荷量を提示して単純構造を目指すのに対し、項目反応理論の文脈では IRCCC を描いたり TIC を描いたりして各反応カテゴリの特徴を精査します。

第 12 章

構造方程式モデリング

本項では**構造方程式モデリング (Structural Equation Modeling, SEM)** について説明していきます*1。この手法はこれまで学んできた因子分析や回帰分析を統合したもので、分散共分散行列に方程式を作り込んでいくものです。我々が調査で得たデータから得られる情報は、分散共分散行列がすべてですから、その変数間関係を細かく作り込んでいき、行列の隅々まで考え尽くすという意味で、究極の手法といえるでしょう。実際、SEM は因子分析や回帰分析の他にも、多くのモデルをその下位モデルとして包含します。言い換えれば、SEM が登場する前に考えられてきたさまざまなモデルが、SEM の枠組みで理解・表現できるようになりました。であれば、この究極の技さえ知っておけば、非常に広範囲に拡張できるわけですから、こんなに便利なことはありません。

このように技術的には非常に高度なものでありますが、これを使うのは意外と簡単にできます。モデルをイメージで表現できる、パスダイアグラムの考え方から入っていきましょう。

12.1 パスダイアグラムの書き方

変数間の関係を表す図を**パス図**あるいは**パスダイアグラム (Path Diagram)** と言います。

分析する際の変数には**観測変数 (Observed Variables)** と**潜在変数 (Latent Variables)** があります。観測変数はその名の通り、観測された変数であり、データとして数値化されたものになります。これに対して潜在変数とは、モデルで仮定された変数のことです。たとえば因子分析やテスト理論では、因子や学力といったものを想定します。それが潜在変数です。潜在変数はデータとして数値化されているのではなく、抽象的な概念ですが、図にする場合はそれも表現しなければなりません。そこで観測変数を矩形(長方形)で、潜在変数を楕円で表現することにします。

また、変数同士の関係を表現する必要がありますが、関係には**因果関係**と**相関関係**があります。調査データから因果関係を見出すのは原理的に不可能ですが、モデル上は一方が他方の原因になっている、と考えることがあります。回帰分析はその典型例で、説明変数が原因となって、被説明変数のあたいが結果的に変わる、という関係ですね。こうした関係を一方向の矢印で表現します。他方、相関関係は因果関係よりも緩やかで、何らかの関係がある、ということの意味しているに過ぎません。因果の向きが定まっているのではなく、どちらの向きにも影響しうる、ということで双方向矢印でもってこれを表現します。

図 12.1 にこの表記方法をまとめました。準備は実はこれだけです。この表記方法を使って、さまざまなモデルが統合的に表現できます。たとえば**回帰分析**を考えてみましょう。説明変数も被る説明変数も観測されている変数を使います。一方が他方の因果的関係にあるので、パスダイアグラムで表現すると図 12.2 の上図

*1 この手法は別名**共分散構造分析**という名前でも呼ばれています。書籍や論文を検索する時は、こちらの名称もキーワードにしてみてください。

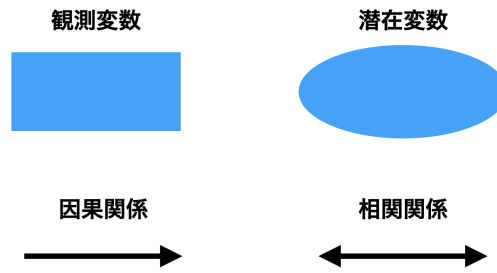


図 12.1 変数と関係の表記法

2513 のようになります。ここで**残差**はデータになく、モデルの計算結果から想定される潜在的な変数なので楕円で
 2514 描かれていることに注意してください。また**重回帰分析**は図 12.2 の下図のようになります。説明変数が複数
 2515 に増えただけですので、表記上はこのように変わるわけですね。

2516 なお、この因果関係を表す矢印の上に (偏) 回帰係数を書いて、影響力の大きさを表現することがありま
 2517 す。構造方程式モデリングでは、この矢印によって表される影響力の強さは**パス係数 (path coefficient)**
 2518 という呼び方で統一されます。

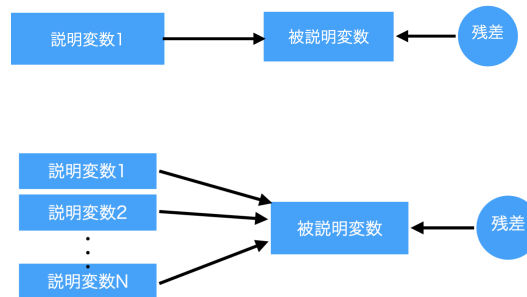


図 12.2 回帰分析と重回帰分析

2519 では**因子分析**をパスダイアグラムで表現するとどうなるでしょうか。複数の項目の背後にある共通要因を
 潜在変数として仮定するので、パスダイアグラムで表現すると図 12.3 のようになります。これは 1 因子モデ

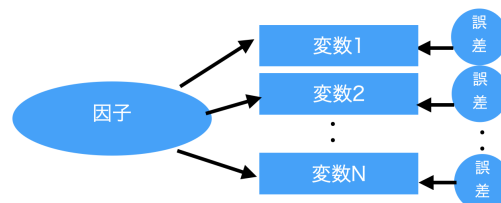


図 12.3 因子分析モデル

2520 ルですが、多因子モデルは基本的に左側にくる潜在変数が複数になるだけです。潜在変数 (因子) からのパ
 2521 ス係数は、**因子負荷量**なのです。また図 12.2 と図 12.3 を見比べると、回帰分析と因子分析も形はよく似て
 2522

2523 いることがわかります。両者の違いは、説明変数が観測変数か潜在変数かであり、因子分析とは説明変数が
2524 潜在変数になった回帰分析モデル、ということもできるのです。

2525 このように、今まで見てきた統計モデルをパスダイアグラムを使って表現できます。

2526 12.2 パスダイアグラムによるさまざまなモデル

2527 パスダイアグラムはさまざまな分析法を統合的に表現するツールであることがわかりました。この方法を
2528 使って他の統計モデルも見てみましょう。

2529 ■主成分分析 このコースではここまで取り上げてきませんでしたが、因子分析とよく似た手法で**主成分分**
2530 **析 (Principle Component Analysis)** というのがあります。主成分分析の目的は観測変数を重みつき
2531 線型結合させ、個々のデータをもっともよく説明するような主成分という合成変数を作り出すことです。たと
2532 えば体力テストなどで複数の競技の記録をつけた後、総合的に誰が一番運動能力があるのか、というときに
2533 この分析が使われたりします。この分析方法をパスダイアグラムで書くと図 12.4 のようになります。

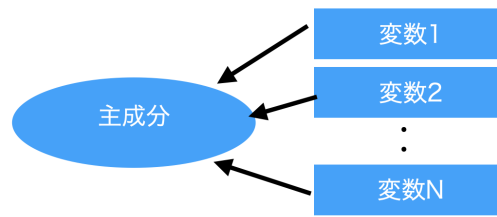


図 12.4 主成分分析モデル

2534 これを見ると、因子分析と似ている、でもちょっと違う、というのがよくわかりますね。因子分析は潜在的な
2535 説明変数を作るのですが、主成分分析は潜在的な被説明変数を作る分析方法なのです。

2536 その違いは、誤差の扱いにも関わってきます。これまでの例にあるように、パスダイアグラムでは基本的に
2537 矢印が入った変数には誤差（残差）がつくのが決まりです。因子分析は説明変数が潜在的だったので、観測
2538 変数にはそれで説明できない残りが潜在変数として構成されました。主成分分析は観測変数が矢印の出所
2539 になりますので、そこには誤差が想定されません。作られた潜在変数は矢印が入ってきた方ですが、これはモ
2540 デルで構成される変数なので誤差を考えることができません*2。主成分分析はデータに誤差を考えないとき
2541 に有用なモデルですから、経済学などデータがはっきりしたものである領域でよく用いられてきました。これに
2542 対して心理統計の領域では、測定されたスコアに誤差が入ると考えるのが当然ですから、因子分析の方がよ
2543 く用いられてきたのです。主成分分析と因子分析は数学的にはよく似た解法で答えを出すのですが、その背
2544 後にある考え方がまったく異なります*3。

2545 ■尺度水準の違いによるモデルの違い たとえば因子分析において観測変数が**順序尺度水準**であれば、
2546 同じパスダイアグラムであってもそれはカテゴリカル因子分析をしたことになります。パスダイアグラムでは変

*2 成分も誤差もモデルから構成されるものなので、モデル上それらを区分する方法がないからです。

*3 もう少し丁寧にいうと、主成分分析も因子分析も計算する上では正方行列の固有値分解をすることになります。ここで、因子分析の場合は項目に誤差を考えますから、固有値分解をするときに相関行列 R ではなく、その対角項に共通性を入れた $R+$ から計算を始めるのでした。主成分分析は誤差を考えませんから、 R から計算したり、単位に意味がある場合が少なくないので分散共分散行列 S の固有値分解でもってパス係数を算出するという違いです。因子分析において共通性の推定を避け、 R を分解する方法をとくに因子分析の主成分分解というのが、初学者を混乱させるポイントになっています。

2547 数の尺度水準を表現することはありませんが、言い換えると尺度水準が違うだけで別の分析モデルになると
2548 いうことです。

2549 たとえば回帰分析において、被説明変数が離散的 (名義尺度水準) なモデルは、かつて**判別分析**
2550 **(Discriminant Analysis)** と呼ばれていました。同様に被説明変数が**順序尺度水準**で得られた場合は、
2551 **順序ロジットモデル (Ordinal Logit Model)** や**順序プロビットモデル Ordinal Probit Model** と呼ば
2552 れたモデルです (2 水準ならロジットモデル, 3 水準以上はプロビットモデル)。

2553 また回帰分析において、説明変数が離散的 (名義尺度水準) なモデルは、**分散分析 (ANalysis Of**
2554 **VAriance; ANOVA)** と呼ばれていたのです。とくに説明変数が 2 カテゴリの場合は **t 検定 (t-test)**
2555 として知られています。これらは平均値差の検定という文脈で説明されてきたかと思いますが、いずれも回帰
2556 分析、もとい**線形モデル**の一種であるというはご存知の通りです。

2557 色々な分析法, 分析モデル名がありますが、**パスダイアグラム**で表現すれば同じである、というのがポイント
2558 です。**構造方程式モデリング**はその下位モデルとしてこれらの分析をすべて含む、あるいは統合的な表現
2559 であるというのはそういう意味です。このように統合される以前は、それぞれの分析方法をそれぞれのソフト
2560 ウェアや関数で実行する必要があったのですが、今は SEM が扱えるソフトウェア、関数 1 つで分析できるの
2561 です。大変ありがたいことではないですか！考えるべきポイントもすべてパス係数や適合度といった統合的
2562 指標でよいのですから。

2563 ■自由なモデルへ 因子分析と回帰分析、あるいはその他のさまざまな統計モデルが、パスダイアグラム
2564 によって統合的に表現できるようになりました。その利点は表現が統一化されただけでなく、同じプラットフォーム
2565 フォームで表現できるので、たとえば図 12.5 のような表現もできるということです。

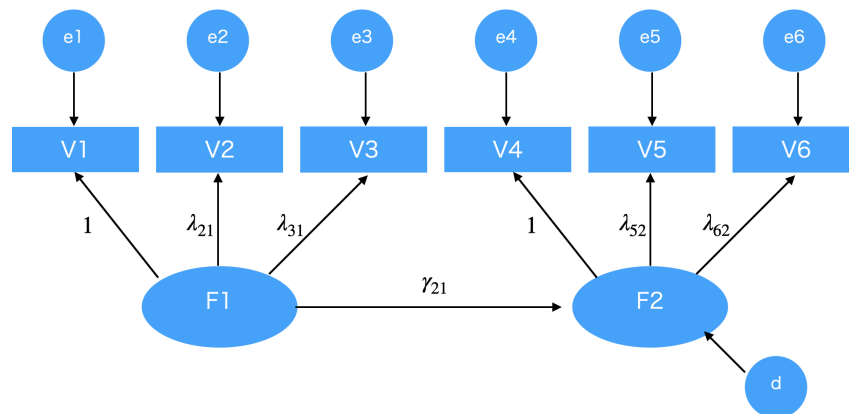


図 12.5 構造方程式と測定方程式。記号の意味は後述

2566 図 12.5 には因子分析がふたつと、因子と因子の関係を矢印で結んだパスが描かれています。潜在変数の
2567 間に回帰分析を組み込んだようなものですね。イメージとしては、ある心理変数が別の心理変数に影響を与
2568 える (ex. 外向性が抑うつ傾向に影響する、など) といったものを表しているわけです。心理学者は心理的変
2569 数がどう影響しあっているか、関係しあっているかをみていきたいわけですから、やりたいことは要するにこ
2570 う複合的なモデルによる検証だったのだ、と言っても過言ではないでしょう。このような分析も、SEM のお
2571 かげで簡単にできるようになりました。ここで観測変数から潜在変数を構成する箇所をとくに**測定方程式**、潜
2572 在変数同志の因果関係などそれ以外の関係を表現する箇所を**構造方程式**と呼んで区別することがありますが
2573 が、**構造方程式モデリング**はその両者からなるモデル全体のこと、ということができるでしょう。

2574 構造方程式モデリングを使うと、調査などから得たデータ全体の見取り図を描いて分析できるのです。言
 2575 い換えると、実際に調査を実施する前に、どの変数とどの変数がどのような関係にあるのか、という仮定を考
 2576 えて調査票をデザインし、過不足なく測定してモデルを検証する、という**仮説検証型**の分析ができるのです。
 2577 それまでの調査は、どういう変数がどういう関係になっているか細かくは分からないけれども、何か関係しそ
 2578 うだからまとめて調査項目にしておいて、因子分析や回帰分析を繰り返して、変数間の関係を**探索的に**研究
 2579 するということがよく行われていました。もちろん今でも探索的研究はあるのですが、既存の尺度や先行研究
 2580 がある場合はより具体的に調査を設計できるようになりました。

2581 因子分析も、因子数がわからないのでスクリープロットなどを描いて探し出し、因子と項目の関係がわか
 2582 らないのですべての因子がすべての項目に影響すると考えて分析する**探索的因子分析 (Exploratory**
 2583 **Factor Analysis, EFA)** がかつては主流でしたが、今ではどの因子がどの項目で測定されるかを明確に
 2584 パスダイアグラムで表現して分析する、**検証的因子分析 (Confirmatory Factor Analysis, CFA)** が
 2585 行われるようになってきています。

2586 12.3 構造方程式モデルによる未知数の推定

2587 12.3.1 未知数は本当に推定できるのか

2588 さて、さまざまなモデルをパスダイアグラムで表現できる、ということがわかりましたが、「方程式」という言
 2589 葉はどこからきているのでしょうか？実は、統計モデルはパスダイアグラムでも表現できますが、同じことを方
 2590 程式でも表現できます。

2591 たとえば図 12.1 で表した回帰分析モデルは、次のような方程式で表せるのでした。

$$y = b_0 + b_1x + e$$

2592 図 12.2 に表した因子分析モデルも、同様に次のような方程式で表せるのでした (項目が 3 つの場合)。

$$\begin{cases} x_1 = \lambda_1 f + e_1 \\ x_2 = \lambda_2 f + e_2 \\ x_3 = \lambda_3 f + e_3 \end{cases}$$

2593 ここで f は因子得点、 λ_j は因子負荷量、あるいはパス係数を表しています。

2594 この因子分析モデルを例に、この方程式をどのように解くかをみていきましょう。そのために、いくつかの記
 2595 号と特徴を確認します。まずベクトルの平均を関数 E で表すことにします。誤差の平均は 0、という特徴は
 2596 $E[e] = 0$ と表すことができます。因子得点も標準化されていると仮定しますので、 $E[f] = 0$ です。つぎ
 2597 にベクトルの分散を関数 V で表すとします。データベクトル x が平均 0、あるいは平均偏差、あるいは標準
 2598 化されていると考えると、 $V[x] = E[(x - E[x])^2] = E[xx']$ は分散共分散行列 Σ だと考えることがで
 2599 きます。また因子は標準化されているので $V[f] = 1$ 、誤差同士は関連しないものの、誤差には分散があり
 2600 ますので $V[e] = E[ee'] = \Psi$ とします。ここで Ψ は対角に誤差分散が入った対角行列です。

2601 さて、さきほどの 3 項目 1 因子モデルを行列で表現してみましょう。各要素を次のようにベクトルでまとめ
 2602 て表現します。

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}, \mathbf{\Lambda} = \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{pmatrix}, \mathbf{e} = \begin{pmatrix} e_1 \\ e_2 \\ e_3 \end{pmatrix}$$

2603 すると式は次のようにあっさり表現できます。

$$\mathbf{x} = \mathbf{\Lambda} \mathbf{f} + \mathbf{e}$$

2604 ここからデータの分散共分散行列 Σ を考えてみましょう。

$$\begin{aligned} V[x] &= E[xx'] \\ &= E[(\Lambda f + e)(\Lambda f + e)'] \\ &= E[(\Lambda f + e)(f\Lambda' + e')] \\ &= E[\Lambda f^2 \Lambda' + \Lambda f e' + e f \Lambda' + e e'] \\ &= \Lambda \Lambda' + \Psi \end{aligned}$$

2605 これをエレメントワイズで書き出すと次のようになります。

$$\begin{aligned} \begin{pmatrix} s_1^2 & s_1 s_2 & s_1 s_3 \\ & s_2^2 & s_2 s_3 \\ & & s_3^2 \end{pmatrix} &= \begin{pmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \end{pmatrix} (\lambda_1 \quad \lambda_2 \quad \lambda_3) + \begin{pmatrix} \sigma_1^2 & & \\ & \sigma_2^2 & \\ & & \sigma_3^2 \end{pmatrix} \\ &= \begin{pmatrix} \lambda_1^2 + \sigma_1^2 & \lambda_1 \lambda_2 & \lambda_1 \lambda_3 \\ & \lambda_2^2 + \sigma_2^2 & \lambda_2 \lambda_3 \\ & & \lambda_3^2 + \sigma_3^2 \end{pmatrix} \end{aligned}$$

2606 この式の左辺はデータから計算された分散共分散行列、右辺はモデルから計算された分散共分散行列で
2607 す。行列の下三角が消えているのは対称行列で同じ情報が入るからで、データからオリジナルに得られるの
2608 は3つの分散 (s_1^2, s_2^2, s_3^2) と3つの共分散 ($s_1 s_2, s_1 s_3, s_2 s_3$) の情報になります。一方右辺での未知数はパ
2609 ス係数 (因子負荷量) の $\lambda_1, \lambda_2, \lambda_3$ と誤差分散 $\sigma_1^2, \sigma_2^2, \sigma_3^2$ ということになります。未知数の数と既知数の数
2610 が同じなので、この方程式は解けます。とくに、未知数と既知数の数が一致しているため、**丁度識別 (just**
2611 **identification)**, あるいは単に**識別可能 (identifiable)** といいます。

2612 分散共分散行列で提供される情報の増え方に比べて、モデルの未知数の増え方は小さいので、一般的な
2613 SEM のモデルのほとんどは**過剰識別**になります。つまりいくつかの答えがあり得る、ということになるので、
2614 その時は尤度が最も高くなるように、といった基準を加えて答えを一意に定めます。逆に未知数の数が多い
2615 場合**識別不可能** (解が求められない) ということになり、その場合はどこかのパラメータを値として入れてや
2616 る (**制約**をかける、といいます) 必要があります。

2617 12.3.2 複雑なモデルの場合

2618 では次に、図 12.5 のような複雑なモデルでも方程式で表せることを示しましょう。

2619 ここで図 12.5 の係数について説明します。ギリシア文字 λ や γ がパス係数を表しています。パス係数の
2620 添字は、被説明変数の変数番号 j と、説明変数の変数番号 k をつかって λ_{jk} のようにかき、これで $k \rightarrow j$
2621 の影響力の大きさを表していると思ってください。測定方程式の係数を λ で、構造方程式の係数を γ で表現
2622 しました。また $\lambda_{11}, \lambda_{42}$ になるべきところが 1 になっていますが、これは「潜在変数から出るパスのうち、1 つ
2623 は必ず 1 にしないと推定できない」という数学的特徴があるからです。決まりごとだと思ってください。

2624 それを踏まえて、モデルの方程式を書くと、次のようになります。

$$\begin{cases} V_1 &= f_1 + e_1 \\ V_2 &= \lambda_{21} f_1 + e_2 \\ V_3 &= \lambda_{31} f_1 + e_3 \\ V_4 &= f_2 + e_4 \\ V_5 &= \lambda_{52} f_2 + e_5 \\ V_6 &= \lambda_{62} f_2 + e_6 \\ f_2 &= \gamma_{21} f_1 + d \end{cases}$$

2625 ややこしいだけで、書けるのは書けましたね！そして当然、これは行列表現した方が楽ではないか、というア
2626 イデアが浮かぶと思います。その通りで、エレメントワイズにまずは書いてみましょう。

$$\begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \\ V_5 \\ V_6 \\ f_2 \\ f_1 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \lambda_{21} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \lambda_{31} \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \lambda_{52} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & \lambda_{62} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & \gamma_{21} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \\ V_5 \\ V_6 \\ f_2 \\ f_1 \end{pmatrix} + \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ d \\ f_1 \end{pmatrix}$$

2627 左辺のベクトルの最後に f_1 が入っているところがポイントで、これは方程式的には $f_1 = f_1$ というだけの
2628 式なのですが、これを入れたおかげで全体を整合的に表現できましたね。この左辺のベクトルを v として、式
2629 全体を $v = Gv + e$ と考えてみましょう。すると次のように展開できます。

$$\begin{aligned} v &= Gv + e \\ v - Gv &= e \\ (I - G)v &= e \\ (I - G)^{-1}(I - G)v &= (I - G)^{-1}e \end{aligned}$$

2630 ここで $P = (I - G)$ とすると、

$$v = P^{-1}e$$

2631 ということになります。これをエレメントワイズで書くと次のようになります。

$$\begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \\ V_5 \\ V_6 \\ f_2 \\ f_1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & -\lambda_{21} \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & -\lambda_{31} \\ 0 & 0 & 0 & 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & -\lambda_{52} & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & -\lambda_{62} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & -\gamma_{21} \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ e_5 \\ e_6 \\ d \\ f_1 \end{pmatrix}$$

2632 さて、変数からなるベクトル v ではありますが、観測変数は V_1 から V_6 までしかありませんから、そこだけ
2633 を取り出す行列 $F = \{I, O\}$ を考えます。ここでは次のような行列です。

$$F = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{pmatrix}$$

2634 これを使って、次のように関係を書き直すことができます。

$$g = \begin{pmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \\ V_5 \\ V_6 \end{pmatrix} = FP^{-1}e$$

2659 ■GFI・AGFI GFI(Goodness of Fit Index) や AGFI(Adjusted GFI) は、0 から 1 までの値
2660 を取る適合度指標で、因果モデルがデータの何 % を説明したかの指標になります。1 に近いほど良いモデル
2661 で、GFI で 0.95,AGFI で 0.90 以上の数字であると良い適合と判断されます。

2662 ■情報量規準 情報量規準 (Information Criteria) とは、一般的な統計モデルを評価するための指標
2663 で、AIC(Akaike's Information Criterion) や BIC(Bayesian I.C.) などがそれにあたります。
2664 これは今検証しているモデルが、データによくフィットする「本当の」モデルからどの程度離れているかを示す
2665 指標です。本当のモデル、というのは分かり得ませんから、複数のモデル A,B,C を同じデータに当てはめた
2666 とき、AIC や BIC が相対的に小さいものが、より良く適合していると判断します。相対的な指標ですから、
2667 違う変数を扱うモデル間の比較はできないことに注意が必要です。

2668 ■RMSEA モデルがデータとどの程度乖離しているか、を直接表現した指標が RMSEA(Root Mean
2669 Square Error of Approximation) です。近似した誤差に対する平方平均平方根、ということですね。
2670 一般に 0.05 より小さければ、そのモデルは当てはまりがよく、0.1 以上であれば当てはまりが悪いと判断し
2671 ます。

2672 ■CFI/TLI CFI(Comparative Fit Index) や TLI(Tucker-Lewis Index) は、観測変数間に
2673 まったく相関がないという非現実的なモデル(独立モデルという)に比べて、当該のモデルがどの程度よいも
2674 のかを指標化したものです。いずれも 1.0 に近ければ近いほど良いモデルとされています。一般にこの数値
2675 が 0.9 以上になることを目指してモデルを改訂していきます。

2676 ■SRMR SRMR(Standardized Root Mean square Residual) は標準化された残差平方平均
2677 平方根を表します。これはモデルで説明できなかったものがどれほどあったか、を表しますので、0 に近けれ
2678 ば近いほど当てはまりの良いことを表します。

2679 ■修正指数 最後に、修正指数 (modification index) を紹介しておきます。上で述べたように、検証し
2680 ようとした統計モデルがどの程度当てはまっているか、ということの数値化できるようになったわけですが、常
2681 に最初にたてた仮説モデルが最適である、ということは滅多にありません。データから仮説の訂正を余儀なく
2682 されることがある、あるいはもっと良くなったり新しい発見があったりすることが、統計モデリングの醍醐味で
2683 もあります。修正指数は、どこをどう改良すれば指標がよくなりますよ、というヒントを与えてくれるものです。
2684 もっとも、適合度を上げることだけが目的なのではないことに注意してください。適合度を上げるために、意味
2685 のわからないモデルを作り上げたところで、モデルの意味がないのですから。

2686 構造方程式モデリングは、非常に複雑な方程式を解いているんだな、ということは理解していただけかと思
2687 います。次回はそんな複雑なモデルでも、R で簡単に分析できるよということを、演習を交えて理解してい
2688 きましょう。

2689 12.5 課題

2690 2 項目 1 因子モデルの方程式、およびそれを行列形式で表したものを書きましょう。

2691 また行列の要素で書いた場合、既知数と未知数はどちらが多くなるでしょうか。言い換えれば、このモデル
2692 は識別できるでしょうか。式を展開して確認してください。

第 13 章

R による構造方程式モデリング

これまでの流れと同じで、統計技術の理論を知っただけではなく、自分で実際に計算できる演習を経てこそ理解が深まります。本講では R をつかって実際に構造方程式モデリングを解くことを演習的に学んでいきましょう。構造方程式モデリングを実装するパッケージは複数あるが、最も応用範囲がひろい lavaan パッケージを用いることにします。授業を始めるにあたって、まずは lavaan パッケージをインストールしておいてください*1。

13.1 モデル式の入力

13.1.1 観測変数だけのモデル

まずは観測変数同士の関係をパスでつなぐモデルをみてみましょう。観測変数同士のパスですから、(重) 回帰分析をやるようなものと同じです。今回のデータとして lavaan パッケージに含まれている HolzingerSwineford1939 データセットを使います。これは (Holzinger and Swineford, 1939) のデータで 301 人に対して行われた 15 の能力テストスコアの一部が入ったデータです。一部を 13.1 に示しました。

表 13.1 Holzinger and Swineford(1939) のデータ

id	sex	ageyr	agemo	school	grade	x1	x2	x3	x4	x5	x6	x7	x8	x9
1	1	13	1	Pasteur	7	3.33	7.75	0.38	2.33	5.75	1.29	3.39	5.75	6.36
2	2	13	7	Pasteur	7	5.33	5.25	2.12	1.67	3.00	1.29	3.78	6.25	7.92
3	2	13	1	Pasteur	7	4.50	5.25	1.88	1.00	1.75	0.43	3.26	3.90	4.42
4	1	13	2	Pasteur	7	5.33	7.75	3.00	2.67	4.50	2.43	3.00	5.30	4.86
5	2	12	2	Pasteur	7	4.83	4.75	0.88	2.67	4.00	2.57	3.70	6.30	5.92
6	2	14	1	Pasteur	7	5.33	5.00	2.25	1.00	3.00	0.86	4.35	6.65	7.50

変数の意味は以下の通りです。x1 から x3 が空間的知覚，x4 から x6 が言語的能力，x7 から x9 が移動する物体の認識力を測るものになっています。

id 被験者 ID
sex 性別 (男性=1, 女性=2)
ageyr 生まれ年

*1 R のコンソールに `install.packages("lavaan")` と入力するか、RStudio のパッケージタブから入力しましょう。

2711 agemo 生まれ月
 2712 school 所属校
 2713 grade 学年
 2714 x1 視覚的知覚
 2715 x2 立方体
 2716 x3 菱形
 2717 x4 段落の理解
 2718 x5 文章完成
 2719 x6 言葉の意味
 2720 x7 加速
 2721 x8 点を数える
 2722 x9 直線・曲線文字の区別

2723 さて、ではまず回帰分析をやってみることにしましょう。変数 x4 を x5,x6 で回帰する重回帰モデルを
 2724 lavaan を実行するには code:13.1 のように書きます。

code : 13.1 重回帰モデル

```
2725
2726 1 model1 <- "
2727 2 x4_~_x5+_x6
2728 3 "
2729 4 result1 <- sem(model1, data = HolzingerSwineford1939)
2730 5 summary(result1, fit.measures = T)
2731 6 # 比較
2732 7 result1.1 <- lm(x4 ~ x5 + x6, data = HolzingerSwineford1939)
2733 8 summary(result1.1)
2734
```

2735 ■コード解説

2736 1-3 行目 モデルの記述
 2737 4 行目 関数による推定
 2738 5 行目 結果の出力
 2739 6-8 行目 従来の lm 関数による当てはめ例

2740 ここでモデルを書くところは、ダブルクォーテーションで括られていることに注意してください。このようにす
 2741 ることで、R に複数行からなるモデルを渡すことが可能になります。実際のモデルは 2 行目だけで、これは
 2742 lm 関数に書いているモデル式の形と同じであることがわかります。モデルを書いて、モデルとデータの組みを
 2743 sem 関数に渡すと結果が帰ってくる、これだけです。結果を出力するときに、オプションとして fit.measures
 2744 を書いてあるところだけが違います。

2745 出力結果の一部を出力 13.1 に示します。実際はもっと色々でているのですが、CFI , TLI, AIC,BIC,
 2746 RMSEA, SRMR などそれぞれ適合度指標を荒らしています。推定結果は Regression: の Estimate
 2747 のところを見てください。これらは、lm 関数と同じ係数になっていることが確認できます。

R の出力 13.1: 観測変数だけのモデル

(中略)

User Model versus Baseline Model:

Comparative Fit Index (CFI)	1.000
Tucker-Lewis Index (TLI)	1.000

Loglikelihood and Information Criteria:

(中略)

Akaike (AIC)	673.134
Bayesian (BIC)	684.255
Sample-size adjusted Bayesian (BIC)	674.741

Root Mean Square Error of Approximation:

RMSEA	0.000
-------	-------

(中略)

Standardized Root Mean Square Residual:

SRMR	0.000
------	-------

Parameter Estimates:

Standard errors	Standard
Information	Expected
Information saturated (h1) model	Structured

Regressions:

	Estimate	Std.Err	z-value	P(> z)
x4 ~				
x5	0.423	0.047	8.958	0.000
x6	0.390	0.056	7.002	0.000

Variances:

	Estimate	Std.Err	z-value	P(> z)
.x4	0.537	0.044	12.268	0.000

2748

2749 このように、重回帰モデルをするのであれば記述方法は同じです。違いはさまざまな観点からの適合度指
2750 標モデルということ。加えて、このモデルはどんどん拡張していくことができる点にあります。

2751 たとえばパス解析 (Path Analysis) というのがあります。これは影響力のルートが $x \rightarrow y \rightarrow z$ のよう
2752 に、次から次へと繋がっていくモデルです。SEM ができるまでは、まず $x \rightarrow y$ を回帰分析し、次に $y \rightarrow z$ を
2753 分析する、という方法でした。しかしこれでは R^2 など適合度を逐一確認する必要があり、またモデル全体の
2754 評価ができないという欠点があります。しかし SEM ではコード code: 13.2 のように書くだけです。

code : 13.2 パス解析のモデル

```

2755 1 model2 <- "
2756 2 x4~x5
2757 3 x5~x6
2758 4 x6~grade
2759 5 "
2760
2761 6 result2 <- sem(model2, data = HolzingerSwineford1939)
2762 7 summary(result2, fit.measures = T, standardized = T)
2763

```

このモデル式 (2-6 行目) にあるように, $grade \rightarrow x6 \rightarrow x5 \rightarrow x4$ という一連の影響力を分析し, 一気に適合度を表現してくれています。出力結果には適合度指標のほか, すべての変数を標準化した標準化係数を出すオプションを追加しました。

これを図にしたのが図 13.1 です。ちなみにこの図も R のパッケージで書きました。パス図を自動的に書いてくれるパッケージとして `semPlot` パッケージ (Epskamp, 2021) や `tidySEM` パッケージ (Van Lissa, 2019) などを使います。ここでは `tidySEM` パッケージの `graph_sem` 関数を使いました。

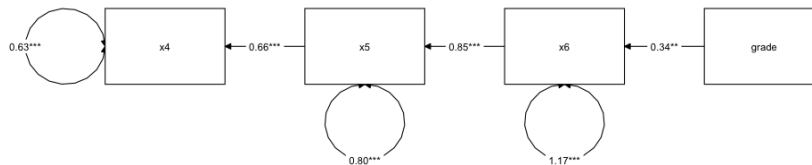


図 13.1 パス解析の図

2770 13.1.2 測定方程式を入れたモデル

2771 今度は測定方程式を入れたモデルを考えてみましょう。回帰モデルは従属変数と独立変数をチルダ (~) で
2772 つなぎましたが, 潜在変数を作る場合は ~ で右辺に観測変数, 左辺に潜在変数をセットします。コード例を
2773 code:13.3 に示します。

code : 13.3 因子分析モデル

```

2774 1 model3 <- "
2775 2 visual~x1+x2+x3
2776 3 textual~x4+x5+x6
2777 4 speed~x7+x8+x9
2778 5 "
2779
2780 6 result3 <- sem(model3, data = HolzingerSwineford1939)
2781 7 summary(result3, fit.measures = T, standardized = T)
2782

```

2783 このモデル図は図 13.2 のようになります*2。モデル上とくに指定はしてありませんが, 潜在変数同士の
2784 相関係数も自動的に推定されています。パス係数は潜在変数からでてきていますから, **因子負荷量**のように
2785 考えることができますね。

2786 このモデルは因子分析の一種ですが**探索的因子分析**とはいくつかの点で違います。1 つは因子数が 3 で
2787 ある, と最初から決めていた点。2 つ目は因子負荷量ですが, 探索的因子分析の場合ほどの項目に対しても
2788 共通因子からの因子負荷量が計算されていました。今回の例では, x4 から x9 の変数にすいても visual 因

*2 この図は小野島 (2021) の関数を使っています。tikz を使って綺麗なモデルが描ける素晴らしい関数の提供に感謝。

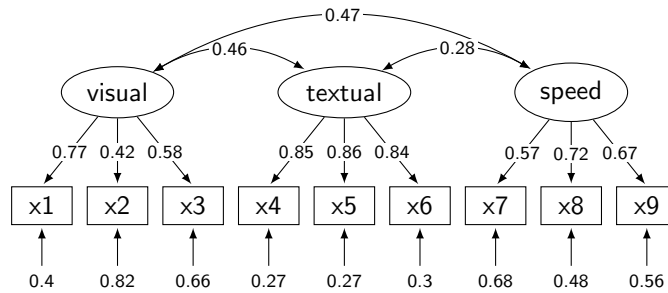


図 13.2 潜在変数モデル

2789 子からのパスが出ている、というモデルだったのです。今回のモデルではそれらのパスはなく、visual 因子
 2790 は x1,x2,x3 にしか影響していません。言い換えると、x4 から x9 に対する visual 因子のパス係数は 0 であ
 2791 る、と特定したようなものです。

2792 探索的因子分析の場合、理想的には「因子が関係する項目には十分影響しているが、関係しない項目
 2793 にはまったく影響しない」という**単純構造の原理 (Principle of Simple Structure)** が理想とされて
 2794 いました。SEM を使うとこの理想的なモデルが合っているかどうかを検証する = モデル上で制約をかけ
 2795 てその適合度を評価する、という使い方ができます。このような方針の因子分析のことを**検証的因子分析**
 2796 (**Confirmatory Factor Analysis**) と呼んで、探索的因子分析と区別します。

2797 検証的因子分析をすることは、関係のない項目への影響を 0 とする、という強い仮定を置いていることに
 2798 もなりますが、このことによって**因子的妥当性**や**収束的妥当性**、**弁別的妥当性**を検討している、と考えること
 2799 もできます*3。もしこの理論的なモデルの当てはまりが悪ければ、異なる因子負荷のパターンを考えなければ
 2800 なりません。あるいは因子数を変えてモデルを組まなければならないかもしれないのです。因子数は同じだけ
 2801 ど因子負荷量が違う、といったモデルを考えることもできますし、因子負荷量も固定して「この値に違いない」
 2802 としたモデルを検証する、ということもできます。モデルを当てはめるデータは違っても構いません。むしろ違
 2803 うデータに同じモデルが当てはめられ、十分な適合度があればそれはそのモデルの普遍性があるということ
 2804 で、いいことなのです。

2805 このように、SEM では同じモデルを男性と女性、学生と社会人、日本人データと外国人データといった複
 2806 数のデータに当てはめて、どこまで同じでどこから違うか、と言った比較をできます。こうしたモデル比較のこ
 2807 とを**多母集団同時分析 (Multi-group analysis)** といいます。

2808 13.1.3 構造方程式モデルへ

2809 それでは潜在変数同士の関係に、更なる仮定を入れたモデルを考えてみましょう。

code : 13.4 構造方程式モデル

```
2810 1 model4 <- "  
2811 2   visual =~ x1 + x2 + x3  
2812 3   textual =~ x4 + x5 + x6  
2813 4   speed =~ x7 + x8 + x9  
2814 5   textual =~ visual + speed
```

*3 因子的妥当性とは、因子が十分に大きく項目に影響していることです。収束的妥当性とは、因子が関係しない項目に影響しないこと、弁別的妥当性は因子同士が別のものであると区別できることを表しています。

```

2816 6  visual~grade
2817 7  "
2818 8  result4 <- sem(model4, data = HolzingerSwineford1939)
2819 9  summary(result4, fit.measures = T, standardized = T)
2820 10 modificationindices(result4) %>%
2821 11   as_tibble() %>%
2822 12   arrange(-mi)
2823

```

2824 ■コード解説

2825 1-7 行目 モデルの記述。visual, textual, speed の 3 因子をつくり, textual は visual と speed から説明されるモデル。さらに speed 因子は学年によっても説明される。

2827 8 行目 関数による推定

2828 9 行目 結果の出力。適合度指標と標準化係数の出力オプションをつけて。

2829 10-12 行目 分析結果の修正指数を出力。ただし指標 mi の大きい順に並べ替えるために 11 行目で出力を tibble 型にし, arrange 関数で並べ替えている。

2831 まずは今回のモデルについても, 図でみたほうがわかりやすいでしょう。結果を合わせて図 13.3 に示してみました。

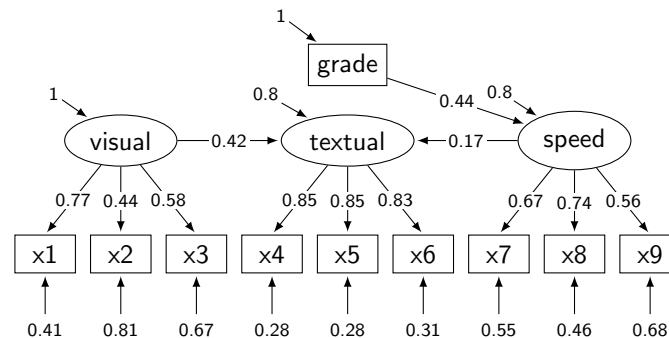


図 13.3 構造方程式モデルへ

2832

2833 この図の出力と, R での出力はどこが対応しているか, それぞれ確認しましょう。

2834 まずは測定方程式のところです。因子負荷量, すなわち潜在変数からのパスはすべての変数を標準化したところの数字ですので, Std.all の列を参照します。この因子で説明できなかった独自成分の大きさは, Varianvces のところに表されています。

2837 構造方程式として, textual 因子が visual 因子と speed 因子に説明されており, そのパス係数が Regressions: の Std.all 列に示されています。speed 因子はさらに grade 変数にも説明されているので, そのパス係数も確認できます。

2840 SEM の基本として, パスが入った (影響を受けた) 変数は, 影響を受けた部分と受けなかった部分に分割されます。影響を受けなかった部分が残りの分散として, Variances: のところに示されています。

2842 visual 因子は説明変数になっていますが, パスが入ってきませんので, この分散は 1.00 です。図では visual 因子の楕円の上に 1 からの影響が入っていますが, これは分散が 1 だということの意味です。同じく

2844 grade 変数も説明されない変数*4ですので、この分散が1(標準化されています) になっているのです。

R の出力 13.2: 構造を入れたモデルの出力結果 (一部)

```

Latent Variables:
      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
visual =~
  x1          1.000
  x2          0.573    0.109    5.249    0.000    0.515    0.437
  x3          0.724    0.124    5.846    0.000    0.651    0.576
textual =~
  x4          1.000
  x5          1.109    0.067   16.646    0.000    1.085    0.851
  x6          0.924    0.056   16.349    0.000    0.903    0.834
speed =~
  x7          1.000
  x8          1.022    0.130    7.837    0.000    0.744    0.737
  x9          0.782    0.108    7.273    0.000    0.569    0.564

Regressions:
      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
textual ~
  visual      0.453    0.094    4.841    0.000    0.416    0.416
  speed       0.226    0.093    2.435    0.015    0.168    0.168
speed ~
  grade       0.645    0.105    6.147    0.000    0.887    0.443

Variances:
      Estimate Std.Err z-value P(>|z|) Std.lv Std.all
.x1          0.554    0.132    4.187    0.000    0.554    0.407
.x2          1.121    0.103   10.841    0.000    1.121    0.809
.x3          0.851    0.097    8.769    0.000    0.851    0.668
.x4          0.370    0.048    7.706    0.000    0.370    0.279
.x5          0.448    0.059    7.635    0.000    0.448    0.276
.x6          0.358    0.043    8.268    0.000    0.358    0.305
.x7          0.658    0.080    8.178    0.000    0.658    0.554
.x8          0.466    0.072    6.478    0.000    0.466    0.457
.x9          0.695    0.069   10.021    0.000    0.695    0.682
visual       0.807    0.161    5.029    0.000    1.000    1.000
.textual     0.764    0.096    7.918    0.000    0.798    0.798
.speed       0.425    0.083    5.130    0.000    0.804    0.804

```

2845

2846 最後に出力 13.3 にある、修正指数 (modification index) の出力を見てみましょう。一行目にあるの
2847 は、visual 因子は変数 x9 へのパスをつけると、適合度がぐっと上がるよ、ということを意味しています。

*4 説明されない変数のことをとくに外生変数 (exogenous variable) といいます。これに対して説明されることがある変数は内生変数 (endogenous variable) といいます。

R の出力 13.3: 修正指数 (一部)

```

# A tibble: 64 x 8
  lhs      op      rhs      mi      epc sepc.lv sepc.all sepc.nox
  <chr>   <chr> <chr>   <dbl> <dbl> <dbl>   <dbl>   <dbl>
1 visual ==~   x9      41.9 0.442  0.397   0.394   0.394
2 visual ~     speed  25.3 0.507  0.411   0.411   0.411
3 visual ~     textual 25.3 2.24   2.44    2.44    2.44
4 textual ==~   x1      16.6 0.490  0.479   0.411   0.411
5 visual ~~    speed  13.6 0.184  0.314   0.314   0.314
6 speed  ~     visual  13.6 0.228  0.282   0.282   0.282
7 visual ~     grade  12.8 0.445  0.495   0.247   0.495
8 grade  ~     visual  12.8 0.137  0.123   0.247   0.247
9 speed  ==~   x1      11.1 0.313  0.227   0.195   0.195
10 x1    ~~    x9      11.0 0.169  0.169   0.272   0.272

```

2848

2849 実際に visual 因子が変数 x9 に影響している (x9 から因子が構成されている) というモデルを作り、適合度指標を比べてみましょう。表 13.2 に代表的な適合度指標の修正前の値と修正後の値を示しました。い

表 13.2 修正前後での指標の変化

index	before	after
CFI	0.895	0.945
TLI	0.856	0.923
AIC	7479.335	7433.224
BIC	7557.115	7514.707
RMSEA	0.100	0.073
GFI	0.919	0.944
AGFI	0.866	0.905

2850

2851 れの適合度指標においても、大きな改善が見られています。さて、これをどう考えたらいいでしょうか。

2852 13.2 実践上の注意点

2853 ここまでで、さまざまなモデルを表現できること、それを統一的に評価できることが理解できたかと思いま
2854 す。とくに、統一的に評価できることでさまざまな**モデル比較**も簡単になり、さらに適合度が良いモデルにする
2855 ためにはどうすれば良いかについて、指標も出てくることがわかりました。

2856 分散共分散行列の隅々まで記述できる大きな力を手に入れたことは間違いないのですが、大事なのは、
2857 我々は構造方程式モデルを使う側であって、それに使われる側であってはならない、ということです。学会誌
2858 に掲載されるような論文で、構造方程式モデルを見ると、その適合度は CFI, TLI, AGFI が 0.9 を超えてい
2859 るような、とても適合度の良いモデルがほとんどです。適合度が悪いモデルだと、これはデータとモデルがあっ
2860 ていないということですから、モデルの改善を求められたり掲載されなくなったりするということがあるでしょ
2861 う。しかし大事なのは、モデルを使って主張したい心理学的内容のほうであって、適合度が良いだけでの中身
2862 のないモデルではないはずなのです。

2863 今回も機械的に visual 因子が x9 変数から構成されるようにすると、適合度が上がるという提案を受けて

2864 実施してみたところ、確かに適合度は上昇しました。ところでこの x_9 変数というのは直線や曲線のスピードを
2865 認識するテストであって、(静的な) 空間認知能力とは違うものではないでしょうか？ そもそもそういうつも
2866 りで作ったものではないのに、機械に指摘されて「ひよっとしたらそういう側面があったのかも。そうだ、最初
2867 からそう思っていたに違いない」というように、自分を無理やり納得させてはいないでしょうか？

2868 心理学の場合は変数の多くが関係しあっていますから、「移動する直線を見るというのは空間認知とも考
2869 えられるのだ。少なくともデータはそう示している」という理屈が成り立つかもしれません。しかし結果を見て
2870 から考え方を考えるのは、適切な研究実践法ではないでしょう。これらの問題は結果の再現性が担保されな
2871 いという心理学の危機の引き金になりました*5。自分に都合の良い結果やモデルを出すことが目的ではなく、
2872 心がどのような状態になっているのかについての理論的積み重ねや発展こそ、目的であったことを忘れては
2873 いけません。

2874 とくに構造方程式まで使えるようになると、心理学的実体同士の関係を描写し、モデル化できるということ
2875 が魅力的に映るかもしれません。心理学者は、これこそがやりたかったことなのかもしれないですね。しかし
2876 データを超えての解釈はご法度ですし、何より構成された潜在変数が心理的な実在であるかどうかは、改め
2877 て考えなければならないのです。これらはあくまでも分散共分散行列から算出されるでしかなく、「文章読解
2878 力」「空間把握力」といった次元に貼り付けた自分たち自身の命名法に酔って、思考停止するようなことがあ
2879 てはなりません。

2880 この潜在変数同士の影響力が心理学的にどういう意味があるのか、をしっかりと考えてから、モデルでの検
2881 証に進まなければならないことに注意して使ってください。

2882 13.3 そのほかの統計パッケージ

2883 構造方程式モデリングの利点の 1 つは、モデルを可視化したことにあります。皆さんもモデルの図を見た方
2884 が、出力結果を見るよりも理解が進んだ気がするでしょう。

2885 こうした構造方程式モデリングを実行するソフトウェアは、R の専売特許ではありません。たとえば **Amos**
2886 という IBM 社が出しているソフトウェアでは、マウスをクリックしながら統計モデルを作り分析していくことが
2887 できます。モデルの構成から GUI でできるのは大変な利点です。

2888 また、構造方程式モデリングはさまざまな分析手法の統合的ツールです。言い換えるとかなり複雑なモデ
2889 ルであっても、ゴリゴリ計算し分析してくれます。現在考えられているさまざまなモデルの、ほぼすべてにつ
2890 いて計算できるソフトウェアが **Mplus** です。このソフトウェアで分析できない SEM はない、と言っている
2891 ほどその用途は広く、また計算スピードも爆速です。カテゴリカルな変数にももちろん対応していますから、
2892 IRT/GRM のような出力もできます。

2893 R の利点は商用ソフトと違ってフリーで手に入るところですが、有用・有償・商用パッケージでも良いので
2894 あれば、これらのソフトも活用することを考えてもいいでしょう。また R でも **lavaan** パッケージの他に、**sem**
2895 パッケージというのがあります。ツールは色々あって、ユーザがそれを選べるようにことが理想的ですから、皆
2896 さんも興味があれば色々試してみましょう！

2897 13.4 課題

2898 今日の授業でおこなったすべての次の計算をする R コードを提出してください。ファイル形式は R スクリ
2899 プトか Rmd とします。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱い

*5 さきほどの良い結果しか雑誌に掲載されない問題のことを、**出版バイアス (publication bias)** の問題といい、今も問題視されています。

2900 になります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を
2901 受けてください。

第 14 章

双対尺度法

さて前回 SEM を学んだことで、分散共分散行列をとことんまで分析し尽くす方法が手に入ったことになります。SEM は統合的な表現方法ですから、観測変数でも潜在変数でもいいですし、順序尺度水準以上の大小関係が表現できる数値データであれば、あらゆる表現ができるわけです。

ではこれで多変量解析はすべて理解したことになるのでしょうか。いえ、もちろんそうではありません。分散共分散行列の分解はできるようになりましたが、変数間関係の表現は分散共分散行列だけではありません。ここまで、データの持つ情報は分散がすべて、変数間関係は共分散で表されるから分散共分散行列がすべてだ、と言わんばかりに話をしてきましたが、その枠を外すとどうなるのでしょうか。

14.1 直線的ではない関係

ご存知の通り、分散共分散行列を標準化した行列は相関行列といいますが、相関係数が 1.0(あるいは -1.0) の状態を考えれば明らかなように、分散共分散行列 (や相関行列) で表現されるのは変数の直線的な関係性に限った話です。心理学の場合は中庸が良いようなシーンも少なくありません。たとえば血圧と健康度の関係で言えば、低血圧でも高血圧でも不健康であり、ちょうどいい血圧が一番健康的だというのはすぐにわかります。このように中庸が良い場合は、散布図が U 字型に現れてきます。散布図が U 字型であれば、相関係数としては 0 近くになりますが、血圧と健康が無関係だとは誰も言えないでしょう。

ごくシンプルかつ具体的な例をあげてみましょう。血圧と頭痛の頻度について調査したとします。血圧は高い、普通、低い の 3 段階。頭痛の頻度は「ない」、「たまに」、「ときどき」、「いつも」の 4 段階です。表 14.1 のような結果を見ると、この集計表に直線的な関係は確かなさそうですね。でも関係ないわけではない、と思いませんか。

表 14.1 血圧と頭痛

血圧	ない	たまに	ときどき	いつも
高い	0	0	3	2
普通	5	0	0	0
低い	0	2	3	1

このような場合は、「血圧が高い人か低い人は、頭痛がある」という傾向が見て取れるはずですが。しかしこのデータ、機械的に血圧が高いを 1、普通を 2、低いを 3 とし、同様に頭痛もない～いつもを 1,2,3,4 として相関係数を計算すると、 $r = 0.165$ になります。相関関係では傾向をうまく読み取れていないのです*1。

*1 相関係数でお話ししましたが、分散共分散行列でも同じです。共分散で表現すると単位のせいで関係が直接的には理解できま

2925 その根本的な原因は、もちろんスコアリングの方法にあります。共分散を計算するには間隔尺度水準以上
2926 の情報が必要で、今回のような 3,4 件法では相関係数を計算して良い数字ではありません。ではポリコリッ
2927 ク相関係数のように順序尺度水準の相関係数なら良いか、といたいところですが、それでもまだ不十分で
2928 す。なぜなら、血圧が高い方から低い方まで、順番が保存されてしまっているからです。

2929 関係を導き出すには抜本的な改革が必要です。たとえば表 14.2 のようにすればどうでしょうか。こうすると
2930 左下から右上にかけての直線的な関係がまだ見えてきます。これは血圧の順番を「高い」「低い」「普通」に並
べ替えたものです。また、頭痛の順番も「いつも」と「ときどき」をひっくり返しました。順番を並び替えてしまっ

表 14.2 並び替えられた「血圧と頭痛」

血圧	ない	たまに	いつも	ときどき
高い	0	0	2	3
低い	0	2	1	3
普通	5	0	0	0

2931

2932 たというのは、尺度水準的には数字を無視して「高い」「低い」「普通」というカテゴリとして扱ったということ
2933 になります。つまり**名義尺度水準**レベルにまで落として考えたのです。

2934 このように並び替えて、血圧のスコアを高い → 3, 低い → 2, 普通 → 1 とし、また頭痛の頻度をいつも
2935 → 3, ときどき → 4 と付け替えて相関係数を計算すると、今度は $r = 0.810$ になりました。こちらの方が線形
2936 性は高いことが数字でも確認できました。

2937 ここで行ったのはこのように、すべての反応カテゴリをただの言葉だと考えて、付与された数字を無視して
2938 並び替えたことになります。そうすることで、左下から右上にかけて線形性を高めることができました。しかし
2939 「高い」と「低い」、「低い」と「普通」の配置も等間隔である必要はなく、もっと直線性がはっきりするように配置
2940 してやってもよいのでは、というアイデアが浮かびます。

2941 ここで考え方が反転したことに注意してください。一般的なリッカート尺度では、反応カテゴリに (シグマ法
2942 などで) 数字をつけて、変数間の線形性を算出して意味を考えるのでした。ここでは逆に、変数間の線形性
2943 を最大にするように反応カテゴリに数字をつけてやろう、という考え方です。データの直線性というのはデー
2944 タの特徴を最も強調し解釈しやすい形です。そのようにデータを整えるためには、反応カテゴリにどういう数字
2945 を付与すれば良いか、と考えるのです。この方法を**数量化の理論 (Quantification Methods)** とい
2946 います。

2947 反応カテゴリに数字を与えるとき、名義尺度水準にまで落として、すなわち数字の意味を無くして並び替え
2948 るような作業をします。もちろん分析対象が、初めから名義尺度水準であっても構いません。たとえば出身県
2949 を調査したようなデータがあったとして、他の変数との線形関係を最大にするように並び替え、数字を付与す
2950 ることもできます。数量化の考え方は、名義尺度水準の変数に解釈しやすい数値を与えることも言えるの
2951 です。

せん。

14.2 林の数量化理論

数量化の研究をしたのは、日本の偉大な統計学者、林知己夫 (はやし ちきお)^{*2}という人です。その名を冠して林の数量化理論と呼ばれることがあります。

林はさまざまな研究成果を挙げており、論文ごとにデータに必要な分析方法を開発するというようなスタイルでした。その膨大な研究業績を、弟子である鮑戸弘^{*3}が分類して、同じような分析方法ごとに I 類, II 類, III 類, IV 類, と呼んでいきました。

表 14.3 林の数量化

手法	外的基準	データ	目的	関連する手法
数量化 I 類	量的変数	質的変数	外的基準の予測	重回帰分析
数量化 II 類	質的変数	質的変数	外的基準の判別	判別分析
数量化 III 類	なし	質的変数	変数間の関係の要約と記述	正準相関分析
数量化 IV 類	なし	対象間の類似度	対象間の関係の要約と記述	多次元尺度法

表 14.3 にあるように、基本的には質的変数、すなわち名義尺度水準や順序尺度水準のデータが得られたときに、それを解釈するためにはどのような数値を割り振ってやれば良いか、という発想から生まれたものになっています。

さきほどの表 14.1 のようなデータの並べ替えについては、数量化でいうところの III 類に該当します。ただこの数量化 III 類はおもしろいもので、同時期に独立にこの手法がフランス、カナダでも発展しており、2 つの別名を持っています。フランス学派が開発した手法は**対応分析 (correspondence analysis)**といい、カナダ在住の日本人、西里静彦^{*4}が開発した手法は**双対尺度法 (Dual Scaling)**といいます。どの分析も、基本的にはカテゴリカルな変数について直線性を最大にする値を割り振ることを目的にします。

カテゴリカル変数なので、分析の応用範囲は多岐にわたります。どのような変数でも名義尺度水準に落とすことができるからです。表 14.2 には一般的なデータセットの形を示しました。ID があって、Q1, Q2... と項目が列方向に並ぶ形です。一行が一人の反応を表しています。Q1 がたとえば性別などの名義尺度水準の変数であっても、男性 → 1, 女性 → 2 のようにコード化するルールを決めて入力します。Q2 はたとえばリッカート法で当てはまる、やや当てはまる...などの 5 件法だったとしましょう。その場合はシグマ法によるスコアを入れる、あるいは簡便的に当てはまるを 5, やや当てはまるを 4, といったように数字を割り振って入力しますね。

これをカテゴリカル変数と見なしてデータセットにした例が表 14.2 です。ID は分析対象ではありませんから横に置くとして、Q1 が 2 列に増えています。ID=1 の人は男性なのですが、この人は Q1:Male=1 かつ Q1:Female=0 というようにコード化されています。同様に、Q2 には「どちらとも言えない」を選択しているのですが、これを 3 とするのではなく 00100 と 5 列にわたってコード化しているのです。このようにすることで、

^{*2} 林 知己夫 (はやし ちきお, 1918 年 6 月 7 日 - 2002 年 8 月 6 日) は、日本の統計学者。正四位勲二等、理学博士。統計数理研究所第 7 代所長。社会調査・世論調査におけるサンプリング方法の確立を始め、数量化理論 (Hayashi's Quantification Methods) の開発とその応用で知られる。1990 年代以降、データの科学 (Data Science) を提唱し、その研究・思想は現在へと引き継がれている。Wikipedia より。

^{*3} 鮑戸 弘 (あくと ひろし, 1935 年 3 月 14 日 -) は、日本の社会学者、東京大学名誉教授。専門は、社会心理学、コミュニケーション論。Wikipedia より。

^{*4} 西里 静彦 (にしさと しずひこ, 1935 年 6 月 9 日 -) は、カナダの行動計量学者 (計量心理学)。学位は Ph.D. (ノースカロライナ大学・1966 年)。トロント大学名誉教授、アメリカ統計学会フェロー、日本行動計量学会名誉会員。北海道十勝郡浦幌町出身 (札幌市生まれ)。Wikipedia より。

2977 すべての変数をカテゴリとして扱うことができますようになります。

表 14.4 一般的なデータセット

ID	Q1	Q2	...
1	1	3	...
2	2	4	...
⋮	⋮	⋮	

表 14.5 カテゴリカル化したデータセット

ID	Q1:M	Q1:F	Q2:5	Q2:4	Q2:3	...
1	1	0	0	0	1	...
2	0	1	0	1	0	...
⋮	⋮	⋮				

2978 また表 14.1 を並べ替えた時のように、こうした名義尺度のデータの線形性が最大になるように、行だけで
2979 なく、列も並べ替えます。データの中で線形性が最大になるように、反応カテゴリと回答者に数字を割り振る
2980 のです。ちなみに表 14.1 は集計されたデータであり、今回の表 14.2 は集計前のデータです。カテゴリカルな
2981 変数の分析の場合は、どちらでも良いのです。行と列の関係が表されている数字であれば、「ある/ない」のよ
2982 うな二値反応でも、集計された度数でも構いません。さらにいえば、行と列の関係が記述されていればなん
2983 もいいのです。

2984 これまでの回帰分析、因子分析から SEM に至るまでの流れは、変数間関係だけを考えてきました。N 行
2985 M 列のデータ行列の計算の途中で、行に関する情報は平均化して潰されてしまい、最終的には $M \times M$ サ
2986 イズの正方行列だけを扱うことになったのでした。そして M 個の変数に**因子負荷量**などの重みをつけて考察
2987 してきました。今回のカテゴリカルなデータは、 $N \times M$ 行の矩形行列をそのまま分析し、行と列の両方に重
2988 みをつけます。正方行列でも矩形行列でも、変数と変数、変数と回答者がどのように関係しているかという
2989 ところを見るという意味では同じです。SEM では分散共分散行列や相関行列が、双対尺度法ではクロス集計
2990 表や素データがその分析対象になります。数学的な面から説明すると、正方行列の場合は**固有値分解**をし
2991 て、**固有値と固有ベクトル**を求め、固有ベクトルが変数の重みになるのでした。矩形行列の場合は**特異値分
2992 解 (Singular value decomposition)**と呼ばれ、行および列に対応する**特異ベクトル**をその重み、座標
2993 と考えることになります。本質的には同じようなものだと思っていただければと思います。

2994 14.3 双対尺度法による分析

2995 それでは実際の分析例を見てみましょう。次のコード code: 14.1 は、MASS パッケージに含まれるサンプ
2996 ルデータ caith を対応分析で分析するものです。

code : 14.1 双対尺度法による分析

```
2997 1 library(MASS)
2998 2 caith
2999 3 result <- corresp(caith,nf=min(nrow(caith),ncol(caith)-1))
3000 4 result
3001 5 plot(result)
3002
3003
```

3004 ■コード解説

- 3005 1 行目 パッケージ MASS のよみこみ
- 3006 2 行目 サンプルデータ caith を表示させる
- 3007 3 行目 対応分析関数 corresp で分析した結果を result オブジェクトに入れる
- 3008 4 行目 結果の出力

3009 5 行目 結果のプロット

3010 データ caith はスコットランドの Caithness 地方に住んでいた人に対してなされた調査で、髪の毛の色と
3011 目の色の関係を集計したものです。列方向に髪の毛の色、行方向に目の色が入っていますね (表 14.6)。

表 14.6 データ caith の中身

	fair	red	medium	dark	black
blue	326	38	241	110	3
light	688	116	584	188	4
medium	343	84	909	412	26
dark	98	48	403	681	85

3012 この行・列のカテゴリにはとくに順序性などありませんが、分析することで最も線形性の高いウェイトをつけ
3013 ることができます。もちろん 1 次元で表現できない可能性があり、その場合は第二、第三と次元数を増やして
3014 行きます。データの行数あるいは列数の小さい方マイナス 1 次元まで求めることができます。それを表現して
3015 いるのが、関数 `corresp` のオプション `nf` で、行数 `nrow(caith)`、列数 `ncol(caith)` の小さい方 (`min`
3016 関数) マイナス 1 を指定しています。

3017 結果は出力 14.1 のようになります。行および列にスコアがついていますね。この値を尺度値として使うと、
3018 データの相関係数が最大になるというような指標化がなされたのです。

R の出力 14.1: 対応分析の結果

```
First canonical correlation(s): 4.463684e-01 1.734554e-01 2.931691e-02 1.134031e-16
```

Row scores:

```

      [,1]      [,2]      [,3] [,4]
blue -0.89679252 0.9536227 2.1884132 1
light -0.98731818 0.5100045 -1.0837859 1
medium 0.07530627 -1.4124778 0.1894089 1
dark 1.57434710 0.7720361 -0.1482208 1

```

Column scores:

```

      [,1]      [,2]      [,3]      [,4]
fair -1.21871379 1.0022432 0.4271282 -0.8692696
red -0.52257500 0.2783364 -4.0268545 -1.3400421
medium -0.09414671 -1.2009094 0.1103959 -0.8453208
dark 1.31888486 0.5992920 0.3450676 -1.2251588
black 2.45176017 1.6513565 -1.5736976 1.1609621

```

3019
3020 結果はプロットされたものを見た方がわかりやすいかもしれません。図 14.1 にプロットを示しました。たと
3021 えば横軸、dim1 にそって目の色を見ていくと、light - blue - medium - dark の順に並べた方が良く、と
3022 いうことがわかります。とくに light と blue は近いのであまり大きな意味的違いはないことがわかります。同
3023 様に髪の毛の色は、fair - red - medium - dark - black の順に並べられることになります。具体的な数字は
3024 それぞれ出力の一行目の通りです。この dim1 だけでは表現できない違いが、dim2 で表されており、それ
3025 が図では上下の広がりとして示されていますね。

3026 この並びが、新しく構成された次元だと考えることができます。そしてこれを見ると、たとえば髪の毛と目の

3027 色がどちらも dark である場合, この二点は近くにプロットされています。この二点が近くにあるということは,
 3028 両者の結びつきが強いということです。髪の毛が暗い人は目の色も暗い人が多い, といった関係の強さが見
 3029 て取れます。

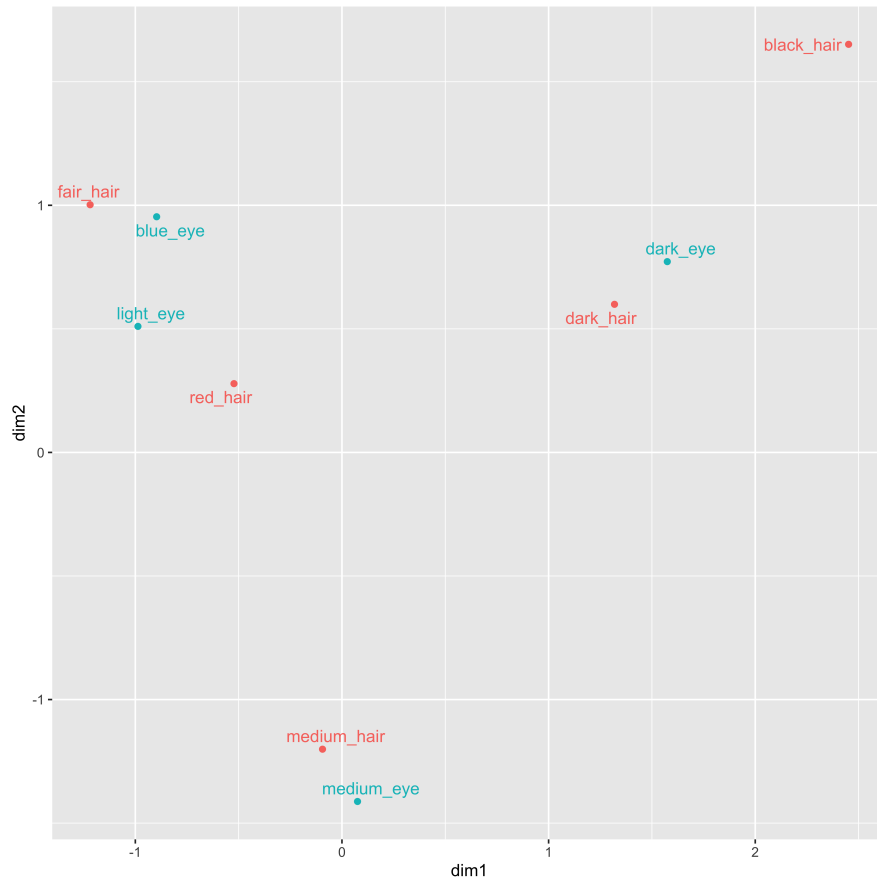


図 14.1 分析結果のプロット

3030 ちなみにこれは対応分析独特の表示法で, 行と列の変数が 1 つの画面にうつされていますね。厳密には,
 3031 行ベクトルが作る次元, 列ベクトルが作る次元は別物です。ですから, 列の要素を行の空間に写像するた
 3032 めには変換が必要で, 逆もまた然り, なのです。双対尺度法の場合はこの違いを厳格に捉えます。対応分析の
 3033 場合は, 相互に写像しあった座標を一枚の図にすることで情報を圧縮しています。数学的には同じモデルで
 3034 あっても, 表現の仕方や表示の仕方に, モデル作成者の考え方が反映されているとも言えるでしょう。

3035 14.4 テキストマイニングへの応用

3036 さて, 本講で紹介した分析方法は, 名義尺度水準を対象にしている, かつ外的な基準がなく内的な構造を
 3037 明らかにしようとするものだということがわかってきました。いわば名義尺度水準における因子分析ですね。
 3038 名義尺度水準を対象にしたということは, およそ言語化できたものはすべて分析の対象にできる, とも言えま
 3039 す。最も低次元で一対一対応した言葉の世界の数字だからです。

3040 応用例として, **テキストマイニング (Text Mining)** を紹介しておきましょう。別名自然言語処理とも言

3041 われませんが、これはテキスト、すなわち普通の「言葉」に潜む関係を分析する手法です*5。文章を対象にします
3042 から、小説や新聞記事などはもちろん、日記や逐語録なども分析の対象にしてしまおうというものです。心理
3043 学の分野では面接法などにおいて、クライアントがどのように語るかをすべて記録することがありますが、これ
3044 をみて「うーん、こういうことを考えているんだな」と読み取るのは主観的な判断によるものです。ここにテキス
3045 トマイニングを用いれば、機械的にどう言った言葉の使われ方をしているかを分析できます。あるいはツイッ
3046 ターやフェイスブックなどの SNS でどのような言葉、記事が流行しているかと言ったことを分析する、というよ
3047 うな使い方もできます。分析対象が一気に広がりますね。

3048 このテキストマイニングがやっていることは二段回に分かれ、第一段階が**形態素解析 (morphological**
3049 **analysis)**、第二段階が多変量解析です。第一段階は、「今日はいい天気ですね」といった平文を「きょう」
3050 「は」「いい」「てんき」「です」「ね」といった要素に分解することを指します。英語のような分かち書きがされる
3051 言語であればこれは簡単なのですが、日本語の場合は分かち書きされていないことに加え、漢字、ひらがな、
3052 カタカナなど表記方法もさまざまですから、この分析をするための特別な解析エンジンが必要です。幸い、す
3053 すでに開発されているものがフリーソフトウェアとして利用できます。時間がかかりますがこれらを使うと、品詞
3054 ごとに分解し、その原形 (活用する前の形) や活用形は何か、と言ったことを一覧にしてくれます。

3055 自然言語ですので大量の分割がなされますが、それを言葉同士の関係を表す行列の形で表現します。使
3056 われている品詞の原形ごとに誰が何回発言したかとか、各要素が 1 回の文章の中で同日に使われた回数を
3057 カウントするなどして、言葉と人、言葉と言葉の関係をデータ化するのがです。データ化できればあとはこっちの
3058 もです。数量化できるのですから、その言葉群のなかでどの言葉とどの言葉が近いのか、他のどの変数と
3059 関係するのか、と言った分析をできます。

3060 テキストマイニングについては R でもできますし、専門的なソフトウェアがあります。詳しくは樋口 (2020)
3061 を参照してください。

3062 14.5 課題

3063 今日の授業でおこなったすべての次の計算をする R コードを提出してください。ファイル形式は R スクリ
3064 プトか Rmd とします。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱い
3065 になります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を
3066 受けてください。

*5 マイニングとは鉱脈を掘る、という意味です。

第 15 章

多次元尺度構成法

今回は**多次元尺度構成法 (Multi-Dimensional Sacling; MDS)**について解説します。ここまで分散共分散行列を分解するところから、クロス表のような**名義尺度水準 (の)** データを分解するところまでやってきました。

多次元尺度構成法は、分散共分散行列を扱う線形モデルよりはやや仮定が緩く、また数量化のように解析者が数字を与えることを目的にするというよりは、その名の通り尺度を作ろうとしているという意味で心理学的・心理測定的なモデルだといえるでしょう。

私たちが単位もない不確かなものを対象にしながらもそれを測定するモノサシをつくることができるのは、2つの対象にたいしてその比較をして一方が他方よりも大である、ということが言えるからでしょう。記号を使って言えば、 x と y を比べて $x \succeq y$ である、ということから、尺度上の値 $p \geq q$ を対応させるということが、尺度を作るということです*1。このとき比較する2つが重さや広さのような物理的なものであればわかりやすいですし、「痛み」や「喜び」といった心理的な要素であっても構いません。あるいは「より賛成」といった態度表明のようなものでも良いかもしれません。この比較から出てくる関係は、分散共分散や相関係数のように強い線形の過程を置かなくても、より緩やかに**距離 (distance)** という考え方で表現されます*2。

MDS は距離行列を固有値分解するモデルです。それではこの方法について内容を見ていきましょう。

15.1 多次元尺度構成法

R にはサンプルデータとして eurodist というヨーロッパ各地の都市間距離のデータが用意されています。表 15.1 にその一部を示します。

元のデータは 21×21 のサイズです。表から、たとえばアテネ (Athens) とバルセロナ (Barcelona) の距離が 3313km、ということが読み取れます。表の右上が空白になっていますが、これは**距離行列が対称行列**なので、ムダな情報を表示しないようにしているからです。アテネとバルセロナの距離は、バルセロナとアテネの距離に等しいですからね。

さて距離行列はこのように、**正方行列**ですから、**固有値分解**をすることで**基底**を求めることができます。基底は座標を形成する基本単位ですから、それを使って各変数の位置にあたる座標を計算できます*3。このように距離行列から座標を求めて、変数をプロットすることで変数間関係を可視化する手法のことを**多次元尺**

*1 経済学では学問の最初の段階で、財を源とした集合に対する二項関係 $x \circ y$ を定義し、それを効用関数 u で実数領域に写像して $u(x) \geq u(y)$ を考える、といった基本的な原理を抑えます。心理学ではなぜか、「集合」「元」「二項関係」という基本的な比較の要素について考えることなく、その分析ツールだけがどんどんと発展してきています。

*2 後述しますが、相関係数も距離の一種と考えられます。

*3 厳密には距離行列 D そのものではなく、それを二重中心化した行列の固有値分解になりますが、本質的には変わりません。詳しくは岡太・今泉 (1994) などを参照のこと。

表 15.1 ヨーロッパ都市間距離データの一部

	Athens	Barcelona	Brussels	Calais	Cherbourg	Cologne
Athens	0					
Barcelona	3313	0				
Brussels	2963	1318	0			
Calais	3175	1326	204	0		
Cherbourg	3339	1294	583	460	0	
Cologne	2762	1498	206	409	785	0

3093 **度構成法 (Multi-Dimensional Scaling: MDS)** と言います。

3094 試しに `eurodist` データを MDS で 2 次元プロットしてみましょう。これを実行するコードは簡単で、
3095 `eurodist` データのようにデフォルトで組み込まれている `cmdscale` 関数を使います。

code : 15.1 計量 MDS の実践と描画

```

3096 1 library(ggplot2)
3097 2 # MDS を実行
3098 3 result.MDS1 <- cmdscale(eurodist, k=3)
3099 4 # y 軸反転させつつ描画
3100 5 g <- result.MDS1 %>%
3101 6   as.data.frame() %>%
3102 7   dplyr::mutate(label = rownames(.)) %>%
3103 8   ggplot(aes(x = V1, y = V2, label = label)) +
3104 9     geom_point() +
3105 10    geom_text_repel() +
3106 11    xlim(-2500, 2500) +
3107 12    ylim(2500, -2500) +
3108 13    xlab("dim_1") +
3109 14    ylab("dim2")
3110
3111

```

3112 ■コード解説

3113 1 行目 `ggplot2` でラベルをプロットするときに、綺麗な配置にしてくれる `ggrepel` パッケージを使います。
3114 このコードを実行するときに持っていない人は、インストールしておいてください。

3115 3 行目 `cmdscale` 関数に `eurodist` データを与えています。`k=3` は 3 次元解を出すように指定していま
3116 す。地球は球体ですから、地球上の地理は 3 次元で表現できますよね。

3117 5-14 行目 `ggplot2` による描画です。結果オブジェクトである `result.MDS` をデータフレームにし、行の名
3118 前になっていた都市名を変数として格納したのち、散布図のようにプロットしています。

3119 図 15.1 をみると、上の方にストックホルム (Stockholm) があって、右下にアテネが、中央にパリ (Paris)
3120 が・・・といった配置になっています。地球上の位置と完全に一致しているとは言えませんが、それでも概ねう
3121 まくプロットできていますね*4。このように、距離関係だけから地図を作ることができるというのが、MDS と
3122 いう手法なのです。

*4 完全に一致しない理由は、地球が球面であるのに対し平面にプロットしたから、と言うのもあります。

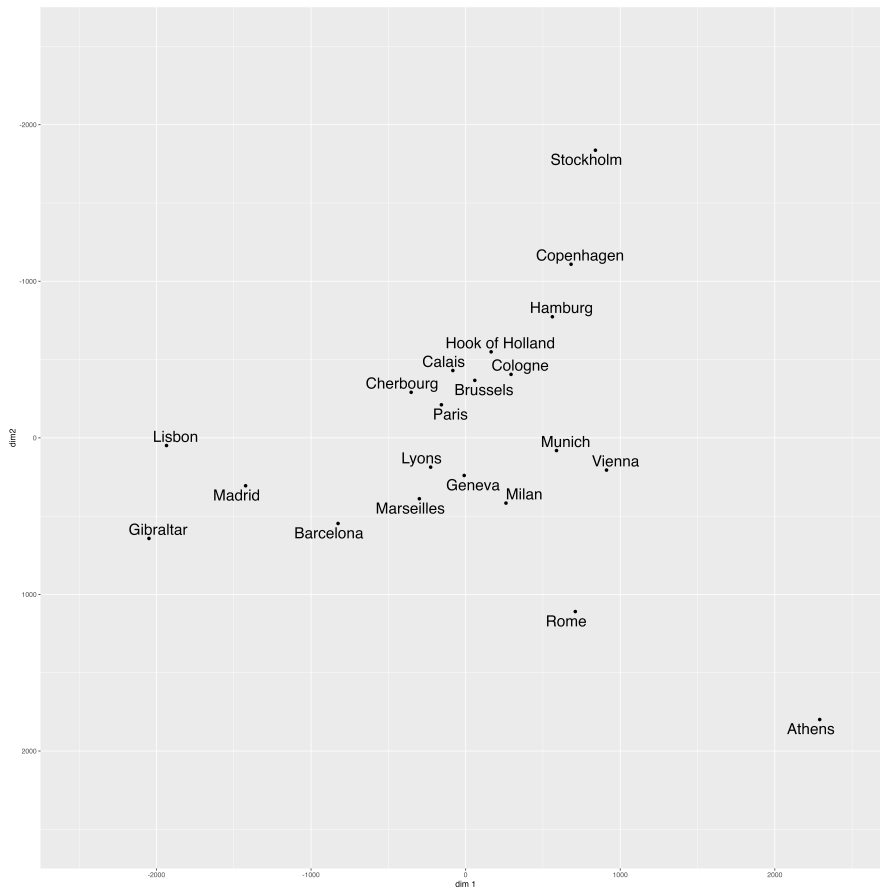


図 15.1 MDS のプロット

15.2 距離と心理学のデータ

MDS は距離関係だけから地図を作る方法でした。因子分析は相関関係から次元を作る方法でした。この 2 つの手法はとても似ています。というのも、相関関係と言うのが変数と変数の近さ・遠さ、つまり距離を表しているとも言えるからです。あらためて距離 (distance) とは何かを考えてみましょう。2 点 x, y の距離を $d(x, y)$ とすると、距離とよばれる数字の条件は次のようになります。

非負性 $d(x, y) \geq 0$

非退化性 $d(x, y) = 0 \Leftrightarrow x = y$

対称性 $d(x, y) = d(y, x)$

三角不等式 $d(x, z) + d(z, y) \geq d(x, y)$

もっとも一般的に使われるのはユークリッド距離で、2 次元座標 $(x_1, y_1), (x_2, y_2)$ があつた時のユークリッド距離は次のように計算します。

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

3 次元座標 $(x_1, y_1, z_1), (x_2, y_2, z_2)$ の場合は項を増やすだけです。

$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2}$$

3135 このようにして計算される距離ですが、上の条件を考えると何も二乗してルートを取らなくても絶対値を足し
3136 合わせるような方法でも構いません。たとえば 2 次元座標の場合は、次のような計算でもいいのです。

$$|x_1 - x_2| + |y_1 - y_2|$$

3137 現にこのような距離のことをマンハッタン距離 (Manhattan distance) といいます。二乗ではなく n 乗し
3138 て n 乗根をとる、という形で一般化することもできます。

$$d(x, y) = \left(\sum_{i=1}^n |x_i - x_j|^p \right)^{\frac{1}{p}}$$

3139 これはミンコフスキー距離 (Minkowski distance) という名前がついています。他にもマキシム距離、
3140 バイナリ距離、チェビシェフの距離、キャンベラ距離などがあり、いずれも R の dist 関数のオプションで選ぶ
3141 ことができますので、ヘルプなどを参照してみてください。

3142 また、相関係数 r_{ij} は $-1 \leq r_{ij} \leq +1$ の範囲にあります。この絶対値から $1 - |r_{ij}|$ とするとこれも距
3143 離の条件に当てはまります。どの程度類似しているかというも、距離と考えることができます。

3144 ここまでは数学的な距離のバリエーションでしたが、次に心理学的な意味に目を向けてみましょう。距離と
3145 は類似度でもありますから、何を距離と見なすかによって、さまざまな心理学的刺激がデータを形作ることに
3146 なります。

3147 たとえば評定尺度で、n 個の項目で何らかの評価をしてもらったとします。 x_{ij} を i 番目の項目における

3148 対象 j の評定値だとすると、 $d(j, k) = \sqrt{\sum_{i=1}^N (x_{ij} - x_{ik})^2}$ とすれば対象の類似度が計算できます。ほかに

3149 も何らかの刺激 A, B について、A か B かの判断をさせた時に混同してしまった混同率や、単語と単語の連
3150 想価、刺激の汎化勾配、反応潜時、ソシオメトリックなデータなど、いろいろなものが「類似しているかどうか」
3151 の指標として使えます*5。尺度評定よりも、具体的な 2 つの刺激が似ているかどうかの反応の方がやりやす
3152 い、というのは誰も実感としてわかることではないかと思います。

3153 類似度のデータが得られれば、それを距離と見做して MDS にかければ、変数間の関係を地図に描くこと
3154 ができるわけです。このように応用可能な領域が非常に広いことも、MDS の利点であると言えるでしょう。

3155 15.3 非計量多次元尺度法

3156 さて心理学的なさまざまな刺激が、MDS によって可視化できるということがわかってきました。距離行列
3157 が構成できれば、後の分析は何とでもできるわけです。ただし、この場合の距離行列とは、間隔尺度水準以上
3158 であることが必要です。しかし心理学的な刺激に対する反応をデータ化する時は、同時に被験者はそこまで
3159 鋭敏に反応しているのか、いいかえれば本当に間隔尺度水準以上の精度で判断できているのかな、という人
3160 間側の問題が気になりますね。そこまで人間は鋭敏無反応をしていないかもしれません。

3161 でも大丈夫。MDS は仮定を緩めたモデルがあります。非計量的多次元尺度構成法 (Non-Metric
3162 Multi-Dimensional Scaling) と呼ばれる手法がそれです。非計量 MDS では、対象 j と k の類似度
3163 を δ_{jk} とし、分析によって埋め込む多次元空間での距離を d_{jk} とすると、

$$\delta_{jk} > \delta_{lm} \text{ ならば } d_{jk} \leq d_{lm}$$

3164 のように、イコールではなく順序関係だけ保持して座標を求めます。この手法では、元データが順序尺度水
3165 準程度の情報しか持っていないなくても地図を描くことができるのです。人間の判断はせいぜいが順序尺度ぐら

*5 これらデータの例に関しては高根 (1980) の Pp.14-27 を参照してください。

3166 いですから、「より類似している ($\delta_{jk} > \delta_{lm}$)」のであれば「より近くにある ($d_{jk} \leq d_{lm}$)」というぐらいの配置
3167 の方が良いかもしれません。

3168 表 15.1 に示したような物理的距離であれば、間隔尺度水準の情報であることは間違いありません。その
3169 場合には、最初に示した固有値分解による手法を使います。これは非計量 MDS に対して、計量的多次元尺
3170 度構成法 (Metric Multi-Dimensional Scaling) と呼ばれています。

3171 15.3.1 非計量多次元尺度法の例

3172 心理学ではデータの性質上、非計量 MDS のほうが便利なが多いでしょう。計量 MDS でも非計量
3173 MDS でも、R では簡単な関数で実行できます。ここでは MASS パッケージに含まれる isoMDS 関数を使っ
3174 て実践してみます。

3175 使うデータは M1score2021.csv とします*6。ファイル名からお察しいただけるように、M-1 グランプリの
3176 評定をデータ化したものです*7。2021 年度のスコアは表 15.2 でした。

表 15.2 M-1 グランプリ 2021 の採点結果

演者	巨人	富澤	塙	志らく	礼二	松本	上沼
モグライダー	91	93	92	89	90	89	93
ランジャタイ	87	91	90	96	89	87	88
ゆにばーす	89	92	91	91	93	88	94
ハライチ	88	90	89	90	89	92	98
真空ジェシカ	90	89	92	94	94	90	89
オズワルド	94	95	95	96	96	96	93
ロングコートダディ	89	90	93	95	95	91	96
錦鯉	92	94	94	90	96	94	95
インディアンズ	92	91	93	94	94	93	98
もも	91	90	91	96	95	92	90

3177 M-1 の採点は審査員各自の主観に基づいて行われ、得点の絶対値はそれほど重要ではないかもしれませ
3178 ん。すなわち、松本人志の 80 点が上沼恵美子の 80 点と同じぐらいの面白さを評価しているか、ということ
3179 については真偽判断ができないでしょう。それでも各審査員の中での相対的評価には、一貫性がありそう
3180 す。すなわちある審査員が漫才師 A に 80 点、漫才師 B に 85 点をつけたのなら B の方が面白かったとい
3181 うことでしょうし、漫才師 C が 83 点なら $A < C < B$ という順序はあると思われます。つまり順序尺度水準程
3182 度の質はあると仮定することに無理はなさそうです。

3183 そこでこのデータをもとに、10 組の漫才師の類似度を計算します。類似度はユークリッド距離を用いるこ
3184 とにします。ユークリッド距離は既に述べたように差分の二乗を総和して平方根を取ったものですが、具体的
3185 な数字で見た方がわかりやすいかもしれませんので、表 15.3 を用意しました。表 15.3 にモグライダーとラン
3186 ジャタイの二組だけ取り出し、これで計算例を見てみます。それぞれの得点の差分、その二乗を計算し、それ

*6 ファイルはシラバスのサイトからダウンロードできます。

*7 念のために解説しておきますが、M-1 グランプリとは 2001 年から始まった漫才の賞レースの 1 つで、年末に年間チャンピオンが決定します。開催年ごとにルールが少し変わることもありますが、基本的には予選を勝ち抜いた 10 組の漫才師が 4 分間のネタを披露し、6-7 名の審査員が 100 点満点で採点します。点数の上位 3 組が決勝戦を行い 2 本目のネタを披露、投票によりチャンピオンが選出されるという流れです。2021 年はオール巨人、富澤たけし、塙宣之、立川志らく、中川礼二、松本人志、上沼恵美子が審査員で、最終的には錦鯉がチャンピオンになりました。このデータセットは 1 本目のネタについての採点を 2001 年から集めたものになります。

表 15.3 距離の計算

	巨人	富澤	塙	志らく	礼二	松本	上沼	総和
モグラライダー	91	93	92	89	90	89	93	637
ランジャタイ	87	91	90	96	89	87	88	628
差分	4	2	2	-7	1	2	5	9
差分の二乗	16	4	4	49	1	4	25	103

3187 を総和したところ 103 という値になっています。これの平方根を取ったもの、すなわち $\sqrt{103} = 10.14889$ が
 3188 この二組の距離、すなわち非類似度ということになります。この計算を全ての組み合わせについて計算してく
 3189 れるのが、`dist` 関数なのです。

code : 15.2 距離行列の計算

```

3190
3191 1 dat <- read_csv("M1score2021.csv")
3192 2 dat.mat <- dat %>%
3193 3   dplyr::filter(年代 == 21) %>%
3194 4   arrange(ネタ順) %>%
3195 5   dplyr::select(-年代, -ネタ順) %>%
3196 6   pivot_longer(-演者) %>%
3197 7   na.omit() %>%
3198 8   pivot_wider(id_cols = 演者,
3199 9     names_from = name,
3200 10    values_from = value) %>%
3201 11   as.matrix()
3202 12 rownames(dat.mat) <- dat.mat[, 1]
3203 13 dat.mat <- dat.mat[, -1] %>%
3204 14   dist()
3205

```

3206 ■コード解説

3207 1 行目 データファイルを読み込み、`dat` オブジェクトに格納します。

3208 2-9 行目 必要なデータだけに絞り込む操作です。流れを解説しますが、他のやり方でもいいですしできあ
 3209 ったものが何かだけわかれば結構です。

3210 3 行目 2021 年のデータだけに絞り込みます。

3211 4 行目 ネタ順に並び替えています。

3212 5 行目 年代変数とネタ順変数はもういらないので削除してしまっています。

3213 6 行目 ロング型に変換しています。これで演者-審査員-採点のデータセットができます。

3214 7 行目 欠損値を除外しています。実はこれがこの操作の目的で、というのも過去の審査データも
 3215 入っているものですから、過去の審査員も大量に変数として含まれていて、それらが欠損値に
 3216 なってしまっていたのです。

3217 8-10 行目 元のワイド型に戻しています。

3218 11 行目 以下の行列処理のため、`data.frame` 型から `matrix` 型に変換しています。

3219 12 行目 変数として一列目に演者名が入っていますが、これを行列の行名に入れています。`matrix` 型は
 3220 行名・列名をデータの外に持つのです。

3221 13-14 行目 変数としての演者名を除いて、パイプで `dist` 関数に入れ、距離行列を作っています。

3222 できた距離行列は、表 15.4 のようになっています。モグライダーとランジャタイの距離が、先ほどの例で計算
 3223 した値と一致していることを確認してください。またこれは**正方対称行列**ですから下三角だけ表示しておいま
 3224 す。また、対角が 0 になっています。自分自身との距離はゼロだからです。

表 15.4 演者の非類似度行列 (演者名は略記)

	モグ	ラン	ゆに	ハラ	真空	オズ	ロン	錦鯉	インデ	もも
モグ	0.000									
ラン	10.149	0.000								
ゆに	4.583	9.165	0.000							
ハラ	7.937	12.806	7.616	0.000						
真空	8.660	7.483	7.071	11.832	0.000					
オズ	12.490	15.652	12.207	14.933	11.000	0.000				
ロン	9.381	11.446	6.403	9.110	7.416	9.487	0.000			
錦鯉	8.485	15.264	8.307	10.909	10.247	7.071	7.874	0.000		
インデ	9.381	14.107	8.062	8.660	9.950	8.124	4.472	6.325	0.000	
もも	10.100	9.110	8.307	12.207	3.606	8.718	6.782	9.592	8.718	0.000

3225 あとはこの距離行列を isoMDS 関数に渡すだけです。isoMDS 関数は引数として、何次元の解を求めるか
 3226 を設定できます。地理データであれば 2,3 次元から作られていることは明らかですが、この評価が何次元か
 3227 は事前にわかりません。何次元にするかの指標として、Kruskal (1964b) は Stress と呼ばれる値を次のよ
 3228 うに定義しました。

$$Stress = \sqrt{\frac{\sum (d_{ij} - \delta_{ij}^2)^2}{\sum d_{ij}^2}}$$

3229 つまり実際のデータの距離 d_{ij} と、MDS で作られる空間上の座標から計算される距離 δ_{ij} の全体的な距
 3230 離を当て嵌まりの指標と考えていることになります。これを目安に、次元数を 1 から 7 まで変化させながら、
 3231 Stress 値がどうなるかを表したのが図 15.2 になります*8。

3232 まるで因子分析のスクリープロットの様ですね。横軸に次元数、縦軸に Stress 値を置いた折れ線グラ
 3233 フですが、当然のことながら反映させる MDS 空間の次元数が増えるとデータとの距離が縮んでいきます。
 3234 Kruskal (1964a) の基準によれば、Stress 値は表 15.5 のように評価できます。この基準でいくと、今回は 2
 3235 次元で 4.9%(0.0049534) ですから、2 次元解で OK としましょう。

3236 演者をプロットしたのが図 15.3 です。東西南北というか、上下左右の軸に特に意味はありませんから、因
 3237 子分析のように因子軸の解釈や命名をすることはありません。近くにプロットされた対象は評価が似ていた
 3238 んだな、ということがわかりますし、相対的な位置関係から、「ハライチとももは芸風が真逆だな」とか「ロング
 3239 コートダディは中心に近いから中庸的な笑い、悪く言えばキャラが立ってないんだな」といったことが読み取れ
 3240 ます。この図 15.1 や図 15.3 のように MDS で作られた座標のことを特に**布置 (configure)** といいます。
 3241 特に非計量 MDS は優劣・大小関係という順序尺度水準の評定だけからでも、その空間的な特徴を描くこと
 3242 ができますから、心理尺度のもつ仮定や測定モデルを考える必要がないため、今後その重要性が再発見され
 3243 ていくのではないかと思います。

*8 どうして 7 までか、というと評定者が 7 名だからです。7 名それぞれの評価次元があると考え、この距離空間の中には最大 7 次元あるはずですから。あるいは 10 組の演者がいますから、組み合わせの自由度から考えて 9 次元まで試してもいいと思います。

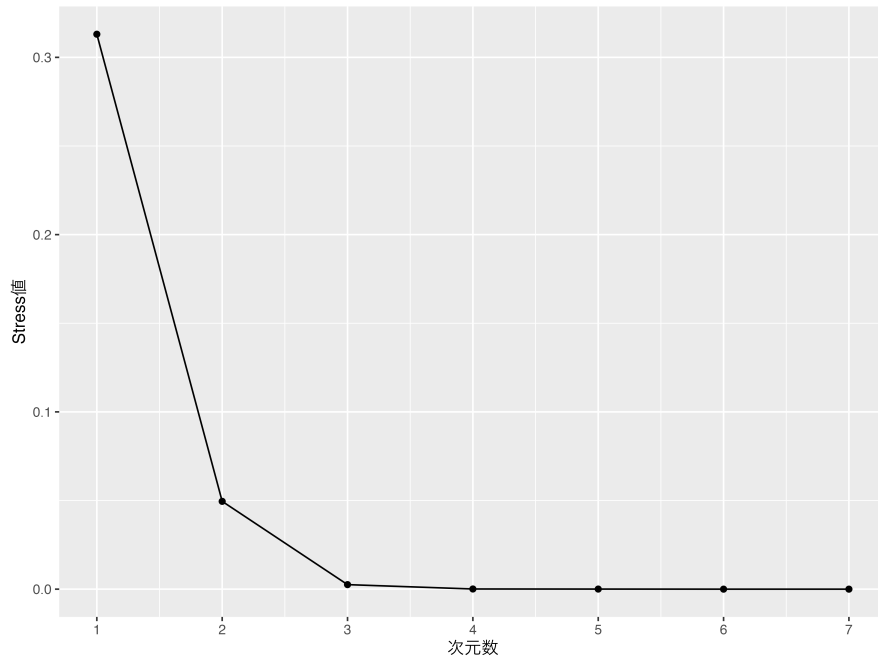


図 15.2 Stress 値の減衰

表 15.5 Stress 値の評価

Stress	Goodness of Fit
20%	poor
10%	fair
5%	good
$2\frac{1}{2}\%$	excellent
0%	“perfect”

3244 15.4 多次元尺度法の展開

3245 多次元尺度構成法で作られた地図は、対象をプロットした地図です。地図には、その上に何か書き込んだり、
3246 地形の図に天気図を重ねるように複数の地図を重ねて表現したりできます。多次元尺度構成法にも、この
3247 ような応用モデルがいろいろ考えられています。ここでいくつかの発展的な MDS モデルを見てみましょう。

3248 ■prefmap 類似度空間の上に、個々人の理想点を追加する方法です。個人の好み preference をマッピング
3249 する方法なので、**Preference Mapping()** と呼ばれています。個人 i が対象 j について、好みの程度
3250 を s_{ij} と評価したとします。対象 $1, 2, \dots, j, \dots, M$ は別途類似度評定によって、MDS のつくった地図上にプ
3251 ロットされているとします。ここで個人 i と対象 j との地図上の距離 d_{ij}^2 に対して、次の回帰式を考えます。

$$s_{ij} = a_i d_{ij}^2 + b_i + e_{ij}$$

3252 つまり、個人 i と対象 j の距離を使って、好みの程度 s_{ij} を予測するモデルを作り、誤差がもっとも小さくな
3253 る点に個人 i をプロットするのです。こうすることで、対象とそれを評価する人を一枚の地図に表現すると言

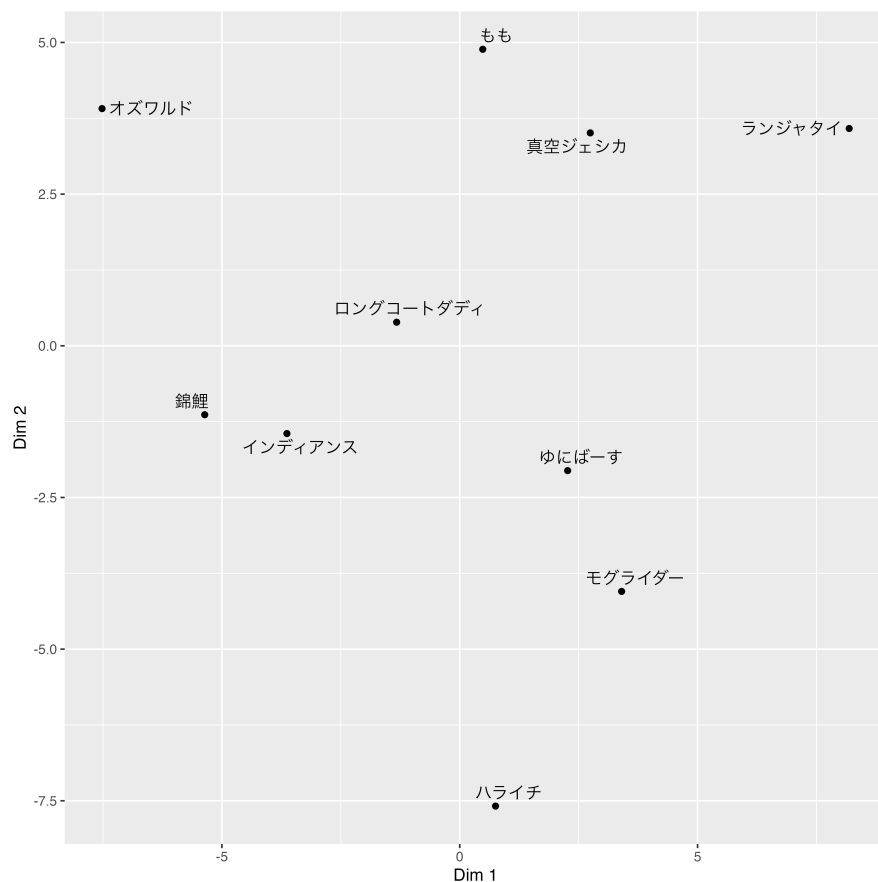


図 15.3 非計量 MDS のプロット

3254 うことができます。

例えば先ほどの M-1 の例ですが、私の採点では表 15.6 のようになりました*9。この評定値を使って著者

表 15.6 著者の評定

演者	採点
モグライダー	90
ランジャタイ	60
ゆにぼーす	85
ハライチ	92
真空ジェシカ	83
オズワルド	89
ロングコートダディ	85
錦鯉	83
インディアンズ	82
もも	88

*9 ランジャタイは何が面白いのかわからなかった。モグライダーはもっと評価されるべき。ハライチも良かったですね、ちょっと時間オーバーしたっばいけど。

3255

3256 の理想点を書き足したのが図 15.4 です。低く評価したランジャタイからは遠く、高く評価したハライチやオズ
 3257 ワルドに近いところに著者の理想的な笑いの点があり、錦鯉やインディアンも近くにありますが、今回の優
 勝にはまあ納得、と言ったことがわかります。皆さんも自分の理想の点を書き加えてみませんか？

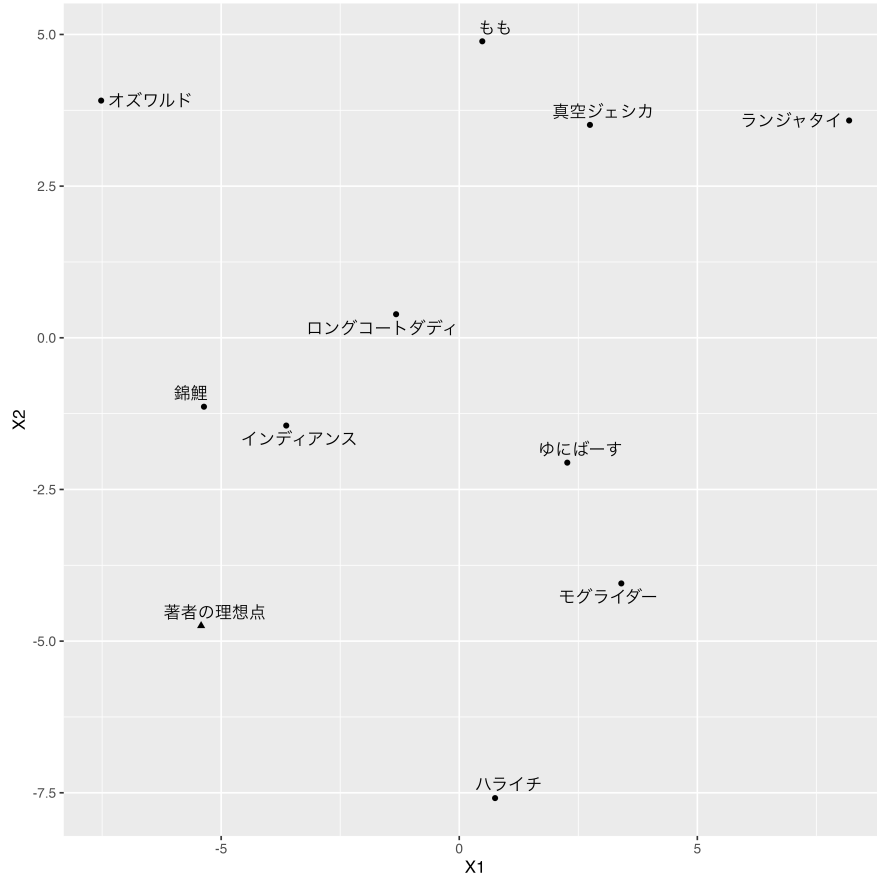


図 15.4 著者の理想点プロット

3258

3259 ■Abelson Map これは (Abelson, 1954) の考えた手法で、これも prefmap と同じく布置される対象に別
 3260 の力 (選好度でもなんでもよい) があると考え、地図空間上に力の場をプロットして等高線を引くことで表現
 3261 するモデルです。

3262 この方法では、各点 P にかかる力 $V(P)$ を次のように定式化します。

$$V(p) = \sum_{j=1}^M \frac{V(j)}{1 + d_{pj}^2}$$

3263 ここで j は各点、ここで言えば漫才師のことで、 $V(j)$ が漫才師 j に与えられた評価点です。 d_{pj}^2 は任意の
 3264 点 p と対象 j の距離ですから、任意の点 p にかかる力は各漫才師が発する笑い力ですね。

3265 この方法で、先ほどの M-1 プロットに、著者の評価を Abelson Map の方法で描き加えてみましょう。図
 3266 15.5 の右にある、ランジャタイの周りにたくさんの線が引かれていますが、いわばこれは低気圧のようにここ
 3267 のパワーが弱い (と著者は思っている) ところになります。実は、真空ジェシカの左側に通っているラインが平
 3268 均点のラインですから、ここを境に著者は「面白い」と「面白くない」を区分しているとも言えます。ハライチや

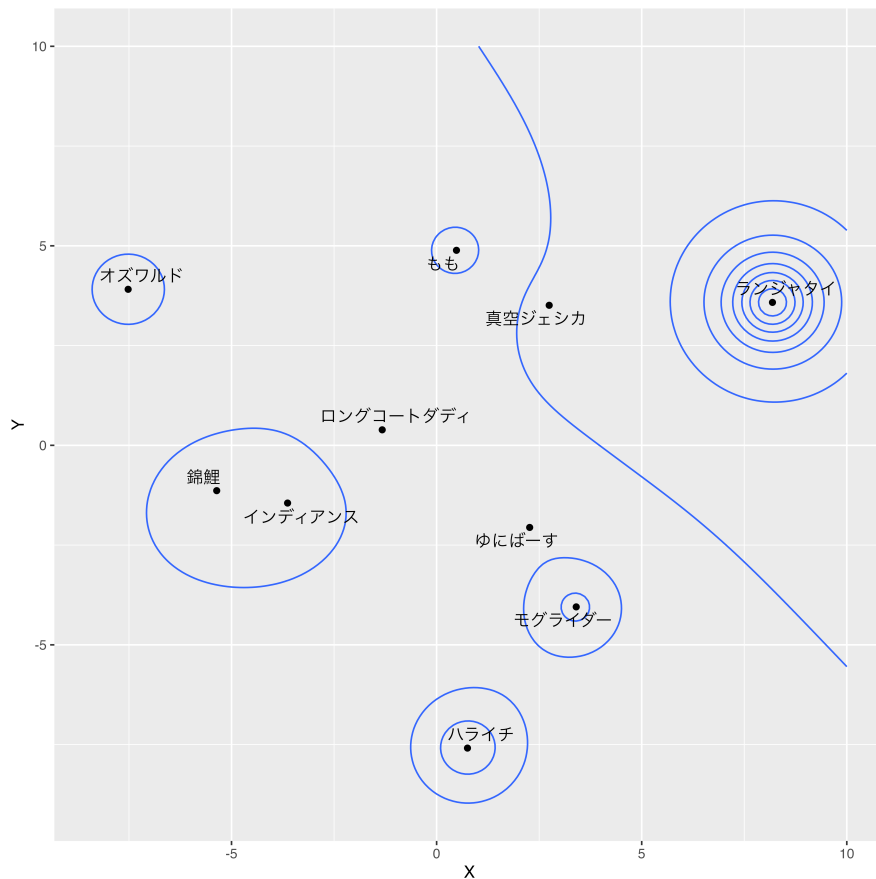


図 15.5 Abelson Mapping によるお笑い力の等高線図

3269 オズワルドの周りにある等高線は山の高さ、パワーの強さを表しているラインですね。

3270 力のメタファーにすぎないと言われたらそうですが、このように色々なものを可視化できるのも MDS の面白ところですよ。

3272 ■非対称多次元尺度構成法 距離データは対称でなければならない、というものの、たとえば好きな人に振り向いてもらえない＝片想いとか、都市間の人口の流入・流出、国際貿易の赤字・黒字などを考えると、対称間の関係が対称でないことは少なくないわけです。そうした非対称な関係を、図に加えようと言うのが非対称 MDS です。非対称情報を対象部分 + 歪対象部分に分割し、この対称でない要素を図の中に書き加えるモデルが一般的で、プロットする対象の周りに縁や楕円で表現するとか、矢印で表現するとか、vonMises 分布と呼ばれる確率分布で表現するとか、さきほどの Abelson Map で高さとして表現するなどのモデルがあります。より数学的なモデルとして、プロットする空間を複素空間にするモデルもあります。詳しくは千野・岡田・佐部利 (2012) を参考にしてください。

3280 ■多次元展開法 リッカート法は心理尺度でもっともよく使われる方法で、「非常に当てはまる」から「まったく当てはまらない」までの段階的な反応を被験者に求める方法です。これは段階反応モデルで分析されることから明らかなように、項目に対する反応段階が順番的に強くなっていくことを仮定しています。しかし、実際に被験者が丸をつけるときは、「自分の感覚にもっとも近いカテゴリーを選ぶ」ということをしているはずですよ。つまり、被験者と反応カテゴリーの距離が、尺度に反映されているはずなのです。この考え方から、

3285 尺度に対する反応を距離とし、項目とそれに回答した被験者の両方をプロットする地図を書く方法がありま
3286 す。Coombs が考えた**多次元展開法 (Multi-dimensional unfolding method)**と呼ばれる手法が
3287 それで、この手法から態度測定の尺度化を考えることもできます。発想としては数量化理論に近く、また清水
3288 (2018) はこれを確率モデルに展開しています。

3289 ここであげたいいくつかの例にあるように、多次元尺度法もさまざまな角度から人の判断や反応から、**尺度**
3290 (**scale**) を作る方法です。多次元尺度構成法は因子分析モデルよりも制約や仮定が少なく、より直接的に心
3291 理的反応をモデル化しようとしています。

3292 本講で学んだように、心理学ではさまざまな角度から「数値化するルール」を考えてきていますが、それぞ
3293 れ仮定や目的、何を良い尺度と考えるかという観点がことなります。我々心理学者は、統計的なツールのユー
3294 ザに過ぎません。データは分析できればそれでいい、分析は機械がやってくれるから深く考えなくていい、と
3295 いう割り切り方もあるかもしれませんが、「そもそも何をデータとするか」「どのように数値化するか」という点
3296 については、そのような態度は取れないはずで、なぜならそもそも「心はいかにして数値化できるか」という
3297 ことが明らかでないならば、その後の分析はすべて嘘っぽい数字を使った統計ごっこに過ぎないからです。**心**
3298 **理学者であるために、われわれは何をどのように数値化しているかについて、しっかりと理解する必要**
3299 **があるのです。**

3300 15.5 課題

3301 計量 MDS, 非計量 MDS をそれぞれ一例ずつ、実践してみてください。データはどのようなものでも構い
3302 ません。提出ファイル形式は R スクリプトか Rmd とします。なお提出されたコード単体でバグがなく動くこと
3303 が確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別 TA
3304 か小杉までメールで連絡し、指導を受けてください。

第II部

3305

心理学データ解析応用 2

3306

第 16 章

プログラミングの基礎

さて、ここから始まる統計の話は、統計学の中でも最近使われるようになってきたベイズ統計を駆使し、色々な心理学的現象を分析していこうというものになります。

このアプローチはデータに合わせて分析モデルを考えていくことになるので、ボタンをポチポチ押すと答えが出るような定型パターンの分析ではありません。個別のケースに合わせて分析モデルを考えていくことになります。その時に必要になってくるのがプログラミング技術ですので、まずはウォーミングアップということでプログラミングの話を始めたいと思います。中にはうんちくのようなところもありますので、興味があるところだけ飛ばして読んでいただいても構いません。最後にあげた課題ができるのであれば、知識は後からでもいいと思います。

16.1 プログラミングの基礎

プログラミングとは、コンピュータに計算をさせるその仕様書を作ること、だと思ってください。お料理のレシピのように、計算レシピを書くわけですね。ただ相手は計算機ですので、言葉の端々まで正確に伝えないと理解してくれません。よく「思った通りに動いてくれない！」と不満を訴える初心者がいますが、それは当然で、プログラムは**思った通りに動くのではなく、書いた通りに動く**のです。思った通りに動かないのであれば、それは仕様書・レシピの方に間違いがあります。

今の第一原則に加えて、プログラミングを進める上での注意点をあげておきます。簡単なことだと思うかもしれませんが、基本的なルールをしっかり守ることが上達への近道です。

1. プログラムは思った通りには動かない。書いた通りに動く。大文字、小文字、スペルの違いに注意する。
2. 書き間違えないための工夫は美しき。綺麗に書くことが大事。
3. 一言一句すべてに意味がある。ただの写経ではなく、意味を考えながら書く。
4. 「遊び心」をもって！ここを変えたらどうなるか、を少しずつ「やってみる」が大事
5. 変更は少しずつ。一気に変えようとどこが変わったかわからなくなるから。

16.1.1 綺麗に書こう

綺麗に書く、というのはコードの可読性を上げるために必要です。綺麗な書き方として、松浦 (2016) は次の 5 点をあげています。

- インデントは必ずする
- データを表す変数の先頭の文字は大文字。パラメータを表す変数の先頭の文字は小文字。

- 3335 • 各ブロックの間は1行あける
- 3336 • camelCase ではなく snake_case で。
- 3337 • 「」や「=」の前後は 1 スペースあける

3338 1 つめのインデントとは、字下げのことです。行の頭を少し凹ませることで、違う人まとまりであることを明示
3339 するのです。たとえば次の 2 つのコードは同じ働きをしますが、16.2 の方が見やすいと思いませんか？

code : 16.1 インデントのないコード

```
3340
3341 1 data {
3342 2   int<lower=0> N;
3343 3   array[N] y;}
3344 4 parameters {
3345 5   real mu;
3346 6   real<lower=0> sigma;}
3347 7 model {
3348 8   for(i in 1:N)
3349 9     y[i]~normal(mu, sigma);}
3350 10 }
3351
```

code : 16.2 インデントのあるコード

```
3352
3353 1 data {
3354 2   int<lower=0> N;
3355 3   array[N] y;
3356 4 }
3357 5
3358 6 parameters {
3359 7   real mu;
3360 8   real<lower=0> sigma;
3361 9 }
3362 10
3363 11 model {
3364 12   for( i in 1:N ){
3365 13     y[i] ~ normal(mu, sigma);
3366 14   }
3367 15 }
3368
```

3369 インデントのある 16.2 は、data というブロック ({ と } で囲まれている領域) の中に、2 つの行が入っている、ということが明確にわかります。また model というブロックの中には、for から始まる分がありますが、これも複数行に渡るブロックを構成する文なので、その中身はインデント (字下げ) されています。このように、
3370 インデントすることでどこの列がどこまでブロックを組んでいるのかがわかりやすくなるのです。ちなみにこの
3371 インデントを作るのは TAB キーを使います。TAB キーは普段、どう使うのかわからないものだと思われがちですが、このインデントをしてくれるためのものです*1。コードエディタによっては、カッコで括ったときに自動的に閉じるカッコを用意し、また改行ごとに字下げ位置を合わせてくれるものがあります。これらの機能を使って是非わかりやすいコードを書いてください。

3372
3373
3374
3375
3376 松浦 (2016) の指摘の 2,3 番目についてはお好みで、4 番目もそれほど強い・広く浸透した決まりではあ

*1 ここには流派があって、TAB でインデントする派とスペースキーを 4 回または 8 回押してインデントする派があります。どちらでも働きは同じなのですが、個人的には数回の字下げを 1 回のキーで行ってくれる TAB のほうが良いと思っています。

3378 りません*²。ただし、5 番目のスペースの前後に少し空白を取るの、見やすさのためにも是非実行してもら
3379 いたいところですよ。

3380 これらの書き方は、**可読性**を上げるためのものです。プログラムは仕様書・レシピですから、あとで読み直
3381 した時やほかの人が読むときも、意味がわかるようにしておいた方が良いでしょう。自分だけわかれば良い、と
3382 思うかもしれませんが、その自分ですら何があるのかわからなくなってしまう可能性があります。そのような
3383 読みにくさはすぐにバグにつながりますから、綺麗に書くことを常々心がけておいて欲しいのです。

3384 16.1.2 意味を考えて書こう

3385 プログラミング言語は、ほとんど英単語のようなものです。プログラミング上の都合から、短い言葉に省略
3386 されていたり、アンダースコアやピリオドでつながって書いてあったりしますが、比較的わかりやすい言葉であ
3387 ることに違いはありません。

3388 プログラムを習得し上達するためにすべきことは、最初は写経と呼ばれる、テキストや指示にしたがって書
3389 き写すことです。もちろんネット上にあるサンプルコードなどをコピー&ペーストしても良いのですが、自
3390 の分析コードを書くためには手を入れる必要があったりします。ですから、最初はなるべく自分でキーボード
3391 を叩いて、コピー&ペーストではなく入力する経験を積んでください。もちろん間違えてしまうことが多いで
3392 しょうが、転ばずに歩けるようになった人が一人もいないように*³、ミスをするのでどういうミスだったかを
3393 考え、自覚することではじめて、すこしずつですがミスが減っていきます。失敗する経験も時には必要な
3394 です。

3395 ということで、時折自分なりにコードを書いてもらうのですが、その時に「ただの記号列」と思わないでくだ
3396 さい。すでに書いたように、英語あるいはそれを省略した表現になっているので、見慣れないものかもしれま
3397 せんが必ず意味があります。プログラムを始める最初の頃は、ミススペルが多く、あちこちでエラーが出て気
3398 持ちが萎えることもあるかもしれません。エラーが出る時は目を凝らして、どこに問題があるかを考えて修正
3399 することになります*⁴。ここで難しいのが、 x や s などの大文字と小文字の区別がつきにくいもの、 1 と l のよ
3400 うに違いがわかりにくいものがあるということです。私は職業柄、多くの人に教え、多くのバグを見つけてきま
3401 したが、その中でも特別見つけるのに苦労したのが、`norma1` というミススペルでした。みなさん、これのどこ
3402 がミススペルかわかりますよね？でもこれ、書く時に「正規分布のことだな」と思っていれば、そもそも入力時
3403 にそんな入れ方はしないとと思うのです。ただの字だ、と心を無にしてしまうと、後々見つけにくいエラーを作る
3404 ことにもなりますから、この文字・この名前・この関数はどういう意味だろう、と少し考えながら取り組んでみて
3405 ください。

*² ちなみにキャメルケース Camel Case、スネークケース Snake Case とは変数名のつけ方のルールです。R などプログラミング言語ではあらゆるものに名前をつけて管理しますが、名前のつけ方は任意です。命名はわかりやすい方がいいですが、長いものになると区切りを入れたいくなるかもしれません。さて Camel とはラクダの意味で、ラクダのコブのように区切りのところだけ大文字にする、という記法です。Snake は蛇のことで、区切りのところにアンダースコア_を使うというものです。たとえば野球チーム、という変数名をつけたい時に、キャメルケースのやり方だと `BaseballTeam` のように書きますし、スネークケースのやり方ですと `Baseball_team` のような書き方になります。ちなみに R は日本語での命名も許しますが、そのほかの言語では一般的ではありませんし、何より全角と半角を切り替える時のミスやエラーが多くなるので、半角英数字での命名をすべきという点はどちらでも共通です。

*³ お釈迦様は除く。

*⁴ ちなみにプログラミング歴 30 年以上の私でも、簡単な英単語、たとえば `library` のスペルを間違えるなんてことはしょっちゅうです。

3406 16.1.3 遊び心が大事

3407 プログラミングを進める上では、遊び心が大事です。うまく動くコードがかけたら、もうおかしなことになりた
3408 くない、と手をつけずにまましてしまう人がいます。

3409 ところで、今まで教えてきた経験からいって、プログラミングが上達する人は、うまくコードが書けたらすぐ
3410 に「ここをこう変えたらどうなるんだろう」と試してみる人だと思います。最初はプログラムの意味がまったくわ
3411 からないので、どこをどういじっていいのか見当もつかないかもしれません。しかしたとえば blue という文字
3412 が出てきたら、これは色の青のことじゃないか、と思いますよね。ではその blue を red に変えたらどうなるん
3413 だろう？やってみよう！となるような、そういう遊び心が大事なのです。

3414 もちろん変えてしまったことでエラーになって、動かなくなってしまうことがあるかもしれません。その場合は、
3415 元に戻して元通り動くかどうか、さらに確認すれば良いのです。数字や文字を少しいじっただけで、PC が爆
3416 発してしまうようなことはありませんから、ちょっと遊び心を出して、どうなるのかな？と思う気持ちを大事にし
3417 てください。

3418 16.1.4 変える時は少しずつ

3419 プログラムを色々いじって遊ぶことが成長につながる、という話をしました。ただし遊び方には気をつけま
3420 しょう。まずうまくいくコードができれば、それは保存しておいて、また同じ内容のものを別名で保存し、「遊ん
3421 でいい方」「壊れてもいい方」を作って、そちらで色々変えて遊ぶといいでしょう。最悪なことがあっても、うま
3422 くいくバージョンは常に残っているわけですから。

3423 そして色々試す時のコツは、一箇所ずつ変更していくことです。たとえば blue を red に変え、line を
3424 box に変え・・・と複数箇所を同時に変更して、うまくいかなかった場合、どちらが原因になっているのか把握
3425 できませんよね。これは心理学実験でいうところの交絡です。操作が交絡してしまって、原因が特定できなく
3426 なるのです。

3427 また、実行するときも 1 行ずつやりましょう。とくに R は一問一答型、つまり 1 つの命令について 1 つの反
3428 応を随時返してきます。複数行をまとめて実行することもできますが、まとめて実行している途中でエラーが
3429 出ていて、その後の計算がすべて空回りしているということも少なくありません。エラーがどこで生じているの
3430 か、しっかり特定することが重要です。そのためにも焦らず、1 つずつ確実にできることを積み重ねていくとい
3431 う姿勢が重要です。

3432 統計分析も複雑で細部まで設定し尽くしたモデルを作り上げていくことができますが、いきなり全体が完成
3433 するのではなく、小さなピースの積み重ねなのです*5。

3434 16.2 プログラミング言語

3435 それでは具体的に、プログラミング言語としての R の機能をいろいろ見ていきましょう。

*5 R やプログラミングで統計を学ぶ意義はここにあります。GUI で操作できる統計パッケージも少なくありませんが、それらは画面が出てきた時にデフォルトで設定されている値があるのがほとんどで、自分が何をやっているのか細部まで気づかないまま実行できてしまうのです。そうすると当然、エラーが出たり、うまくいっても誤用している、ということにもなりかねません。自分で理解できていない技術に振り回されないようにするために、しっかりとわかるピースを積み重ねていく必要があるのです。

16.3 プログラミング言語の基本的な働き

プログラミング言語はさまざまありますが、それらのやっていることは基本的に計算であり、その表現方法が違うだけです。どの言語にも共通する、プログラミング言語の特徴的な働きは「代入」と「反復」と「条件分岐」です。以下、この3つの働きについて R を使ってみていきます。

16.3.1 代入

料理でよくあるシーンとして、「野菜を切ってザルにあげておく」など、いったん作業の途中経過を横に退けて別の作業をする、というのがあります。計算プロセスでも、「計算結果をいったん横に置いておく」という操作をしたいことがあります。これが**代入**で、コンピュータとしてはあるメモリ番地に値を書き込んでおく、という操作をすることになります。横に置いておく時に名前をつけることができ、これが R のなかでは**オブジェクト (Object)** と呼ばれるものになります。すでにみなさんもやったことがあると思いますが、計算結果をオブジェクトに代入しておくコードの例をみてみましょう (code: 16.3)。

code : 16.3 代入操作

```

1 a <- 1
2 b <- 2
3 a + b
4 a <- 3
5 a + b

```

■コード解説

1 行目 オブジェクト *a* に 1 を代入

2 行目 オブジェクト *b* に 2 を代入

3 行目 オブジェクト *a* と *b* に + という操作。すなわち $1 + 2$ の計算を指示していることと同じ。

4 行目 オブジェクト *a* に 3 を代入。もともと入っていた情報は上書きされる。

5 行目 オブジェクト *a* と *b* に + という操作。すなわち $3 + 2$ の計算を指示していることと同じ。

非常に単純な例で、今更何を、と思った人もいるかもしれませんが、これが基本中の基本なので改めて示しました。R での代入は `<-` あるいは `=` を使います。前者は矢印のイメージ、後者は `a=3` のようにして、「*a* は 3 だよ」という意味で代入を表しています。ポイントは 3 行目の「オブジェクト名で値を指定していること」と、4 行目の「値の上書き」にあります。オブジェクト名で計算を表現することで、操作の一般化ができます。すなわち、3 行目と 5 行目が同じ式であるのに結果が違うように、「どんなに値が違って同じ動作をする」ようにできるわけです。コンピュータがすごいのは、この手の代入がどれほど複雑に、どれほど繰り返されても決して計算ミスせず、指示通りに動くことができるということです。人間は「ベクトルの 2 番目の要素と 3 番目の要素を足す」という操作を数回やるだけでも計算ミスをするのがあり得ますが、コンピュータは同じ計算を 100 回やっても 1000 回やってもミスなく動作します！

ただし注意が必要なのは、長いプログラムを書いていて、途中で「さっきの分析をやり直そう」と思って部分的に実行し、さらに続けて分析をする、ということをやっている時に上書きが生じることがあることです。途中の部分的なやり直しでできた計算は、改めて冒頭からやり直すとうまくいかないことになったりします。上書きをする時は注意して、また完成版として清書する時は一度プログラムの冒頭から走らせてみて、処理が間違いないか進んでいるかを確認するようにしましょう。

3474 また、ここでの代入例は 1 つの数字だけでしたが、数字がセットになったベクトルや行列もオブジェクトに代
 3475 入できます。セットではなく、別々の要素だけどもまとめて持っておきたいという時は list 型を使います。さら
 3476 に一般的なデータセットは行に個体が、列に変数がある矩形行列ですが、R では list 型の特殊ケースであ
 3477 る data.frame 型と呼ばれます。このように、形式による違いをデータの型といいます。

code : 16.4 さまざまなデータの型

```

3478
3479 1 # 数字をセットで持つvector型, matrix型
3480 2 x <- c(1, 2, 5, 8, 9)
3481 3 A <- matrix(c(1, 2, 3, 4, 5, 6), ncol = 2, nrow = 3)
3482 4 # 型にこだわりなくなんでも収納list型
3483 5 dataSet <- list(name = c("kosugi", "koji"),
3484 6                 v1 = c(1,2,4,5,3,5,6),
3485 7                 v2 = matrix(c(1,2,3,4),ncol=2))
3486 8 # 矩形に整えられたlist型であるdata.frame
3487 9 df <- data.frame(list(name=c("kosugi", "suzuki"),
3488 10                    V1 = c(176,173),
3489 11                    V2 = c(83,55)))
3490

```

R の出力 16.1: データの型による持ち方の違い (コード 16.4 の作ったオブジェクト)

```

> x
[1] 1 2 5 8 9
> A
  [,1] [,2]
[1,]  1   4
[2,]  2   5
[3,]  3   6
> dataSet
$name
[1] "kosugi" "koji"

$v1
[1] 1 2 4 5 3 5 6

$v2
  [,1] [,2]
[1,]  1   3
[2,]  2   4
> df
  name  V1 V2
1 kosugi 176 83
2 suzuki 173 55

```

3491

3492 16.3.2 反復

3493 エビフライをたくさん作ることを考えてみましょう。「えびの殻を剥いて、衣をつけてあげる」という動作を何
 3494 回も繰り返すことになると思います。このような反復捜査もコンピュータは得意とするところ。なにせ、疲

3495 れを知りませんので。反復操作を指示するコードは for 文という書式を使います。反復計算コードの例は次
3496 のようなものです (code:16.5)。

code : 16.5 反復操作 1

```
3497 1 a <- 0
3498 2 for (i in 1:4) {
3499 3   a <- a + 1
3500 4   print(a)
3501 5 }
3502
3503
```

3504 ■コード解説

3505 1行目 オブジェクト a に 0 を代入 (初期化)
3506 2行目 for 文開始。中括弧で括られている 5 行目までの中身を反復計算する。
3507 3行目 オブジェクト a に、もとの a に 1 を加えた値を代入 (上書き) している。
3508 4行目 オブジェクト a を出力させている。
3509 5行目 for 文ここまで。

3510 ここでは 3 行目の「1 を加える」、4 行目の「表示する」という操作を繰り返していることとなります。ポイン
3511 トは 2 行目の書き方ですね。for の後ろの小カッコ () の中身が、繰り返しに使う要素の指定です。後ろの
3512 1 : 4 は R 言語特有の書き方で、「1 から 4 まで」すなわち 1, 2, 3, 4 を意味しています。変数 i がこの順番で
3513 変わるよ、ということの意味しているので、 i は 1 回目, 2 回目, 3 回目, 4 回目, と進んでいくことになりま
3514 す。code:16.6 のような書き方をすると、その振る舞いがよくわかるかと思います。

code : 16.6 反復操作 2

```
3515 1 a <- 0
3516 2 for (i in c(1, 3, 5, 15, 12)) {
3517 3   print(paste(a, "に", i, "を加えます"))
3518 4   a <- a + i
3519 5   print(a)
3520 6 }
3521
3522
```

3523 4 行目には $a <- a + i$ という代入があります。「 a に $a + i$ を代入せよ」という意味であり、数学のイコー
3524 ル記号だと意味がわかりませんが ($2 = 2 + 1$ なんて式はおかしいですもんね), ここでは $a + i$ の計算をし
3525 たものを、新たにオブジェクト a に上書きせよという意味になります。またここでのポイントは、反復
3526 用インデックス i が 1, 3, 5, 15, 12 の順に変わっていくということと、変わっていくインデックス
3527 i の値そのものも計算に使えるということです。2 点目については注意が必要で、反復用インデッ
3528 クスが 1, 2, 3, 4 と変わるような例の場合、途中で `verbi <- 1` などとしてしまうと、いつまでもインデッ
3529 クスが前に進まず永遠に計算が終わらないこととなります。そのようなミスが生じないように、注意してくださ
3530 いね。

3531 また、この for 文は入れ子にして使うことができます。たとえば i が 1 から 5 まで変わり、 j が 1 から 3 ま
3532 で変わりながら計算をする、ということを考えて、次のような表現が可能です。

code : 16.7 反復操作 3

```
3533 1 A <- matrix(1:15, ncol = 3, nrow = 5)
3534 2 for (i in 1:5) {
3535 3   for (j in 1:3) {
```

```

3537 4     print(paste("Aの",i,"行",j,"列目の要素は",A[i,j]))
3538 5     A[i, j] <- A[i, j] + (i * j)
3539 6   }
3540 7 }
3541 8 print(A)
3542

```

3543 ■コード解説

- 3544 1行目 オブジェクト *A* は 3 行 5 列の行列で、1 から 15 までの数字が順に入っている。
- 3545 2行目 for 文開始。*i* が 1 から 5 まで変わる。
- 3546 3行目 for 文その 2 開始。*j* が 1 から 3 まで変わる。
- 3547 4行目 オブジェクトの中身を表示させる。
- 3548 5行目 行列 *A* の *i* 行 *j* 列目の要素に対し、 $i \times j$ の計算結果を加えたものを代入 (上書き) させる。
- 3549 6行目 for 文その 2 がここまで。
- 3550 7行目 for 文その 1 がここまで。
- 3551 8行目 最終計算結果の表示。

3552 行列の行、列それぞれについて順番に、各要素を指定しながら代入していくという計算です。挙動を確認し
3553 ておきましょう。

3554 16.3.3 条件分岐

3555 「卵が 200 円より安かったら買ってくる」というようなお使い指示、ありますよね。これは「200 円以上の時
3556 は買わない」という意味でもあります。この 200 円かどうか、という条件に対して、その後の動きが「買う」「買
3557 わない」に分岐するので、条件分岐と呼ばれます。これをプログラムで書くと次のようになります。

code : 16.8 条件分岐

```

3558 1 egg <- 250
3559 2 if (egg < 200) {
3560 3   print("卵を買います")
3561 4 } else {
3562 5   print("卵を買いません")
3563 6 }
3564
3565

```

3566 ■コード解説

- 3567 1行目 オブジェクト *egg* に 250 を代入。
- 3568 2行目 if 文開始。() 内が条件で、条件が該当すれば次の中括弧で括られている領域までを実行する。
- 3569 3行目 「卵を買います」と画面表示させる。
- 3570 4行目 条件が該当した時の実行内容を閉じつつ、該当しない場合の実行内容を書く領域を展開する。
- 3571 5行目 「卵を買いません」と画面表示させる。
- 3572 6行目 条件が該当しない場合の実行内容を閉じる。

3573 一行目の *egg* にさまざまな数字を入れて検証してみてください。思った通りの振る舞いができているで
3574 しょうか。この if 文は小括弧の中身が条件節ということになります。条件が成立していることを真または

3575 **TRUE** と表現し、成立しないことを偽または **FALSE** と表現します*6。今回は egg 変数が 200 未満である
 3576 ことを条件としています。もし 200 円も含めたいのであれば (200 円以下にしたいのであれば), egg <=200
 3577 と書く必要があります。

3578 ところで、ちょうど 200 円のときに、という一致を表す条件はどう書けば良いのでしょうか？
 3579 if(egg==200){...}としたいところですが、このままでは括弧の中身が「egg オブジェクトに 200 を
 3580 代入せよ」という命令と同じになってしまいます。条件節で使う「同じかどうか判定する」の記号はとくに、==で
 3581 表します。「卵がちょうど 200 円であれば」という条件節は、if(egg==200){...}と書くのが正しい表記で
 3582 す。逆に「同じでないとき」を表現したい場合は、!=と書きます。他にも条件節は、「A かつ B のとき」、「A
 3583 または B のとき」のような表現をしたくなるがよくあります。このような条件節の表記や計算のことを論理演
 3584 算といい、「かつ」は&&、「または」は||で表現します。コード 16.9 に論理演算の例を示しました。これを実行し
 3585 て、R が TRUE か FALSE のどちらで返答してくるか、確認してみてください。

code : 16.9 論理演算

```
3586 1 X <- 2
3587 2 X > 3
3588 3 X < 3
3589 4 X == 3
3590 5 X != 3
3591 6 X > 1 && X < 3
3592 7 X < 1 || X > 3
3593
3594
```

3595 ■コード解説

3596 1 行目 オブジェクト X に 2 を代入しておきます。
 3597 2 行目 X は 3 より大きい?
 3598 3 行目 X は 3 より小さい?
 3599 4 行目 X は 3 と等しい?
 3600 5 行目 X は 3 と異なる?
 3601 6 行目 X は 1 より大きく、かつ、X は 3 より小さい?
 3602 6 行目 X は 1 より小さい、または、X は 3 より大きい?

3603 条件分岐をする場合は、条件が思った通りに設定できているか、分岐した後のルートが間違いなく書かれ
 3604 ているかなどに注意が必要です。というのも、我々が日常言語で使っている「もし~なら XXX する」といった
 3605 表現は曖昧なことがあり、書いた通りにしか解釈しないコンピュータは、思った通りに動いてくれないとい
 3606 ことがあるからです。Twitter 上で有名になったジョークに、次のようなものがあります*7。

3607 ある妻がプログラマの夫に「買い物にいて牛乳を 1 つ買ってきてちょうだい。卵があったら 6 つお
 3608 願い」と言った。夫はしばらくして、牛乳を 6 パック買ってきた。妻は聞いた「なんで牛乳を 6 パックも
 3609 買ってきたのよ！」夫いわく「だって、卵があったから……」

3610 これはプログラミング的思考の例として示されており、気の利かない夫とされていますが、言外の意味を含
 3611 めすぎた妻の条件分岐指示がまずかったといえます。より適切な指示の例は、次のようなものです*8。

*6 TRUE, FALSE が大文字なのは重要で、真偽を表す R の特別な用語です。こうした用語のことを予約語といいます。

*7 <https://twitter.com/beamtetrode350b/status/1406773935069229057>

*8 <https://twitter.com/tak1/status/1127065591380971521>

3612 プログラマを夫にもち長年苦勞してきた妻「1 パック 8 個以上入った 190 円以上 210 円以下の鶏卵
3613 のパックがあったら買ってきて。それがなかったら 1 パック 8 個以上入った 189 円を超えない最も高
3614 い鶏卵のパックを買ってきて。売り場について 15 分以内に見つけられなかったときは私に電話してか
3615 ら指示にしたがって」

3616 これでもまだ指示が厳密でない！条件漏れの可能性がある！というコメントがありますが、「書いた通りに
3617 しか動かない」ことの重要さがお分かりいただけるかと思います。

3618 16.4 まとめ

3619 ここで説明したのは、レゴのピースのようなものです。ピース 1 つ 1 つは小さくて、それだけではなんの形も
3620 作り上げることができませんが、細かなピースでも組み合わせ次第で大きなオブジェを作ることができます。
3621 プログラミングにはたった 1 つの正解というのではなく、どういう形であれいと通りに動くものができればそれで
3622 構いません。書き方は人それぞれでよく、結果が伴うかどうかポイントです。どのピースをどのように組み上
3623 げてもいいので、思い通りのものが作れるようにトレーニングをしておきましょう。

3624 16.5 課題

3625 ■FizzBuzz 課題 1 から 15 までの数列に対して、次の条件にあった出力をするプログラムを書きなさい。な
3626 お、割り算の余りを計算する R の関数は `%%`、文字を出力する関数は `print` です。

- 3627 • その数が 3 で割り切れる場合には、その数の代わりに「Fizz」を出力する
- 3628 • その数が 5 で割り切れる場合には、その数の代わりに「Buzz」を出力する
- 3629 • その数が 3 でも 5 でも割り切れる場合には、その数の代わりに「FizzBuzz」を出力する

3630 ■行列のかけ算 行列 $A = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 2 & 0 \end{pmatrix}$ と $B = \begin{pmatrix} 1 & 5 \\ 2 & 3 \\ 3 & 8 \end{pmatrix}$ を用意し、 AB の掛け算をしたいとします。

3631 R の行列の積を求める記号ではなく、`for` 文をつかって、この計算をするプログラムを書きなさい。

第 17 章

データ生成モデリング

17.1 データ生成モデリング

「心理学統計法」と名のつく講義のほとんどは、記述統計や推測統計、とくに帰無仮説検定を扱うことが中心であるのが現状です^{*1}。心理学において帰無仮説検定を行う理由は、心理学実験の効果を検証すること、それも手元の標本についての記述ではなく母集団においてもその効果があると言えるかどうかを検証するためです。手元のデータに基づいて、目に見えない母集団全体のことに推論するというのは、基本的に不良設定問題です。つまり、確実な正解を導き出すにはヒントが少なすぎる、圧倒的に不利な状況で知恵を絞る必要があるのです。こうした不利な状況下ですから、いくつかの仮定をおいて、確率的に推論するのでした。その推論の方法として、モーメント法 (moment method)、最尤法 (Maximum Likelihood estimation)、ベイズ法 (Bayesian estimation) があるのでした^{*2}。

これらの推定方法が実験計画 (Experimental Design) (要因計画ともいう) と組み合わせで用いられ、帰無仮説検定となるのでした。心理学における実験的なアプローチは、実験群と統制群に無作為に割り当てた集団に対し、介入・処置のあとでの群平均を見ることでその効果を検証します。人間を相手にする研究ですから、当然誤差や個人差がデータには含まれますが、無作為割り当てと平均化によってそれらはキャンセルアウトされ、平均の比較をすることで効果を見ることができると考えるのでした。また標本の平均ではなく、それを用いて母平均を推測することで、結果の一般化を考えます。母平均の推定には点推定と区間推定とがあり、確率的表現を用いた区間推定をつかって慎重に結論を導き出すのです。ただし区間推定は判断基準が明確ではありませんから、帰無仮説と対立仮説という 2 つのモデルを戦わせて決着を見る、というやり方をして一応の決着を見るのでした。こうした「実験計画」+「推測統計学」=「帰無仮説検定」は、長らく心理学のスタンダードとして君臨しています。

こうしたアプローチについて、昨今批判があることについては、今は触れないでおきましょう^{*3}。それよりもこうした問いの立て方や解決策に注目してみたいと思います。心理学は物理学を範として、科学 science の仲間入りをしたい、というモチベーションが強く根付いている世界であり、また人間というのは嘘をついたり間違えたりするものだ、ということが骨身に染みてわかっていますから^{*4}、データを分析する際にも客観性を大事にすることが殊更重視されます。客観性の反対は主観性、つまり本人の思い込みや考え方の癖がもっとも

^{*1} とくに公認心理士に必要な授業として、科目名「心理学統計法」を名乗る必要があり、そこで教える基本的な内容としてこれらが含まれています。なお測定論やそれに関する多変量解析法はそれほど中心的話題ではありません。みんな使うのになあ。

^{*2} これら 3 つの方法について、十分に思い出せない人は一年時の授業資料を振り返ってみてください。

^{*3} Amrhein, Greenland and McShane (2019) などが指摘するように、誤った使い方をされることが多く心理学研究の意義そのものが疑われるほどになっている現状があります。また豊田 (2020) ではさらに辛辣に批判がなされています。

^{*4} どうか心理学の研究というのは、いかに人間がダメで偏った考え方をする生き物であるかを、滔々と明らかにしていくという側面もあります。

3658 邪魔になるのです。そこでデータに対しては真っ白な気持ちで向き合うものだ、という姿勢をとることになり
3659 ます。

3660 言い方を変えると、分析をする前はデータについて何も考えてないよ、という態度をとるわけです。もちろ
3661 んデータは要因計画の結果として得られるものですから、計画を立てる際は主観的な誤謬が微塵も入り込む
3662 ことないように緻密に練り上げるのですが、出てきた結果は結果、数字の羅列に過ぎない考えるのです。そ
3663 の結果を統計的に分析するときは、それが記憶実験の数字であろうが臨床実験の結果であろうが気にするこ
3664 とはなく、淡々と帰無仮説検定の俎上に乗せていくことになります。こうしたアプローチができるからこそ、そし
3665 てそれを許す統計ソフトウェアがあるからこそ、心理学の研究を積み重ねていくことができたのだという一面
3666 があります。こうした研究アプローチは、データ駆動型分析と言えるでしょう。つまりデータができてから、分析
3667 が始まるという考え方です。

3668 これに対して、データがどのように生まれてきたのか、そのメカニズムを考え、その仕組みを明らかにしてい
3669 こうというのがデータ生成モデリングです。松浦 (2016) は統計モデリングを「確率モデルをデータに当ては
3670 めて現象の理解と予測を促す営みのこと」と定義しています。すなわちこのアプローチは、まずデータがどの
3671 ようなメカニズムで生まれてきたのかを、簡潔な数式を使って表現します。そのモデルをデータに当てはめ、こ
3672 のモデルからデータが出てきたと言えるかどうか、他のモデルの方が今のデータをうまく説明するのではない
3673 か、といった比較検証をしていくことになります。データ駆動型分析では、こうしたメカニズムが実験計画の中
3674 に暗黙理に埋め込まれていたと言えるかもしれません。データ生成モデル駆動型の分析は、メカニズムを明
3675 示して検証する、というところが違います。

3676 データが作られるメカニズムを考えるアプローチの利点は、予測にも向きます。メカニズムが正しい、ある
3677 いはうまく現状のデータを再現できるのであれば、おそらく今後も同じようにデータが作られていくでしょう。
3678 であれば将来の予測ができる、という考え方です。たとえば市場の動向の予測、消費者の傾向の予測ができ
3679 ればそのメリットは想像に難くありません。昨今はデータサイエンスという言葉が流行していますが、そこには
3680 こうした研究アプローチの隆盛があるのです。もちろん、心理学においてもモデリングのアプローチは有用で
3681 すし、従来通りの平均値の比較をする上でも、さまざまなメリットがあります。

3682 この授業では、データ生成モデリングのアプローチをとって、心理学的にもどういう意義があるのかを考え
3683 ていきましょう。

3684 これを考える上で重要なのが**確率モデル (stochastic model)** です。言葉の通り、データ生成過程を確
3685 率を使って表現します。なぜ確率を？と思うかもしれませんが、手元のデータは誤差や個人差を含んで微小
3686 に変わる、偶然の成分が必ず入っているからです。こうした偶然をハンドリングし、偶然の中にも理論的な予
3687 測をするという点で確率が必要になるのです。また確率モデルで考えるときに、最尤法よりもベイズ法の方が
3688 よく使われます。確率モデルが複雑になっていくにつれ、最尤法では計算コストが非常に高く、実質的に解が
3689 得られないということも少なくありません。ベイズ法によるアプローチはこの問題をクリアしてくれるからです。
3690 それでは改めて、ベイズ法について復習しておきましょう。

3691 17.2 ベイズ推定の基礎

3692 データ生成モデルが確率の言葉で記述される、というのはすでに述べたところです。推測統計というのがそ
3693 もそも不良設定問題、すなわち少なすぎるヒントから未知なるものを当てるという状況に置かれた学問であ
3694 り、この「わからないこと」を確率で表現するからです。この確率についての考え方として、ベイズの公式という
3695 のがあるのでした。ベイズの定理は次のように書き表されます。

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)}$$

3696 ここで $P(X)$ は X についての確率, という表現です。 $P(A|B)$ は条件付き確率というもので, B が与え
3697 られた時の A の確率, という意味です。

3698 ベイズの定理の用語は次の通りです。まず右辺の分子に注目しましょう。 $P(D|\theta)$ とありますが, ここで D
3699 はデータを表しています。 θ はデータを生む確率のパラメータです。 $P(D|\theta)$ はパラメータ θ のもとでのデー
3700 タ D が得られる程度を表すもの, という意味になります。これは**尤度 (likelihood)** と呼ばれるもので, パラ
3701 メータの関数になっているのでした。たとえば正規分布からデータが生成されると考えるのであれば, 手元の
3702 データがパラメータ θ から出てくる可能性がどれぐらいあるのか, を表現していると言えるわけです。この尤
3703 度関数は, データを生み出すメカニズムの表現そのものです。

3704 右辺の分子の第二項, $P(\theta)$ はパラメータの確率です。正規分布からデータが生成される例で言えば, 尤
3705 度はあるパラメータの値からデータが得られる程度を表現しているのですが, そのパラメータが出てくる確率
3706 はそもそもどの程度であるか, を考えているのです。これは別名**事前分布 (prior distribution)** と呼ば
3707 れます。実際にデータが出てくる前の段階で, そもそもそのパラメータがどれほど出やすいかという確率を表
3708 現しており, これは今回のデータを取る前までの事前のデータや経験に基づいている, ということができます。

3709 右辺の分母, $P(D)$ はデータ全体が得られる確率であり, **周辺尤度 (marginal likelihood)** とか**エビ
3710 デンス (evidence)** と呼ばれます。**正規化定数 (normalized constant)** ということもあります。これに
3711 ついては理解しにくいところもあるかもしれませんが, 後に述べる理由によって, ひとまず深く考える必要はあ
3712 りません。

3713 左辺はこの計算の結果得られる**事後分布 (posterior distribution)** を表しています。 $P(\theta|D)$ はデー
3714 タ D で条件づけられたパラメータ θ の確率です。我々は確率モデルを作るわけですが, そのパラメータがど
3715 うなっているかは事前にはわかりません。が, データを取ることで「データが与えられたのであれば, パラメー
3716 タはこうだよ」というのがわかるわけです。あるいは事前にパラメータはこのあたりにあるのではないかと経
3717 験上考えていたとしても, それがデータを取ることによって更新される (確信が強くなったり, 違うかもしれな
3718 いと思ひ直したり), ということを意味しています。

3719 この式を言葉で表現し直すと, 次のようになります。

$$\text{事後分布} = \frac{\text{尤度} \times \text{事前分布}}{\text{周辺尤度}} = \text{尤度} \times \frac{\text{事前分布}}{\text{周辺尤度}}$$

3720 **尤度**がメカニズムそのものだ, ということはすでに書きました。たとえば**回帰分析**においても, 尤度を計算
3721 できます。普通の回帰分析は, 誤差が確率的に生じると考え, 誤差以外のところは線形モデルで表現します。
3722 すなわち, 従属変数 Y_i に対して, 独立変数 X_i があつたとすると,

$$Y_i = \hat{Y}_i + e_i = \beta_0 + \beta_1 X_i + e_i$$

$$e_i \sim N(0, \sigma)$$

より,

$$Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma)$$

3723 となるのでした。ここで $N(\mu, \sigma)$ は平均 μ , 標準偏差 σ の正規分布を表し, データ Y_i が正規分布から生成
3724 されているというモデルになっています。回帰分析の場合は**最尤法**で未知数 β_0, β_1, σ を求めたりしました
3725 が*5, その名の通りこのモデルが示す尤度を最大にする未知数の求め方を最尤法というのでした。

3726 ですが, 尤度は確率を表すものではありません。確率とは非負の実数で, すべての確率を足し合わせる
3727 (あるいは積分する) と 1.0 にならなければなりません, 尤度は足し合わせても 1.0 にならない数字なので

*5 あるいは**最小二乗法 (Least Square method)** で求めるのですが, これは幸い最尤法の結果と一致します。正規分布を
使った線形回帰モデルの場合, データの記述統計的値と確率モデルによる推定値が一致するので, 気づかないまま推測統計学
に足を踏み入れてしまえるのでした

3728 す。ですから、「手元のデータがパラメータからでてくる**可能性**」という変な言い回しをしていました。「出てく
3729 **る確率**」とは言えないのですね。そこで、尤度を確率に置き換えたい、と考えたときに便利なのがベイズの定
3730 理なのでした。ベイズの定理は尤度に事前分布をかけ、周辺尤度で割るという操作によって、事後分布が得
3731 られる式、と読むこともできます。事後分布は**確率分布**です。尤度にある数字をかけて（事後の）確率分布に
3732 していると考えるといいでしょう。ちなみに事前分布も確率分布ですが、**周辺尤度**は確率ではありません。た
3733 とえば度数を総度数で割った**相対頻度**は確率と考えることができますが、それと同じように、分子をとある数
3734 字で割って、全体を 1.0 に整えるための定数なのです。**正規化定数**というのはそういう意味ですね。定数倍
3735 （正確には定数の逆数）をかけて大きさの調整をしているだけです。ベイズの式はさらに次のように書き
3736 直すことができます。

$$\text{事後分布} \propto \text{尤度} \times \text{事前分布}$$

3737 ここで \propto は比例する、という関係を表しています。最終的には事後分布の形が知りたいのですが、その形
3738 を決めているのは分子の尤度と事前分布だけだ、ということになります。

3739 ここで事前分布が**一様分布 (uniform distribution)**であれば、事後分布の形状には影響せず、事後
3740 分布は尤度そのものの形を反映することになります。従来のデータ駆動型分析では、データに対して事前の
3741 仮定を一切置かないということでしたが、ベイズ流の分析でもそのことは同様に表現できるわけです。

3742 以上がベイズの定理についての簡単な復習でした。しかしベイズの定理自体は、1740 年代には明らかにな
3743 っていたことであり、それがなぜ 21 世紀の今になって見直されてきたのでしょうか。それには色々な理由が
3744 ありますが、その 1 つは最近まで「ベイズ流の分析は絵に描いた餅」だったからです。すなわち理屈ではこの
3745 ような形になることは明らかだったのですが、実際に計算するのは難しいことが多々ありました。その問題を
3746 解決したのが**マルコフ連鎖モンテカルロ法 (Markov Chain Monte-Carlo method; MCMC)**と
3747 呼ばれる方法です。

3748 17.3 マルコフ連鎖モンテカルロ法

3749 ベイズの定理から、事前分布と尤度が分かれば、計算式を解いて事後分布の形を算出できます。しかし**確**
3750 **率分布**の式が複雑になればなるほど、その計算はとて難しくなります。確率関数を複数組み合わせたり、求
3751 めるべきパラメータの数が増えて行ったりすると、とてじゃないけど計算して答えを出すことができない、と
3752 いうことになります。そのせいもあって、ベイズ統計学は長らく実践には向かない手法だったのですが、**マル**
3753 **コフ連鎖モンテカルロ法 (Markov Chain Monte-Carlo method)**、略して MCMC が出てきてか
3754 らその状況が一変しました。

3755 MCMC 法は近似的な答えを探す 1 つの方法です。MCMC という名前は「マルコフ連鎖」と「モンテカル
3756 ロ法」の 2 つのパートからできあがっています。マルコフ連鎖は、目標となる確率分布状態になるような推移
3757 規則を作る方法であり、モンテカルロ法は乱数を発生させるアルゴリズムです。この 2 つが合体することで、
3758 確率分布の形がわからなくてもサンプルが得られるような、乱数発生器を作ることができるようになったの
3759 です。

3760 ベイズ統計のゴールは事後分布を作ることです。事前分布を一様分布にして、尤度はベルヌーイ分布にす
3761 る、といった簡単な場合ですととくに問題ないのですが、ある確率分布のパラメータがあって、さらにそれを生
3762 成する確率分布があるといった、階層的に入れ子になっているようなモデルが出てくると、「最終的な事後分
3763 布の形がわからない」という状況になります。正規分布とかポアソン分布といった、名前がついている確率分
3764 布であればその特徴がはっきりわかるのですが、それらを組み合わせで作られるものは名もなき合成関数に
3765 なり、どんな形状をしているのか、まったく想像つかないことがあります。マルコフ連鎖を使うと、そのなもなき

3766 合成関数の形状をとにかく作り上げることができる、そんな数学的技術です。

3767 モンテカルロ法は乱数発生技術です。乱数というのは規則性のない数ですが、これを作るのはなかなか難
3768 しいものです。人間が 0-9 の数字を適当に書いていくと、知らず知らずのうちに均等でないパターンができ
3769 てしまいます*6。コンピュータに (たとえば R に) 乱数があるじゃないか、と思われるかもしれませんが、コン
3770 ピュータはあくまでも計算機でどこまでも合理的です。乱数によって動いているように見えますが、乱数に見
3771 えるような数字を生成するアルゴリズムがその中には埋め込まれていますから、正確には擬似乱数でしかあ
3772 りません。コンピュータの作る乱数は、もちろん人間が適当に思いつく乱数よりもより規則性が少ないです
3773 が、それでも基本的に

- 3774 • 何らかの関数 g によって、内部状態 S_t をアップデートする; $S_{t+1} = g(S_t)$
- 3775 • 内部状態 S_t から何らかの関数 h によって、実現値が生成される; $x = h(S_t)$

3776 というステップを反復 ($t = 1, 2, 3 \dots$) することで生成するのです。こうした乱数列ができれば、それを正規
3777 分布の形だとか二項分布の形に当てはめて出力することは簡単なのです。このステップ・バイ・ステップの計
3778 算法としてある確率分布に従う乱数発生技術があり、これがマルコフ過程とひっついてできたのが MCMC
3779 法ということになります*7。

3780 MCMC 法は、ですから、どんな形の確率分布であっても乱数のサンプリングは生成できる、という技術
3781 なわけです。事後分布の関数の形、形状がわからなくてもそこからの乱数は作れるという技術であり、コン
3782 ピュータ技術の性能が発展した今では大量の乱数を生成することも瞬時に行われます。

3783 乱数を用いるアプローチはいくつかの利点があります。たとえば**確率分布**の平均である期待値など、分布
3784 の特徴を記述するための計算は積分を含むので解析的に解くのは知識も技術も必要です。しかしその確率
3785 分布から乱数を発生させると、その平均値を求めることでその近似値を得ることができます。確率分布の計
3786 算が記述統計の計算に置き換えられるのであり、また統計ソフトウェアにとって大量のデータの記述統計量を
3787 描くのは造作もないことなのです。乱数による近似値は、あくまでも近似値、近くて似ている値に過ぎませ
3788 ませんが、精度を上げるためにはその乱数の数を増やしてやるだけで良いこととなります。この点もいいですね。

3789 また求めたいパラメータが複数ある場合、たとえば正規分布だと平均と標準偏差が未知数ですし、回帰分
3790 析では平均の中に切片と傾きといった未知数が入っているわけですが、このような場合の事後分布は同時確
3791 率空間ということになります。すなわち平均と標準偏差という 2 つの変数の場合でも、可視化するなら 3 次元
3792 空間が必要です (x 軸に平均、 y 軸に標準偏差、 z 軸に確率密度)。こうした多次元空間において、あるパラ
3793 メータだけについて期待値を計算したい、という場合はそれ以外のパラメータについては**周辺化**といって積
3794 分して全部の可能性をつぶしてしまう必要があるのですが、この計算も解析的にやるには実に大変なもので
3795 す。しかし多次元の事後分布空間から生成された乱数があるなら、注目したい変数だけについての記述統計
3796 をすれば、他の変数を周辺化したこと同じになります。なんと便利なのでしょう。

3797 このように、乱数を使ったアプローチは計算上非常に有利な特徴を揃えています。乱数を大量に発生させ
3798 られるような計算機のパワーは、最近の PC でしたら十分持っていますから、最近になってベイズ統計学や
3799 モデリングアプローチが生きてきたわけですね。さらにありがたいことに、事後分布から乱数を発生させるた
3800 めのプログラミング言語ツールが登場したのも大きいでしょう。古典的には BUGS というソフトウェアが、最

*6 嘘の 538(ゴサンパチ) という標語があって、人間がふと思いつきで数字を作ろうとすると 5, 3, 8 が多くできてしまうと言われて
います。エビデンス (出典) がわからないので本当かどうかわかりませんが...

*7 余談ですが、スマホアプリなどで「ガチャを引く」というのも内部では乱数が生成されていて擬似的に「偶然あたりが出た」のを
装っているに過ぎません。私はゲームなどをやる時、擬似乱数に思いを馳せ「どこかの誰かに遊ばされている」と思うとやる気がな
くなるので、ガチャを引くようなシーンは興奮めしてしまいます。やはりマルチプレイヤーゲームのように、人間が背後にいたほうが
よっぽど意外な行動が見られますし、ゲームよりもリアルの世界の方が擬似的でない本物の偶然を楽しむことができてもいい
と思います。皆さんはどうお考えですか。

3801 近では JAGS や Stan といったのがそれで、**確率的プログラミング言語 (stochastic programming**
 3802 **language)** と呼ばれたりします。これらの言語では、尤度と事前分布をモデルとして表記し、それにデータを
 3803 与えてやることで、事後分布からの乱数を生成します。いわば万能乱数発生器なのです。こうした環境が整っ
 3804 たことで、簡単に分析できるようになりました。

3805 17.4 乱数によるアプローチの例

3806 それでは MCMC を使った実践に入る前に、乱数を使って何ができるのか、R で確かめてみましょう。

3807 17.4.1 乱数による近似の例

3808 まずは馴染み深い**正規分布 (Normal distribution)** の例からいきましょう。

3809 正規分布に感ずる R の関数は `dnorm`, `pnorm`, `qnorm`, `rnorm` などがあります。この `d`, `p`, `q`, `r` は他の
 3810 確率分布を表す関数の前につける接頭語で、`d` は確率密度、`p` は累積確率、`q` はある累積確率になるときの
 3811 確率点、そして `r` が乱数発生を意味しているのです。

code : 17.1 乱数のコードの例

```
3812 1 rm(list = ls())
3813 2 set.seed(12345)
3814 3 dnorm(0, mean = 0, sd = 1)
3815 4 pnorm(0, mean = 0, sd = 1)
3816 5 qnorm(0.6, mean = 0, sd = 1)
3817 6 rnorm(10, mean = 0, sd = 1)
3818 7 set.seed(12345)
3819 8 rnorm(5, mean = 0, sd = 1)
3820
3821
```

3822 ■コード解説

- 3823 1 行目 環境の初期化
- 3824 2 行目 乱数発生開始点の設定
- 3825 3 行目 標準正規分布の、確率点 $x = 0$ における確率密度の計算
- 3826 4 行目 標準正規分布の、確率点 $x = 0$ までの累積確率
- 3827 5 行目 標準正規分布の、累積確率 60% になるときの確率点
- 3828 6 行目 標準正規分布から乱数を 10 個出力する
- 3829 7 行目 乱数発生開始点の再設定
- 3830 8 行目 標準正規分布から乱数を 5 個出力する

3831 このコードを使って確認しておいて欲しいところは、3 つあります。まず `d`, `p`, `q`, `r` をつけたときの意味で
 3832 す。それぞれ正規分布のどういう数字を返しているのか、しっかり確認しておいてください。次に 7 行目にあ
 3833 る乱数の発生です。これを実行すると、標準正規分布にし違う乱数が 10 個出力されます。乱数ですので、数
 3834 字の並びに規則性はありません。バラバラの数字が出ていることを確認してください。最後に 8 行目、9 行目
 3835 の内容です。8 行目は 2 行目と同じく、乱数の開始点を定めています。R で出力される乱数は、規則性がな
 3836 い数字とは言え、計算によって算出している数字ですから規則性は拭いきれず、再現してしまいます。どこか
 3837 ら計算を始めるか、という開始点は**シード値 (seed value)** と呼ばれ、ここに入力した数字が開始点になり
 3838 ます。シード値の設定に全く意味はなく、好きな数字にさせていただいて結構です。ポイントは、任意の数字で

3839 あっても同じ数字に設定すると、乱数はそこから計算して算出されますので、9行目で実行した5つの乱数
 3840 は7行目の最初の5つと同じ数字が再現されているというところでは、規則性のない数字なので再現されて
 3841 は困ると思うかもしれませんが、科学計算のアプローチという意味では再現性が担保できることも重要な
 3842 です。

3843 では乱数を使う例をもう少し見てみましょう。次はベルヌーイ分布 (Bernoulli distribution) について
 3844 考えてみたいと思います。これはコイントスをして表が出るか、裏が出るかという0/1の離散的結果を生み出
 3845 す分布です。残念ながらRにはベルヌーイ分布の関数はなく*8、二項分布の特殊例として使うことにします。
 3846 **二項分布 (Binomial distribution)** とは、N回コイントスをしてK回表が出る確率を表す分布ですが、
 3847 これのコイントス回数が1回であればベルヌーイ分布と同じことになります。

code : 17.2 乱数のコードの例 2

```
3848 1 rbinom(10, size = 1, prob = 0.5)
3849 2 rbinom(10, size = 1, prob = 0.3)
3850 3 rbinom(10, size = 1, prob = 0.7)
3851 4 # パッケージの利用
3852 5 library(extraDistr)
3853 6 rbern(10, prob = 0.5)
3854
3855
```

3856 ■コード解説

3857 1行目 二項分布の乱数を10個発生させる。確率は0.5で。
 3858 2行目 二項分布の乱数を10個発生させる。確率は0.3で。
 3859 3行目 二項分布の乱数を10個発生させる。確率は0.7で。
 3860 4-6行目 extraDistr パッケージによるベルヌーイ乱数例

3861 ここで確認して欲しいのは、コイントスを10回するとして、確率0.5の場合はほぼ半々、0.3の場合は半
 3862 分以下、0.7の場合は半分以上表(1)が出ている、ということです*9。ここで確率を幾つにするかは、われわ
 3863 れが自由に設定できます。また乱数を幾つ発生させるかも自由です。これらの数字を色々変えて遊んでみま
 3864 しょう。

3865 たとえば乱数の数を100回、1000回、10000回と増やしたとします。その結果をすべて表示するのは大
 3866 変ですが、その平均値を計算してみるとどうなるでしょうか。

code : 17.3 乱数のコードの例 3

```
3867 1 N100 <- rbinom(n = 100, size = 1, prob = 0.5)
3868 2 N1000 <- rbinom(n = 1000, size = 1, prob = 0.5)
3869 3 N10000 <- rbinom(n = 10000, size = 1, prob = 0.5)
3870 4 mean(N100)
3871 5 mean(N1000)
3872 6 mean(N10000)
3873
3874
```

3875 ■コード解説

3876 1-3行目 二項分布の乱数を100個、1000個、10000個発生させる。確率は0.5で。

*8 extraDistr パッケージを用いると、bern という関数名でベルヌーイ分布が使えます。

*9 正確には二項分布なので、1が出た回数が出力されているのです。試行数 (size) を1にしているの、1回中1回表が出た=表(1)が出た、と解釈しても、結果的に同じことであるというだけです。

3877 4-6 行目 それぞれの平均値を計算する

3878 これの実行結果は、私の環境では出力 17.1 のようになりました。

R の出力 17.1: 乱数出力の結果

```
> mean(N100)
[1] 0.62
> mean(N1000)
[1] 0.52
> mean(N10000)
[1] 0.5003
```

3879

3880 もともと確率を `prob = 0.5` と設定しましたから、理論的には 0.5、つまり半々の確率で表 (1) が出たり裏
3881 (0) が出たりするはずですが、最初の 100 回では 62 回なのでやや表が多めに出たようです。乱数ですので、
3882 そういうことはあります。ただ、この傾向も、1000 回、100000 回と繰り返していくと、ほぼ偶然による偏りは
3883 なく、半々の比率になっていくのが確認できると思います。このように、乱数が十分多ければ、確率分布の近
3884 似値として使うことができるというわけです。

3885 ちなみに今回は平均をとりましたが、これは確率分布の**期待値 (Expectation)** を計算したことと同じにな
3886 ります。

3887 17.4.2 乱数を用いた確率分布のプロット

3888 乱数を使って近似できる、という例を可視化で見てください。また正規分布を例にしてみたいと思います。

code : 17.4 確率分布の可視化 1

3889

```
3890 1 x <- seq(-4, 4, 0.05)
3891 2 plot(x, dnorm(x))
3892 3 library(tidyverse)
3893 4 ggplot(data.frame(x = c(-4, 4)), aes(x = x)) +
3894 5   stat_function(fun = dnorm)
```

3896 ■コード解説

3897 1 行目 -4 から 4 までの範囲で、0.05 刻みの数列を作る

3898 2 行目 低水準プロット。各確率点の確率密度をプロットする

3899 3 行目 パッケージの読み込み

3900 4-5 行目 ggplot による美しい描画。stat_function は関数の結果を図にするもの

3901 これは理論的な正規分布の形を図示するものです。美しいですね。でも同様のことが、乱数を使ってもでき
3902 ます。

code : 17.5 確率分布の可視化 2

3903

```
3904 1 N <- 1000
3905 2 X <- rnorm(N, mean = 0, sd = 1) %>% as.data.frame()
3906 3 ggplot(data = X, aes(x = .)) +
3907 4   geom_histogram()
```

3909 ■コード解説

3910 1 行目 発生させる乱数の数

3911 2 行目 乱数を発生させ、data.frame 型に組み上げる

3912 3-4 行目 ggplot による美しい描画。geom_histogram になっているところに注意

3913 このコードの結果描かれるカーブは、それほど美しいものではないかもしれませんが。しかし発生させる乱数の
 3914 数が増えればどうでしょうか。どんどん形が似ていくことがわかると思います。ここでのポイントは、ヒストグラ
 3915 ムを描いたらその稜線が**確率分布**の形になること、乱数の数が増えれば十分な近似になることです。よく確認
 3916 しておいてください。

3917 17.4.3 乱数を用いた確率分布の要約

3918 最後に、確率分布の特徴を記述するコードの書き方について練習しましょう。乱数発生によるアプローチ
 3919 は、記述統計量で確率分布の特徴を記述できます。確率分布の期待値は、その算術計算の平均で良いので
 3920 すが、中央値やパーセンタイル、分散や標準偏差も分布を記述する重要な指標です。さらに確率分布におけ
 3921 る、確率密度が最も高くなる点は、「最も生じる確率が高い点」という意味で重要ですが、それを求めるため
 3922 には特別な関数を書く必要があります。

3923 まずは data.frame 型にした正規乱数を使って、記述統計量を算出する例を見てみましょう。

code : 17.6 確率分布の要約

```

3924 1 N <- 10000
3925 2 X <- rnorm(N, mean = 0, sd = 1) %>%
3926 3   as.data.frame() %>%
3927 4   dplyr::rename(val = 1)
3928 5 X %>%
3929 6   summarise(
3930 7     Exp = mean(val),
3931 8     SD = sd(val),
3932 9     Median = median(val),
3933 10    U50 = quantile(val, prob = 0.50),
3934 11    U90 = quantile(val, prob = 0.90),
3935 12    L90 = quantile(val, prob = 0.10)
3936 13 )
3937
3938

```

3939 ■コード解説

3940 1 行目 乱数を 10 万個作ることにします。

3941 2 行目 標準正規分布の乱数発生

3942 3 行目 データフレーム型に整形

3943 4 行目 変数名がつけられていないので、val という名前をつけることに。ここでの 1 は 1 列目を意味して
 3944 いる。

3945 6 行目から データフレーム X を使って計算

3946 7 行目 summarise 関数は要約計算。変数名 = 計算式で表す。

3947 8 行目 期待値を mean 関数で算出

3948 9 行目 確率分布の標準偏差を sd 関数で算出^{*10}

3949 10 行目 中央値を median 関数で計算

3950 11-13 行目 パーセンタイル点を算出。50 パーセンタイル点は中央値に同じ。

3951 また、確率密度が最も高くなる点を算出するための関数を自分で作ってみることにします。次のようにコードを入力してください。

code : 17.7 確率分布の要約

```
3953
3954 1 map_estimation <- function(z) {
3955 2   density(z)$x[which.max(density(z)$y)]
3956 3 }
3957 4 # 試してみる
3958 5 X %>%
3959 6   summarise(
3960 7     MAP = map_estimation(val)
3961 8   )
3962
```

3963 ■コード解説

3964 1 行目 関数の宣言。map_estimation という関数名を作る。この関数は引数 z を取る

3965 2 行目 density 関数は与えられた数列の関数を考え、そのカーネル密度を計算します。which.max 関数はその中でも最も密度の高い値を返すものです。

3966 3 行目 関数の終わり

3967 5-8 行目 さきほどの例で試してみる。

3969 このようにして、確率密度関数の近似値が簡単な関数で求められます。

3970 次回以降は、モデルに基づいて計算された事後分布から得られる乱数について、これらの関数を使ってその特徴を記述していくことになります。関数や指標の意味をしっかりと理解しておいてください。

3972 17.5 課題

3973 次の計算をする R コードを記述し、提出してください。提出は R スクリプトファイルでも Rmd ファイルでも構いませんが、どの課題に対するコードなのかがわかるよう、コメントや説明文を記入しておくこと。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

3977 1. 平均 50、標準偏差 10 の正規分布に従う乱数を 10 万点生成し、その平均値、中央値、標準偏差、15 パーセンタイル、75 パーセンタイル、2.5 パーセンタイル、97.5 パーセンタイルを算出してください。

3979 2. さきほど求めた正規乱数の記述統計量が、理論値の近似値になっていることを確認します。dnorm, pnorm, qnorm などを使って理論値を算出してください。

3981 3. 正規乱数のデータセットのうち、30 以上 60 未満の値が含まれる割合を計算してください。またそれが理論値の近似値になっていることを確認するため、dnorm, pnorm, qnorm などを使って理論値を算出してください。

^{*10} R の sd 関数は不分散の平方根を計算しており、厳密には標本標準偏差の計算で良いが、サンプルサイズが十分に大きいのでほとんど影響しません。

第 18 章

いんたーみっしょん ; Stan の概略と環境の準備について

このセクションでは、確率的プログラミング言語 Stan を導入するにあたっての、周辺知識を解説します。授業中に解説するものではありませんし、ここに書かれていることのすべてを理解していないと Stan を使えないというわけではありませんが、導入や利用にはトラブルやエラーも多く、その際に前提知識、周辺知識があるとないとでは理解度が大きく異なります。「わからなくても、なんとなく動いた」という状態で受講し続けるのは自身のためにならないだけでなく、面白くないです。知識があつてはじめて価値がわかるということもありますので、共用としてご一読いただければと思います。

18.1 はじめに

この授業では統計言語としての R、それを使う統合開発環境としての RStudio、そして確率的プログラミング言語 (stochastic programming language) としての Stan を利用します。大学での PC ルームの利用にはこれらの環境がすでにある程度準備されていますが、最新バージョンではありませんので、自分の PC に環境を準備することを強く推奨します。今後の研究室配属やその後の卒業研究などでも活用することになると思いますので、まずは自身の PC 環境に準備することを考えてみてください。

この付録資料は、これらの環境を準備する方法について解説するものですが、提供されるソフトウェアのバージョンや対応する環境などは日々発展するものですので、必ずしも最適な情報提供になり得ません。基本的な情報は提供いたしますが、執筆時^{*1}での情報であることも多く、より新しい情報についてはインターネットなどでキーワード検索を行なって、より良いものを探してみてください。

PC の操作が苦手に思っている人のために付言しますが、うまくいかなかったからといって癩癩を起こしたり、諦めたりしてはいけません。相手は機械ですから、命じられたことを愚直に実行しているだけです。また文房具のようなものですから、恐れて嫌うのではなく、愛して可愛がってあげることが肝要です。もしみなさんの中に、お持ちの PC に名前をつけていない人がいるようでしたら、今すぐ命名することをオススメします。パソコンが、と思うと腹が立ちますが、わたしの XXX ちゃんが、と思うとエラーも「ちょっとご機嫌斜めなのかしら」と受け入れやすくなります^{*2}。

また、困った時は教員、TA、先輩などに相談することが重要ですが、その際には症状や経緯を正確に伝えることが必要です。ペットを病院に連れていくときに「何もしてないのに勝手に病気になった。治してほしい」

*1 この記事は 2022 年 11 月 10 日に加筆修正されました。

*2 ちなみに私が初めて手に入れたノート PC には「秀吉」という名前をつけました。Mac に変わってからは代々数学者の名前をつけることにしています。私はいま、Gauss ちゃんや Hermite ちゃんを持ち歩いています。

4011 といわれても獣医さんは困ると思います。普段の様子や症状についての丁寧な報告を心がけてください。と
 4012 くに PC は機械ですので、何もしていないのに壊れるということはありません。自分が自覚していないことで
 4013 あっても、「なにかをした」から様子がおかしくなるのです。自分はそんな大それたことはしていない*3、と思っ
 4014 ても必ず何かトリガーがあったはずなのです。どこをクリックしたか、どういうコマンドを入れたかという履歴
 4015 をしっかりと把握し、丁寧に報告することを心がけてください。

4016 18.2 Stan の位置付け

4017 **確率的プログラミング言語 (stochastic programming language)** とは、確率モデルを記述し、デー
 4018 タと合わせることで統計的推論をしてくれるコンピュータ言語のことです。Stan ができる前は、BUGS や
 4019 JAGS(と) いうものがありました。BUGS は開発が終わってしまいました*4。今、こうした言語は Stan が
 4020 最先端だといって間違い無いでしょう。これらの言語は、具体的には確率モデルとデータの関係性を記述するこ
 4021 とで、事後分布からの乱数を生成するというものです。以下ではこの言語の利用形式について解説します。

4022 18.2.1 コンパイラとインタプリタ

4023 プログラミングはコンピュータを動かすための仕様書を書くことです。ここで「コンピュータ」というのは計算
 4024 機という意味だと思ってください。もちろんコンピュータは数字の計算だけでなく、音楽を奏でたりゲームを楽
 4025 しませてくれたり、と色々なことができるのですが、その背後にあるのはとにかく 0/1 の数値演算です。0 か
 4026 1 かという数字がどうして映像や音声、通信になるのか、と疑ってしまうほど異なっているようですが、それ
 4027 もすべて数字の計算から構成されています。

4028 コンピュータができ始めたごく初期の頃は、これを動かすのに 0 と 1 からなる数字の羅列で仕様書を用意
 4029 していました。流石にそれではわかりにくく表現力にも乏しいので、次にできたのがマシン語とよばれる 16 進
 4030 数による表現でした。この段階でも、知らない人にとっては暗号か、意味のある連なりに思えない文字列にす
 4031 ぎません。画面に線を引いて欲しい時に line という命令で伝えられるようになって、やっと普通の人間にも
 4032 意味がわかるレベルになってきます。このように、人間がわかるレベルでコンピュータに命令が伝えられるよう
 4033 な言語のことを**高級言語**と言います。高級言語は人間寄りで、直接コンピュータが理解できませんので、この
 4034 高級言語を機械の言葉に翻訳するためのアプリケーションを介させます。それがいわゆる**プログラミング**
 4035 **言語 (programming language)** と呼ばれるものです。

4036 プログラミング言語は、古くは BASIC, PASCAL, FORTRAN などと呼ばれる書式のものが、その後
 4037 出てきた C 言語、その改良版である C++ 言語*5などが有名です。他にも JAVA や Object C, Python
 4038 などが有名ですね。R も言語の一種で、統計解析に特化したものです。また、この授業で扱う Stan という**確**
 4039 **率的プログラミング言語 (stochastic programming language)** は、確率モデルの分析に特化したも
 4040 のです。みなさんがデータサイエンス業界に進むのであれば、Python や R を使えるようになっておくとい
 4041 いでしょう。これらの言語の特徴の 1 つは、パッケージを追加することで機能が拡張できることにあります。と
 4042 くに**機械学習系 (AI など)** のパッケージが豊富なのが Python です。またこれらの言語は基本的にフリーソ
 4043 フトウェアであり、導入にお金がかからないところもポイントですね。タダで始められるおもちゃみたいなもの

*3 たとえば計算途中だけと時間が来たので電源をオフにした、といったことでも、内部ファイルにアクセスしているときの中断であれば、十分に PC を破壊することができます。

*4 JAGS はまだ開発が続いているようで、R から使うこともできます。

*5 この ++ という書き方は、C 言語特有の「変数に 1 加える」という表記法であり、C 言語の改良版という意味で C++ と命名されています。読み方はシープラスプラスですが、シブラブラの愛称でも知られています。

4044 です！*6*7*8*9

4045 プログラミング言語の分類には、その設計思想や書き方などを基準にすることもできますが、実行方法を基
4046 準にするものとしてインタプリタ型とコンパイラ型とに分けることができます。インタプリタ型は、interpret、
4047 つまり翻訳型です。これは毎回の命令を機械語に翻訳して実行していきます。R はインタプリタ型言語で、毎
4048 回の命令を逐一翻訳して実行していきます。R を実行するとき、コンソールに>というマークが出ていますよ
4049 ね。これは R が聞き耳を立てている状態、入力待ちの状態なのです。ここに命令を書く（たとえば 2+3 と書
4050 く）、R はすぐさま答えを返してきます（たとえば [1] 5 と返す）。基本はこのように一問一答、ひとつひとつの
4051 命令を毎回機械語に翻訳して実行し、その答えを出力するという手続きをとっています。もちろん私たちの分
4052 析プログラムは、一行で書ける簡単なものではありませんから、複数行に渡る長々としたファイルを書くことが
4053 多いでしょう。こうしたファイルはプログラムともいわれますし、一行一行の命令のことをスクリプト (script)
4054 ということもあります。プログラムファイルあるいはスクリプトファイルを開いて、一行ずつ実行するのが R の
4055 基本的なスタイルであり、RStudio はスクリプトファイルを編集するエディタ (editor) がセットになってい
4056 るのが便利な点なのでした。

4057 これに対して、インタプリタ型は、スクリプトファイルをまとめて機械語に翻訳 (コンパイル) し、実行ファイル
4058 というのを別途作ります。その上で、実行ファイルを実行すると計算が進められるのです。どうしてこんな手間
4059 がかかることをするのでしょうか。ひとつには、ひとつひとつの命令分を逐一翻訳しているのでは、実行スピー
4060 ドが遅くなるということが挙げられます。命令を逐一聞き耳を立てて待ち、全ての命令を聞き終わるまで途中
4061 の指示を記憶しておいて、命令文が終わってはじめて行動に移す、というのでは記憶容量も時間もかかるわ
4062 けですね。これに対して、一連の計算命令文書を一括で渡すことができれば、機械は内部的に効率よく計算
4063 手順を配置して実行できます。コンパイラ型は計算スピードが速いのです。ちなみに実行ファイルは機械がわ
4064 かる言葉に書き換わっているの、人間が見ても意味がわかりません。また、OS や CPU に理解できる形に
4065 書き換えていますので、MacOS の実行ファイルを Windows で実行する、ということではできません。翻訳後
4066 の言語が違うからです。スクリプトファイルは環境を超えて共有できますが、それを翻訳したり最適化したりす
4067 る仕組みは、個別の環境に依存します*10。

4068 Stan はコンパイラ型です。Stan のファイルを読み込んで、機械にわかるように「事後乱数生成命令文」に
4069 翻訳して実行します。確率モデルから乱数を生成するのは非常に高度な計算機能ですから、コンパイルして
4070 まとめておくことでスピードアップの効率が断然良くなるわけです。コンパイラ型の利点はこのスピードにあ
4071 りますが、欠点は「命令に変更やミスがあれば翻訳し直さなければならない」というところです。あり得ない命
4072 令を出しているようであれば、コンパイルの段階で「おかしな文法だよ」とエラーを出して翻訳を止めてくれま
4073 す。翻訳が通れば良いかというそうではなく、翻訳はあくまでも文法上の手続きであって、数値や内容に間
4074 違いがあっても翻訳自体は完了させることができるのです。翻訳 (コンパイル) には少し時間がかかりますか

*6 余談ですが、Scratch や任天堂 Switch の「ナビつき! つくってわかる はじめてゲームプログラミング」などにみられる、GUI ベースのプログラミング言語もあります。これはプログラミングの要素がアイコンやキャラクターになっていて、それらを組み合わせることで計算させるという方法をとっています。その操作のしやすさ (ミスペルがない!) から、小学生などにも導入できる素晴らしい取り組みです。同様の発想は、LOGO 言語などにもみられました。

*7 Perl や JavaScript など、Web ブラウザ上で動くプログラミング言語もありますが、これはブラウザ上での操作に限定されており、数値計算にはあまり向かないのでここには取り上げていません。さきほどあげた JAVA と JavaScript は別物なので注意してください。

*8 GUI アプリケーションを作ることに特化した言語もあります。数値計算と違って、マウスがどこでクリックされたか、と言った「動き」を契機としてプログラムが走る、イベントドリブン型の言語で、Microsoft 社の Visual Basic などが有名です。これは Microsoft Office 製品の中に埋め込むこともでき (VBA)、これを使うと Excel でも高度な統計解析ができたりします。Excel を使った統計ソフトウェアとして代表的なものに HAD があります (清水, 2016)。他にも Delphi などの言語がありました。

*9 (シ, 2016) には、数え方にも癖がありますが、117 ものプログラミング言語が紹介されています。

*10 ちなみに R や Rstudio はもちろん、Word や Excel などのオフィスアプリ、果ては OS そのものも、コンパイルされた実行ファイルを PC が計算して実行しているに過ぎません。

4075 ら、文法上も内容上も間違いのない命令分をしっかりと準備しておくことが重要になってきます。

4076 18.2.2 ファイルのやり取りについて

4077 ここで、R/RStudio をつかって Stan を使う際の、ファイルのやり取りを見ておきましょう (図 18.2)。

4078 Stan のコードは R のコードの中にも書くこともできますが、別のファイルに保存しておくのがベターです。
4079 Stan ファイルは拡張子を `.stan` とすることが一般的です。このファイルの中身はプレーンテキストでコードを
4080 書いていますから、一般的なエディタ^{*11}で編集できます。RStudio のエディタ画面も必要十分なエディタ機
4081 能を持っていますので、RStudio をエディタとして利用するのがいいでしょう。

4082 RStudio で File > New File として新しいファイルを開くときに、Stan ファイルとして開くことが可能で
4083 す。Stan ファイルとして開くと、初心者向けの配慮からか、最初からちょっとしたコードがすでに書かれていま
4084 す。コードの書き方を示すサンプルコードなのですが、実際にこのコードを使って何かすることはありませ
4085 ぬで、中身は全部削除してしましましょう。それでも RStudio はその画面が Stan ファイルのものだということ
4086 は認識していますから、そこで書いたファイルは Stan ファイルとして扱ってくれます。ここでまちがって新しく
4087 R Script で開いてしまった、というときに、保存するときだけ `.stan` をつけておけばいいか、という対応をす
4088 ると失敗します。RStudio では拡張子をあえて表示していませんので、ファイル名として `hoge.stan` と書く
4089 と実際には `hoge.stan.R` というファイル名になっています。拡張子は最後のピリオドで判断されますので、
4090 これは R ファイルなのです。このミスを防ぐために、ファイル名にピリオドは使わないこと、ファイルの種類を変
4091 える時はエディタ画面の右下でファイル種別を選択してください (図 18.1)。

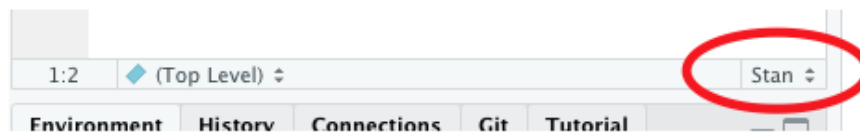


図 18.1 ファイル種別の変更

4092 RStudio を使って Stan を書くとき便利なことがいくつかあります。ひとつは強調表示機能で、ブロッ
4093 クや関数名など Stan で使うことが決まっている専門用語は色やフォントが変わって表示されます。
4094 `transformed parameters` のように長い専門用語を書くときはミススペルが心配ですが、書いた後で強調
4095 されなければスペルミスがある、ということがわかります。もうひとつの利点は文法チェックの機能です。コン
4096 パイル型なので、Stan に書かれた命令文は実行前に一括で翻訳されますが、スペルミスや型の違いなど、文
4097 法的に間違っているところがあればこれも自動的にチェックして警告記号がでます。また Stan ファイルが開
4098 かれていたペインの左上に Check というボタンがあり、これを押すことでファイル全体の文法チェックをして
4099 くれます。もし文法上の問題がなければ、「hogehoge.stan is syntactically correct.(hogehoge.stan という
4100 ファイルは文法的には正しいですよ)」というメッセージがコンソールに表示されます。逆にいうと、これが表示
4101 されないというときは何か間違っている、コンパイルに進む前に修正しましょう。

4102 さてこのようにして Stan ファイルを準備します。また分析用のデータセットが外部ファイルにある場合など
4103 は、これまで同様、当該プロジェクトフォルダ内に置いておくといいでしょう。これらはいずれも、R のコードで

^{*11} エディタとは ASCII ファイル、いわゆるプレーンテキストを編集するアプリケーションの総称で、OS にデフォルトで「メモ帳」などの名称で含まれています。デフォルトのアプリは文字を読み書きできるだけで、もう少し発展的、あるいは便利な拡張機能を持っている専用のアプリを使うことが一般的です。筆者は最近もっぱら VS Code というエディタを利用しています。他にも macOS でしたら、「mi」というアプリを好んで使っていました。Windows でしたら「秀丸」というアプリが有名です。Ubuntu には gedit というアプリがデフォルトで入っていますし、歴史やユーザの多さでは vim や emacs などがあります。Linux 界隈で、vim 派か emacs 派かの話題は、きのこたけのこ戦争ばかりに激しい戦いになります。

4104 呼び出して使うことになります (図 18.2 の 1. と 2. のステップ)。R では、Stan ファイルやデータファイルを
 4105 読み込んで、これを PC 内部にある Stan に渡します (図 18.2 の 3. のステップ)。Stan は R と違う言語で
 4106 あり、Stan 独自の計算機能で計算してくれるわけですから、R と Stan の橋渡しが必要になります。R では
 4107 これを `cmdstanr` や `rstan` というパッケージを経由して行います。これらのパッケージが Stan を呼び出し
 4108 てくれるわけです*12。

4109 データと確率モデルの関係を記した Stan ファイルと、データのそのものを受け取った Stan は、コンパイル
 4110 して PC 内に乱数発生装置を作ります (図 18.2 の 4. のステップ)。事後分布の関数を書き出すのではなく、
 4111 その形状とそれを代表する乱数を作る状態をもつわけです。R の指示をうけて、そこから必要な数だけ乱数
 4112 を作り出します (図 18.2 の 5. のステップ)。R はこれを受け取って、統計解析を始めるわけです。R は統計
 4113 に特化した言語ですから、集計や可視化などの作業はお手のもの。事後分布から得られた大量の乱数は、事
 4114 後分布の特徴を反映した大量のデータセット、事後分布の具体例たちになっているのです。

4115 分析の進め方は以上ようになります。R でデータを整形し、Stan のファイルも別途用意して、R から
 4116 Stan を呼び出す、Stan が返してきたデータを R で解析する、ということです。R の拡張子は `.R` で、Stan
 4117 の拡張子は `.stan`、データファイルの拡張子は `.csv` などでしょうか。いずれにせよ、複数の種類のファイルを
 4118 やり取りしながら進めるので、相互の関係についてしっかり理解しておく必要があります。

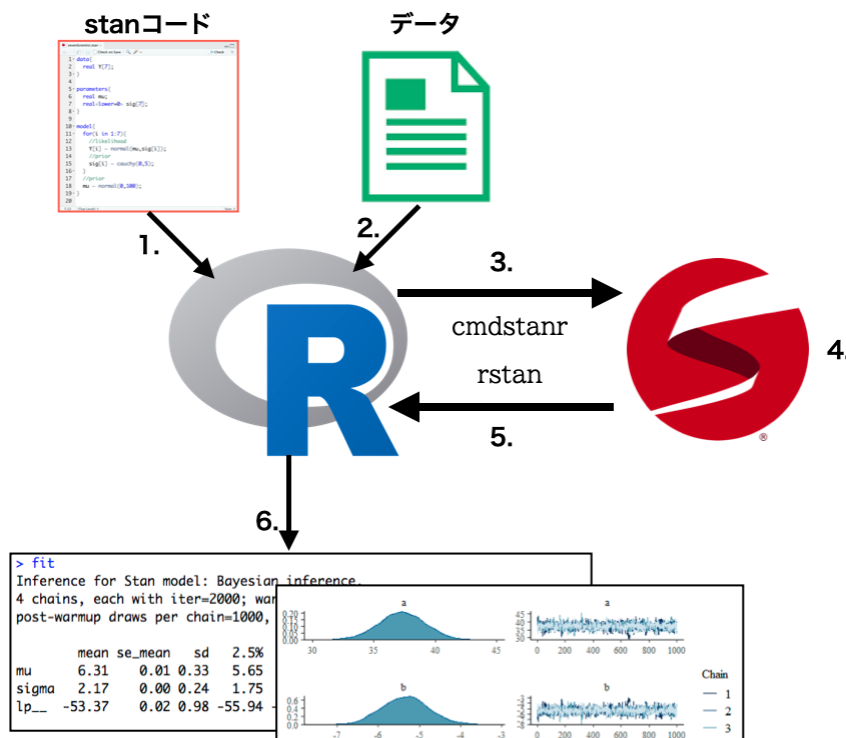


図 18.2 R/RStudio と Stan ファイルのやり取り

*12 このように、Stan は独立した計算環境であり、R 以外のアプリケーションから呼び出して使うことが可能です。Python から使いたい場合は、`PyStan` というパッケージ経由で、Julia から使いたい場合は `StanJL`、コマンドラインから使いたい場合は `cmdstan` といったように、いろいろなルートがあります。

4119 18.2.3 2つのルート

4120 ここでは R から Stan を呼び出す 2 つのルート, 具体的には `rstan` パッケージと `cmdstanr` パッケージ
 4121 について説明します。この 2 つのパッケージは, いずれも Stan を呼び出してつかう, R と Stan の間を取り
 4122 持つインターフェイスの役割を持ったパッケージです。間を取り持つだけなので, Stan ファイル自体は同じで
 4123 あっても構いません。同じ Stan ファイルを `rstan` から呼んでも, `cmdstanr` から呼んでも, 結果は同じで
 4124 す。ただしインターフェイスが違うので, 結果の扱い方が変わってきます。他にも違うところがいくつかありま
 4125 すので順に説明していきますが, 結論からいうとこれから^{*13}は `cmdstanr` のほうを使うほうがお勧めです。

4126 歴史的には `rstan` のほうが古くから存在します。このパッケージも最初の頃は導入に一苦労するもので
 4127 したが, 2016 年ごろから CRAN を通じて配布されるようになりました。要するに, 他のパッケージと同じ
 4128 ように, `install.packages("rstan")` と R で書くだけで導入できるようになったのです。このようにし
 4129 て導入するとわかりますが, `rstan` は多くの依存パッケージがあります。Stan はコンパイルするのに OS
 4130 に応じた C++ コンパイラを借用しますが, R から C++ を呼び出すパッケージや R と C++ の間に入る
 4131 StanHeader と呼ばれる緩衝材など, 関連する多くの環境を整えてやっと R から Stan へのルートが通じる
 4132 のです。ユーザにとってはそのような苦労は知ったことではない, という感じですが, 「間に多くの調整が入る」
 4133 ということがエラーの温床になっていました。つまり, いくつかの内部的ステップのどこかでバグがある, 互換
 4134 性がない, というようなことがあると, Stan ファイルが同じでも, バージョンが変わるとコンパイルできないと
 4135 いうことがよくあったのです。Stan 本体の方も開発が盛んですから, Stan の新しいバージョン, 新しい機能
 4136 を使いたいと思っても, パッケージが対応するのを待たなければなりません。パッケージが対応しても, 途中
 4137 の媒介パッケージが追いついていなければバグになったりします。これらを使う R や OS もアップデートして
 4138 いきますから, その度に「動かなくなる」ということはよくありました。研究の再現性が問題になる昨今ですが,
 4139 自分の研究とデータであっても, 自分の環境で再現できないということが少なくなかったわけです^{*14}。そうな
 4140 ると, 今動いているからそれでいい, と環境のアップデートを控えるようになりますが, ご存知の通り OS やア
 4141 プリのアップデートはバグやセキュリティホールへの修正など, 安全につかうための基本的な機能に関わってき
 4142 ますから放置もできない, という問題があるわけです^{*15}。

4143 このような不安定な環境になる原因は, R と Stan の間にあるさまざまな障壁と, それをなんとかして調整
 4144 しようという媒介パッケージの存在でした。そこでなるべくシンプルに R と Stan をやりとりする方法はないも
 4145 のか, ということで考えられたのが `cmdstanr` です。これはコマンドライン^{*16}から Stan を呼び出す `cmdstan`
 4146 を R から呼び出して使おう, というものです。Stan を R になんとか取り込んで動かす, というのではなく,
 4147 Stan の部分は OS のプリミティブな環境にお任せ, R は結果をもらうだけ, というシンプルな設計にしよう
 4148 というわけです。たとえば `rstan` は Stan の計算結果を R オブジェクトとして取り込みますが, `cmdstanr` は
 4149 実は MCMC の結果は csv ファイルのような形で吐き出しており, それを読み込んで使います。御用聞きや
 4150 仲介業者を介して情報をもらうのではなく, 現地で直接買い付けしてくるようなイメージでしょうか。このよう

^{*13} 執筆時点, 2022-11-14 現在

^{*14} 実際 R や OS のアップデートに関わるエラーは大変で, 筆者は一度 OS をアップデートした後半年ほど Stan が使える環境が
 作れず, 研究が頓挫したことがあります。そのために別途, PC を購入し直したぐらいです。

^{*15} OS のアップデートなど, 通知が来ても「今支えているからいいや」と無視する人もいますが, できれば OS やアプリは最新版であ
 るほうが良いのです。セキュリティホールの問題などを放置しておく, インターネットを介して自分の PC の中身が全部見られて
 公開される危険性があります。Stan しか使わないネットに繋がらない PC というのがあればいいのかもしれませんが, 文房具とし
 ての PC はもっと一般的な用途がありますから, そうもいつてられませんよね。また, 一昔前の OS アップデートは, アップデート
 の失敗で PC が動かなくなるというようなこともありましたから, アップデートが公開されてすぐに対応するのは「人柱」などと揶
 揄されていましたが, 最近はそういったこともほとんどありませんので, 安心して常に最新の状態でしておいてください。

^{*16} MacOS や Linux では端末, ターミナル, Windows ではコマンドプロンプトなどと呼ばれる, PC に直接コマンドで命令を出
 すプリミティブな実行環境です。

4151 に間に挟むものを減らすことによって、不安定な要素をなくしたわけです。

4152 `cmdstanr` は Stan と R の複雑な絡まりを排していますので、Stan の開発が進めばすぐにそれを R に届
4153 けることができます。たとえば 2022/11/14 現在で、公式にリリースされている `rstan` のバージョンは 2.21.1
4154 ですが、`cmdstanr` が呼び出す `stan` のバージョンは 2.30.1 です。`cmdstanr` のほうが新しいものに対応
4155 できていますね。このように、`cmdstanr` がでてきてからはこちらの対応、反応が早く、また動作も安定的で
4156 すから、よりお勧めしやすくなっています。

4157 これまでに出ている R で Stan を使う方法を解説したテキストの多くは、`rstan` をつかっていますので、
4158 今それを使うためには多少の読み替えが必要です。とはいえ、Stan のコード自体はそのまま使えるものがほ
4159 とんどで、R とのインターフェイス、やりとりの方法が違っているだけです。`cmdstanr` には、結果を `rstan`
4160 で得られたオブジェクトに変換する関数も含まれていますから、これらを使って変換すれば、以後のコードは
4161 `rstan` 準拠のものでもエラーになることはありません。このテキストでは第 2 版から `cmdstanr` を中心にす
4162 るように切り替えましたが、以前のコードでも本質的に問題はありません。

4163 18.3 導入の概略

4164 環境の導入は、OS の種類によって変わります。Windows なのか Mac なのか、あるいは Linux なのかと
4165 いった違いに合わせて導入を考えてください。OS の種類で言えば、Linux は全体的に CUI 操作が主であ
4166 り、堅牢かつ合理的な使い方ができるものです。如何せんマイノリティですので、ググってもあまり情報が出
4167 てこないところが玉に瑕です。Mac は Linux ベースの OS になっているので、比較的堅牢かつ合理的な使
4168 い方ができます。ハードウェアも OS も 1 つの会社 (Apple) が作っているので、サポートもしやすいという利
4169 点がありますが、日本では少数派なのが残念なところです。Windows は多数派の OS ですが、OS メーカー
4170 はソフトウェア会社でハードウェアが別会社なことが多く、利用される範囲は広いのですが、その分問題が生
4171 じたときにケースバイケースになりがちです。多くのハードウェアの違いを吸収するために最大公約数的な設
4172 計をするせいか、ソフトウェアデザインの統一性がないところが欠点です。ユーザ数は最も多いので、一生懸
4173 命調べたら同じような問題で困っている同じような機体の人に出会える可能性が高いのがせめてもの救い
4174 です。また光明として WindowsOS も Linux ベースの仕組みを取り入れ始め (WLS)、比較的合理的な振
4175 る舞いができるようになったことが挙げられますが、OS のグレードによって導入のしやすさが異なります。

4176 ちなみに ChromeOS や iOS, Android などタブレット・スマートフォンでの統計環境の利用は限界があり
4177 ます。これらは計算結果を利用するフロントエンドとしては非常に高機能なのですが、インタラクティブに使う
4178 ことには不向きです。いわば書籍のように多くの情報を一方的に与えてはくれるのですが、ユーザが計算する
4179 といった働きかけの補助になるような (ノートとペンのような) 使い方ができませんので注意してください。

4180 さて、環境の導入として 2 つのルートを紹介しましょう。第一のものが最も基本的なアプローチで、自分の
4181 手元の環境を作り上げるというものです。これを**ローカル環境**での構築と呼びます。ローカル環境は自分だ
4182 けの環境ということなので、自分の責任でもってしっかりと環境構築をできます。第二のアプローチとして、あ
4183 る程度できあがった環境を丸ごと取り込む、**仮想環境**での構築という方法があります。PC の OS の中に別
4184 のマシン・別の OS をさらに憑依させて使うようなやりかたで、これができると取り込む環境は完成しているの
4185 で細かい設定をする必要がありません。問題は「憑依させる」準備がいろいろ大変であるということです。ま
4186 た憑依させた PC はブラウザ経由でアクセスすることになることにも注意してください。

4187 18.3.1 導入方法 1 ; ローカル環境での構築

4188 ローカル環境での構築は、R, RStudio の導入から入ります。ここは既にできている人もいるかもしれま
4189 せんが、念のためバージョンが最新のものかどうかを確かめておいてください。ローカル環境構築の方法
4190 がわからない人は、「RStudio」「インストール」などのキーワードでウェブ検索し、自分の環境にあった例を
4191 見つけてそれをみながら進めると良いでしょう。参考までにいくつかあげておきますと、新しいものでは、高
4192 知工科大学柳井先生のこちら <http://yukiyanai.github.io/jp/resources/> とか、Quiita の記事
4193 <https://qiita.com/hujuu/items/ddd66ae8e6f3f989f2c0> が良いでしょう。書籍では、手前味噌で
4194 恐縮ですが小杉 (2019a) などが良いでしょう*17。

4195 次に確率的プログラミング言語, Stan の導入を行います。Stan を R から使うためのパッケージとして
4196 cmdstanr と rstan という 2 種類があるという話はしました。かつては rstan のほうがよく使われていた
4197 のですが、最近は cmdstanr のほうが安定的に動作し、開発も盛んになっているので、2022 年度後期から
4198 は cmdstanr をこの授業のデフォルトパッケージとしたいと思います。このパッケージの導入には、公式サイ
4199 ト <https://mc-stan.org/cmdstanr/articles/cmdstanr.html> を参考にしましょう。「cmdstanr イ
4200 ンストール」で検索すると色々な紹介サイトが出てきますので、自分と同じ環境でなるべく日付の新しいもの
4201 から参考にしていくといいでしょう。

4202 インストールにあたっては、Stan がコンパイルされるときに利用する C++ 言語環境が必要です。これは
4203 MacOS であれば Xcode の導入が、Windows であれば Rtools での導入が必要になります。

4204 ■MacOS ユーザの場合 AppStore で Xcode と検索し、Xcode をインストールすれば OK です。

4205 ■Windows ユーザの場合 Rtools という R 周辺のツールをインストールする必要があります。 <https://cran.r-project.org/bin/windows/Rtools/> に
4206 いて、自分の R のバージョンにあった RTools を
4207 インストールしてください。この他にも Windows ユーザは導入にあたって色々障壁にあたる場合があります
4208 ので、こちらのサイト <https://norimune.net/3609> なども参考にしてみてください。

4209 18.3.2 導入方法 2 ; 仮想環境での構築

4210 第二のルート、仮想環境を構築する場合ですが、これについては我らが国里先生が大変詳しいサイトを
4211 作ってくださっています。日本心理学会のチュートリアルワークショップ、「再現可能な日本語論文執筆入門：
4212 jpaRmd で実現する再現可能で低コストな日本語論文執筆のはじめの一步」で使われた時の資料集がそ
4213 れで、このサイト <https://ykunisato.github.io/jpa2021-tws-jpaRmd/> の事前準備編をよく読むと
4214 導入できます。通信環境にもよりますが、慣れているとものの 10 数秒でシステム導入できるという優れたもの
4215 です。面倒な設定が怖い人は、是非試してみてください。

4216 18.4 導入方法 3 ; 外部サーバの利用

4217 最近、Google 社がオンラインで利用できる分析環境、Google Colaboratory を提供してくれています。
4218 これは Google 社が提供する、ブラウザで実行できるプログラミング環境です。Google が教育、研究用に提
4219 供しているもので、90 分間は無料で利用できます。90 分を超える場合でも、書いたプログラムを保存して

*17 インストールに際しては、RStudio Desktop の Free edition を選んでください。RStudio Cloud や RStudio Server は別のものです。

4220 おけば、新しいセッションを開始することで続けて利用できます。これを用いる利点は、どのユーザにも同じ環
4221 境を提供できることです。

4222 基本言語環境は Python で提供されており、Jupyter Notebook を使いますが、次のアドレス (に含まれ
4223 るコード) を使えば言語環境を R に変えて利用することができます。

4224 <https://colab.research.google.com/notebook#create=true&language=r>

4225 ここでウィンドウに R のコードを入力し、実行していきます。cmdstanr や rstan もインストールできま
4226 す。インストールに際して、install.packages("rstan") のように入力することもできますが、この方法
4227 ですとインストールに随分と時間がかかります。そこで少し例外的ですが、システムコマンドを直接利用して、
4228 system("apt install -y r-cran-rstan") のようにすることで、大幅にスピードを短縮することがで
4229 きます。

4230 18.4.1 導入後の確認

4231 インストールが終わったらサンプルコードを実行して「動くかどうか」を確かめてみてください。サンプルコー
4232 ドは Getting started with CmdStanR のサイト、あるいは RStan Getting Started のサイトに掲載され
4233 ています。どちらのパッケージから実行しても、文字列がズラズラっと出力されればとりあえず動いたと思っ
4234 てください。赤い文字やエラーなどが表示される、あるいは全く表示されない場合は、インストー
4235 ルがうまくいっていないことを疑いましょう。上に戻って、丁寧に説明に目を通して一步一步進めてください。
4236 機械に命令することですので、自分勝手にショートカットしたり変更したりしないことが重要です。

4237 18.5 Stan を使ってみよう

4238 具体的な Stan コードの書き方や統計モデルへの応用については、本書の次の章から順に説明していきま
4239 す。それに先立って、全体的な注意をしておきたいと思います。具体的には、バージョンによる違い、パッケ
4240 ージによる違い、そして本書の以下の章で使う関数の紹介です。

4241 18.5.1 バージョンによる書き方の違い

4242 cmdstanr と rstan はパッケージの使い方以上に、それらが呼び出す Stan のバージョンの違いがありま
4243 す。上で解説したように、cmdstanr のほうが最近が開発が進んでいて、より新しい Stan の機能を使ってい
4244 ることとなります。バージョンが違ってても、基本的な書き方などに違いはないのですが、バージョン 2.26 から
4245 Stan の文法に変更が予定され、導入されました。

4246 データ X が N 人分あったとします。数学的フォームでは、 $X_1, X_2, X_3, \dots, X_N$ というようなものです
4247 ね。これはコンピュータ上では X と名付けられた配列、サイズ N というように考えます。これは、Stan では次
4248 のように書くものでした。

code : 18.1 Stan2.26 以前の書き方

```
4249 1 int n[5];
4250 2 real a[3, 4];
4251 3 real<lower=0> z[5, 4, 2];
4252 4 vector[7] mu[3];
4253
4254
```

4255 これらは上から順に、サイズ 5 の整数変数 n 、サイズ 3×4 の実数変数 a 、下限 0 のサイズ $5 \rightarrow, es4 \times 2$

4256 の配列 z , 長さ 7 のベクトル μ が 3 つ, ということを意味しています*18。

4257 この書き方が, Stan のバージョン 3.32 以降は次のような書き方になります。

code : 18.2 Stan2.32 以降の書き方

```
4258 1 array[5] int n;
4259 2 array[3, 4] real a;
4260 3 array[5, 4, 2] real<lower=0> z;
4261 4 array[3] vector[7] mu;
4262
4263
```

4264 このコード 18.2 はさきほどのコード 18.1 と同じ意味で, 書き方が変わっただけです。一瞥してわかるように,
4265 サイズを先に array という言葉で宣言しましょう, というだけです。

4266 さて, 今現在は Stan2.26 と 3.32 の間, 過渡期になります。ですから, どちらのコードで書いてもエラーに
4267 はなりません。ただし, (ここからが少し面倒なのですが)

- 4268 • RStudio のコードチェック機能は新しいコードの書き方に対応しておらず*19, 新しい書き方をす
4269 るとエディタ上でエラーの警告 (赤いバツェンがつきます) が出ますし, チェックボタンを押しても
4270 SYNTAX ERROR が出ます。
- 4271 • rstan パッケージは使っている Stan が古いこともあって, 新しい書き方のコードをコンパイルしよう
4272 とするとエラーになります。RStudio と同じ, SYNTAX ERROR が出ます。
- 4273 • cmdstanr パッケージは逆に, 新しい Stan の書き方を推奨していますので, 「その書き方は古いよ」と
4274 という警告を出してきます。Declaration of arrays by placing brackets after a variable name is
4275 deprecated and will be removed in Stan 2.32.0. Instead use the array keyword before the
4276 type. This can be changed automatically using the auto-format flag to stanc. という警告が
4277 できますが, これは「配列のカッコを変数の後ろに置く書き方はもうダメで, Stan2.32 以降は無くなりま
4278 す。型の前に array と書くようにしてください。これは auto-format フラグを立てることで自動的に変
4279 更するようになります。」という意味です。

4280 つまり, 本書は cmdstanr を推奨していますが, そうすると RStudio のエディタ上ではエラーが出て文法
4281 チェックができず, コンパイルした後でいざサンプリング, というときにエラーが出たりします。チェックをするに
4282 は, コンパイルしたオブジェクトを使ってサンプリングをする前に, model\$check_syntax() 関数を実行す
4283 る必要があります。rstan パッケージを使うとスマートに文法チェックもできていいのですが, rstan そのも
4284 のが不安定ですし, 今後徐々に開発・発展をしなくなっていくことが明らかですので, どこかで新しい方に舵
4285 を切る必要があります。

4286 本書のコードは新しい書き方に統一していますが, 最初のうちは, そして今しばらくは, 実際にコードを書く
4287 ときは古い方で書いたほうが便利かもしれません。

4288 18.5.2 パッケージによる指示と出力の違い

4289 さて今度はパッケージごとの違いを見ていきましょう。

*18 ベクトルは数字をセット, まとめてとして演算するものです。最後の変数は, 7 つの数字のセットが 3 つという意味で, 7 つセットで 1 つのまとまりですよ, ということを明示的に宣言していることになります。

*19 RStudio のバージョン 2022.07.2 Build 576 で確認しました。

4290 指示の仕方の違い

4291 まずは実行方法です。Stan では事後分布からの乱数を生成しますが、それにあたって「乱数の数」「ウォー
4292 ムアップ期間の長さ」「チェーンの数」などオプションに指定できるものがあります。これは見てもらったほう
4293 が早いかと思えますので、同じことをそれぞれのパッケージで実行するときどのように指定方法を変えるか
4294 を見てみましょう。まずは `rstan` パッケージの書き方からです。

code : 18.3 rstan のスタイル

```
4295 1 fit <- rstan::sampling(model,
4296 2   data = dataSet,
4297 3   chains = 4,
4298 4   iter = 6000,
4299 5   warmup = 1000
4300 6 )
4301
4302
```

4303 ここで指定しているのは、データを `dataSet` として設定し、同時に 4 本のチェーンを発生させ、6000 個の
4304 乱数を作って、そのうち 1000 個はまだ機械が温まっていないので候補から除外する、というものです。複数
4305 のチェーンを作ること、温まっていない部分を捨てる理由などは後ほど説明しますが、結果的に合計 20000
4306 個の乱数が作られることを知っておいてください。この数字は $(6000 - 1000) \times 4 = 20000$ という計算か
4307 らでできます。また、乱数の発生は並列計算させた方が良いので、`rstan` パッケージを使う場合は事前に、
4308 `options(mc.cores = parallel::detectCores())` という一行を入れておきます。

4309 さて、同じことを `cmdstanr` パッケージでやると次のようになります。

code : 18.4 cmdstanr のスタイル

```
4310 1 fit <- model$sample(
4311 2   data = dataSet,
4312 3   chains = 4,
4313 4   parallel_chains = 4
4314 5   iter_sampling = 5000,
4315 6   iter_warmup = 1000,
4316 7 )
4317
4318
```

4319 ここにあるように、並列するチェーンの数をオプションに書き込む必要があります。またサンプルの数は捨てる
4320 部分を別にして (含めずに) 計算させることができます。このコードで 20000 個のサンプルが得られます。細
4321 かいことですが、多少の違いがあるわけです。

4322 出力の違い

4323 さて、今度は出力結果の使い方についての注意です。`rstan` パッケージは、出力結果を `stanfit` オブ
4324 ジェクト型という形にします。`rstan` のもっているさまざまな関数、たとえば要約した結果の出力やプロットな
4325 どは、`stanfit` オブジェクトだとわかるとそれに対応した出力に変えて表現してくれるわけです。`cmdstanr`
4326 の出力はまた別物^{*20}です。

^{*20} `class` 関数で `rstan` の出力がどういうクラスなのかを表示させると、`stanfit` と出てくるのですが、`cmdstanr` の出力を同様にチェックすると `"CmdStanMCMC"` `"CmdStanFit"` `"R6"` という答えが返ってきます。クラスというのはデータの形を定義する方法で、オブジェクト指向プログラミングに必要な概念なのですが、ここではややこしすぎるのでパス。違うということだけわかっ
てもらえれば結構です。RStudio のばあいは Environment タブを見るだけでも違いがわかると思います。`rstan` パッケージの出力はちゃんとしたクラス (Formla Class `stanmodel`) なのですが、`cmdstanr` パッケージの出力は Environment (グローバル環境の要素) として剥き出しの出力が出ている感じがします。

4327 同じモデルをそれぞれのパッケージで出力させて、違いを見てみましょう。まずは `rstan` パッケージの出力
4328 から。

rstan の出力 1: rstan パッケージの出力

```
Inference for Stan model: tttest01.
4 chains, each with iter=6000; warmup=1000; thin=1;
post-warmup draws per chain=5000, total post-warmup draws=20000.

      mean se_mean  sd  2.5%  25%  50%  75%  97.5% n_eff Rhat
mu1   59.96   0.04 5.58  49.01  56.33  59.99  63.60  70.91 16043  1
mu2   40.04   0.05 5.62  28.98  36.38  40.03  43.67  51.17 14330  1
sig   15.54   0.03 3.12  10.86  13.35  15.08  17.23  22.92 12118  1
lp__ -51.58   0.02 1.36 -55.12 -52.18 -51.24 -50.60 -50.05  7115  1

Samples were drawn using NUTS(diag_e) at Tue Nov 15 09:50:17 2022.
For each parameter, n_eff a crude measure of effective sample size,
and Rhat is the potential scale reduction factor on split chains (at
convergence, Rhat=1).
```

4329

4330 `rstan` パッケージが出力しているのは次のような情報です。

- 4331 • モデル名, チェイン数, 反復回数, 最終的に得られた MCMC サンプル数 (draws) などの説明
- 4332 • パラメータの事後分布の記述統計量, すなわちベイズ推定による事後分布の情報。順に平均値
4333 (mean), 平均値の標準誤差 (se_mean), 標準偏差 (sd), 2.5%,25%,50%,75%,97.5% パーセンタ
4334 イル。
- 4335 • MCMC サンプルについての特徴量。有効サンプルサイズ (n_eff) と Rhat(Rhat)
- 4336 • サンプリングに関するあとがき

4337 それぞれの内容についてはあとの章で説明するとして、続けて `cmdstanr` パッケージの出力をみてみます。

cmdstanr の出力 1: cmdstanr パッケージの出力

```
variable  mean median  sd  mad    q5    q95  rhat  ess_bulk  ess_tail
lp__ -51.56 -51.20 1.34 1.07 -54.18 -50.12 1.00    7799    10071
mu1   59.99  60.00 5.63 5.47  50.82  69.16 1.00   16325   12764
mu2   40.05  40.05 5.48 5.18  31.09  49.05 1.00   16511   12466
sig   15.48  15.03 3.07 2.75  11.41  21.07 1.00   14489   12137
```

4338

4339 同じような情報ですが、ちょっと違うところもあるようです。前から順に、事後分布の平均値 (mean), 中
4340 央値 (median), 標準偏差 (sd), 平均絶対偏差^{*21} (mad), 5%,95% パーセンタイルなど分布の特徴^{*22}と、
4341 Rhat や有効サンプルサイズに関する情報 (ess_bulk, ess_tail) が示されています。しかし、まえがき・あ
4342 とがきのようなものがなく、結果だけがシンプルに出力されていますね。

^{*21} 中央値絶対偏差 (Median Absolute Deviation) とは、中央力のばらつきを表す指標で、各データ点の中央値からの偏差の絶対値をとり、その中央値を計算したものです。標準偏差とは違う分布の幅を表現する方法で、中央値を使って計算しますから外れ値に強いという特徴があります。

^{*22} `rstan` パッケージは 2.5% から 97.5% のパーセンタイル, すなわち囲まれる区間の面積が 95% です。`cmdstanr` パッケージは 5% と 95% のパーセンタイル, すなわち囲まれる区間の面積が 90% です。なぜだかわかりませんが、ちょっとした違いがありますね。

出力としては `rstan` パッケージの方が丁寧で良いかもしれません。さて、`cmdstanr` の出力は実は `csv` ファイルとして一時フォルダに保存されているので、これを `rstan` パッケージの関数を使って `stanfit` オブジェクトに変換することができます。関数の使い方の例は次のようになります。

code : 18.5 `cmdstanr` の出力を `rstan` のオブジェクトに変える

```
4346 1 fitR2 <- fitC$output_files() |> rstan::read_stan_csv()
4347
4348
```

コードの中身ですが、まず `fitC` と書いてあるのが `cmdstanr` の出力オブジェクトです。このオブジェクトがさまざまな情報を持っているという形になっており、ここから出力ファイルの場所を聞き出すのが `output_files()` です。このファイルを `rstan` パッケージの `read_stan_csv()` 関数に渡してやると、`csv` ファイルを読み込んで `rstan` のオブジェクトに変えてくれます^{*23}。変えたものをここでは、`fitR2` オブジェクトに保存しています。この方法を使うと、先ほどの `cmdstanr` の出力が `rstan` の出力と同じになります (同じなので再掲しません)。

このコードを覚えておけば、`rstan` パッケージ向けに書かれたものであっても、`cmdstanr` パッケージで実行したあとと同じ関数を適用できるので、使いやすいですね。

18.5.3 本書で利用する準備関数

どちらのパッケージであっても、欲しいものは事後分布からの乱数をデータセットにしたものであり、そのデータセットさえ持っていれば、既存の関数を使わなくとも自分で工夫して加工すれば良いでしょう。

そこで本書では、MCMC サンプルを抜きだしてデータセットに変換する関数を作り、これを利用して話を進めることにします。ここではその関数の中身を解説しておきます。

MCMC サンプルをデータフレームにする関数

まずは MCMC サンプルをデータセットにして渡してくれる関数です。この関数は伴走サイトのコードの冒頭に必ず含まれており、以後はこれを使って結果の要約を示すといった使い方をしていきますので、中身を知っておいてもらった方が良いでしょう。自分の使っているパッケージの方だけでも目を通してください。この関数を経由すると、出てくるオブジェクトは同じ形式になっているところがミソです。

■`stanfit` オブジェクト (`rstan` パッケージ) の場合 `rstan` パッケージを使って `stanfit` オブジェクトとして MCMC サンプルを得た場合、そのオブジェクトから乱数部分だけを抜き出す関数は `rstan::extract` 関数です。これで取り出したものをデータフレーム型にして返す関数として、次のようなものを作りました。

code : 18.6 MCMCtoDF 関数 (`rstan` 版)

```
4370 1 MCMCtoDF <- function(fit) {
4371 2   fit %>%
4372 3     rstan::extract() %>%
4373 4     as.data.frame() %>%
4374 5     tibble::as_tibble() %>%
4375 6     tibble::rowid_to_column("iter") %>%
4376 7     dplyr::select(-lp_) %>%
4377 8     tidyr::pivot_longer(-iter) -> MCMCsample
4378 9   return(MCMCsample)
4379
```

^{*23} 当然のことながら、`rstan` パッケージがないとこの関数も呼び出せませんので、こちらもインストールしておく必要があります。なおこのコードの、`|>` は `tidyverse` パッケージのパイプ演算子、`%>%` と同じ機能を持つ R の演算子で、ネイティブパイプと呼ばれています。ネイティブパイプは R4.1.0 以降に導入されたものです。

```
4380 10 }
4381
```

4382 ■コード解説

- 4383 1 行目 関数を作る宣言。引数として `stanfit` オブジェクトをとります。
- 4384 2 行目 引数で引き受けたオブジェクトを加工していきます。
- 4385 3 行目 まずは MCMC サンプルを抜き出します。
- 4386 4 行目 `data.frame` 型にします。
- 4387 5 行目 `tibble` 型にします。別にしなくてもいいんですが、著者の好みです。
- 4388 6 行目 行番号を表す変数 `iter` を作ります。
- 4389 7 行目 変数の中から `lp__` を除外します。
- 4390 8 行目 行番号変数は残して、`tidy` なデータにし、`MCMCsample` というオブジェクトに保存します
- 4391 9 行目 戻り値として `MCMCsample` を返します。

- 4392 ■`cmdstanr` パッケージの場合 `cmdstanr` パッケージを使って MCMC サンプルを得た場合、そのオブ
- 4393 ジェクトから乱数部分だけを抜き出す関数は `draws` 関数であり、これをオブジェクトに直接作用させます。こ
- 4394 れで取り出したものをデータフレーム型にして返す関数として、次のようなものを作りました。

code : 18.7 MCMCtoDF 関数 (`cmdstanr` 版)

```
4395
4396 1 MCMCtoDF <- function(fit) {
4397 2   fit$draws() %>%
4398 3     posterior::as_draws_df() %>%
4399 4     tibble::as_tibble() %>%
4400 5     dplyr::select(-lp__, -.draw, -.chain, -.iteration) %>%
4401 6     tibble::rowid_to_column("iter") %>%
4402 7     tidyr::pivot_longer(-iter) -> MCMCsample
4403 8   return(MCMCsample)
4404 9 }
4405
```

4406 ■コード解説

- 4407 1 行目 関数を作る宣言。引数として `stanfit` オブジェクトをとります。
- 4408 2 行目 引数で引き受けたオブジェクトに `draws` 関数を作用させ、サンプルだけ取り出します。
- 4409 3 行目 取り出したサンプルをデータフレーム型にする `posterior` パッケージの `as_draws_df` 関数を適
- 4410 用します。
- 4411 4 行目 `tibble` 型にします。別にしなくてもいいんですが、著者の好みです。
- 4412 5 行目 変数の中から `lp__` や他の隠し変数を除外します。
- 4413 6 行目 行番号を表す変数 `iter` を作ります。
- 4414 7 行目 行番号変数は残して、`tidy` なデータにし、`MCMCsample` というオブジェクトに保存します
- 4415 8 行目 戻り値として `MCMCsample` を返します。

- 4416 ■`MCMCtoDF` 関数の結果 出力結果は、どちらも同じく以下のような形式になります。

R の出力 18.1: MCMC サンプルをデータセットにしたもの

```
# A tibble: 60,000 × 3
  iter name  value
  <int> <chr> <dbl>
1     1 mu1   60.3
2     1 mu2   41.0
3     1 sig   12.5
4     2 mu1   62.9
5     2 mu2   37.6
6     2 sig   11.5
7     3 mu1   68.7
8     3 mu2   46.4
9     3 sig   13.7
10    4 mu1   67.7
# ... with 59,990 more rows
# [X] Use `print(n = ...)` to see more rows
```

4417

4418 ここで iter とあるのは MCMC のステップ番号, name とあるのが変数名, value とあるのがその値で
4419 す。tidy な形になっていますから, このまま分析や描画に用いることができます。

4420 MCMC サンプルの結果を要約する関数

4421 MCMC の結果を分析するにあたっては, bayestestR パッケージなど既存のパッケージを使うと便利で
4422 しょう。ですが, こうしたパッケージはどんどん開発が進んでアップデートされたり, rstan か cmdstanr かで
4423 書き方が変わったりするなど, テキストで紹介するには難しいところがあります。そもそも MCMC サンプルを
4424 加工して記述統計を示したり, 描画したりするわけですから, 上で紹介したように MCMC サンプルを抜き出
4425 すことができれば自力でもできるはず。ということで, 自作の関数を用意しました。本書では以下, この関数を
4426 つかって解説しますので中身を知っておくと良いでしょう。

4427 ■MAP 推定関数 確率分布の特徴を報告するときに, その確率分布に従う乱数を使って, 期待値や中央値
4428 を計算するのは簡単です。記述統計としての平均値やパーセンタイルでよいからです。しかし確率分布の確
4429 率密度が最も高くなる場所, すなわち MAP 推定値 (MAP Estimation) の計算はちょっと難しいです
4430 ね。というのも, 関数であれば微分して極値を求めれば良いのですが, MCMC サンプルは関数ではなくて
4431 値なので, ヒストグラムを書くしかなく, 最大の値を求める計算がむずかしいからです。

4432 でも大丈夫, R の描画関数を駆使して対応することができます。MAP 推定値を計算する関数は, 次のよ
4433 うに書きます*24

code : 18.8 MAP 推定する関数のコード

4434

```
4435 1 map_estimation <- function(z) {
4436 2   density(z)$x[which.max(density(z)$y)]
4437 3 }
```

4438

4439 これは R の density 関数でヒストグラムに密度関数を当てがい (カーネル密度推定), その関数の特徴か
4440 ら値を返すようになっています。

*24 このコードは関西学院大学社会学部の清水裕士先生に教えてもらったものです。記して謝意を表します。

4441 ■MCMC サンプルの要約関数 先ほどの MAP 推定する関数を含めつつ、また MCMC サンプルをデー
4442 タフレームにする関数を使いつつ、MCMC サンプルの要約を表示する関数を次のように作りました。

code : 18.9 MCMC の要約を報告するコード

```
4443 1 MCMCsummary <- function(MCMCsample) {
4444 2   MCMCsample %>%
4445 3   dplyr::group_by(name) %>%
4446 4   dplyr::summarise(
4447 5     EAP = mean(value),
4448 6     MED = median(value),
4449 7     MAP = map_estimation(value),
4450 8     SD = sd(value),
4451 9     U95 = quantile(value, prob = 0.975),
4452 10    L95 = quantile(value, prob = 0.025)
4453 11   ) %>%
4454 12   mutate(across(where(is.numeric), ~ num(., digits = 3)))
4455 13 }
```

4458 ■コード解説

4459 1 行目 関数を作る宣言。引数として MCMCtoDF 関数の戻り値をとります。
4460 2 行目 引数で引き受けたオブジェクトを加工していきます。
4461 3 行目 変数名でデータをグループ化します。
4462 4-10 行目 記述統計の計算を通じて、EAP,MED,MAP,SD,95% 区間を返します。
4463 11 行目 出力する数値データは、少数下 3 桁まで表示させるようにします。

4464 この関数は、先ほどの MCMCtoDF とセットで使い、次のような結果を得ます。

R の出力 18.2: MCMCsummary の結果

```
> fit %>% MCMCtoDF() %>% MCMCsummary()
# A tibble: 3 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu1    59.897    59.888    58.867    5.547    48.672    70.847
2 mu2    39.958    39.935    39.748    5.523    28.973    50.943
3 sig    15.519    15.062    14.646    3.074    10.877    22.788
```

4466 ここにあるように、各変数の EAP 推定値、MED 推定値、MAP 推定値、95% 確信区間が表示され
4467 ます*25。

*25 この区間はコードから明らかなように、`quantile` 関数で計算しています。下から 2.5%、上から 2.5% を取り除くと 95% の区間が残ることになりますが、このように分布の両端を均等に切った区間を厳密には**等裾区間 (Equal-tailed interval;ETI)** といふのでした。これに対して、分布の最も信頼できる部分、かつ分布の大部分をカバーし、区間内部のすべての点が、区間外の任意の点よりも高い密度を持っているように推定するものを**最高密度区間 (Highest-Density Interval;HDI)** と呼びます。これらは左右対称の単峰分布であれば同じになるのですが、事後分布の形がぐちゃっとしている時は必ずしもそうなりません。両者の違いについては [Kruschke \(2014\)](#) を参考にしてください。また、HDI を出力する関数は `bayestestR::hdi()` です。

4468 これらの推定値や, そもそも事後分布が何をどのようなことを表しているのか, といった点については今後
4469 の授業で説明していきます。以下は結果をしれっとこの関数で表現していきますので, 何をやっているのかわ
4470 からなくなったらここに立ち戻ってきてください。

第 19 章

ベイジアンアプローチと確率的プログラミング 1

それではいよいよ、データ生成モデリングとベイズ推定による実践例を見ていくことにしましょう。

19.1 7 人の科学者

ここでは Lee and Wagenmakers (2013) より「7 人の科学者」の例を紹介します。そこでのカバーストーリーは次のようなものです。

実験スキルが大きく異なる 7 人の科学者が、全員同じ量について測定を行う。彼らが得た数値結果は次の通り。

$$Y = \{-27.020, 3.570, 8.191, 9.808, 9.603, 9.945, 10.056\}$$

直感的には、最初の二人の科学者はひどく適性を欠いた測定者であり、この量の真の値はおそらく 10 をわずかに下回るくらいであるように思えるのだが・・・？

さあ、これを見て「あれ、どこが統計的な問題なんだ？」と思った人もいるかもしれません。なんらかの測定をして、データのばらつきがあるんだけど、それが測定者によって違うらしい、というのはわかったかと思いますが、これをどうやって統計的な問いにするのでしょうか。まず、ストーリーの背後に「正確な測定値 (真の値) はわからないけど、定数のはず」という前提があることを確認しましょう。測定の基本として測定誤差があるので真の値を特定出来はしませんが、測定を繰り返すと真の値に近づいていくはずではあります。また、今回の測定結果は 7 人それぞれ別の人による測定であることから、測定誤差の出方が測定毎すなわち測定者毎に異なるだろう、という仮説もあります。

これらを踏まえて謎解きをしていきましょう。わからない量を確率の言葉で表現し、データを使って少しでも正解に近づこうとするのが統計のやり方です。まずここでわからない量は、「正確な測定値」と「個人毎の誤差の大きさ」です。誤差は正規分布に従うと思われるから、データ Y は $Y \sim N(\mu, \sigma)$ と表すことができます。データ生成メカニズムという意味では、「正規分布に従ってデータが作られる」と言ってもいいかもしれません。また、正確な測定値はここでは μ ということになります^{*1}。また、データは 7 点あって、それぞれ Y_1, Y_2, \dots, Y_7 , あるいは一般に Y_i と書くことにしましょう。ここで添字の i は第 i 回目の測定でもありますし、測定者 i のことでもあります。この測定者 i 毎に誤差の出方の大きさが変わります。誤差の大小、言い換

^{*1} 正確な測定値 μ に誤差 $N(0, \sigma)$ がついてデータになる、すなわち $Y = \mu + N(0, \sigma)$ ということですが、誤差の平均は 0 で μ の分だけ正規分布がズレる、そして「分布に従う」を表現する \sim を使いたいの、 $Y \sim N(\mu, \sigma)$ となるのです。

4494 えれば測定の精度は誤差の SD で表現されていますから、 σ が個人毎に変わる、すなわち σ_i と書くことがで
4495 きます。これが今回のデータ生成モデルです！

4496 ここまでのことをイメージ図であらわすと図 19.1 のようになります。この図の書き方のポイントは、まずデー
4497 タを一番下に置くことです。次にこのデータがどこから来たのかな、ということを示す矢印で表現して書きます。
4498 データは誤差を伴うなどで、毎回定数にならないでしょうから、**確率分布**をその上に書くことになります。確率
4499 分布にはその形状を定めるパラメータがあるはずですから、それを書き加えれば OK です。

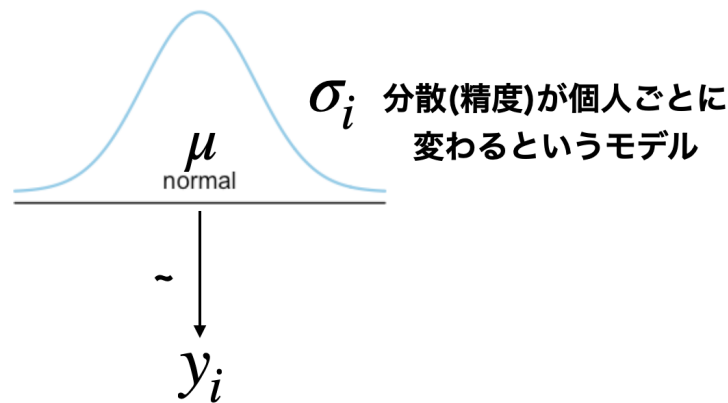


図 19.1 データが生成されるプロセス

4500 しかしここでも、 μ と σ_i は結局なんなんだ、どういう数字なんだということがわからないままです。わから
4501 ない量は確率の言葉で表現するのがベイジアン生き様。ここは未知のパラメータがどうやって出てくるかな
4502 なんてわからないところなので、「わからない」ことを確率で表現する必要があります。パラメータに対して「この
4503 辺りにあるだろう」と事前に想定する確率分布のことを**事前分布 (prior distribution)** ということでした。

4504 ここでは μ は平均 0、SD100 の正規分布からきていることにしましょう。正規分布というのは左右対称で
4505 平均値・中央値・最頻値が一致する単峰の分布です。これは正確な測定値を表しているのでしたから、なんら
4506 かの値を取るとしても 1 つの値だろう、というのは無理のない仮定でしょう。複数の値を取る可能性があるの
4507 なら多峰性の分布を仮定したらいいと思いますが、今回はそうではないだろう、と仮定するのです。平均 0 で
4508 SD100 ということは、-300 から +300 の範囲にある可能性が 99.7% だということでもあります*2。実際の
4509 データが -27 から 10 ぐらいの範囲に入っているの、-300 から +300 の間というのも十分過ぎるぐらい
4510 広めの範囲にしてある数字と言えるでしょう。このように、うっすら情報を持っているけどほとんど意味がない
4511 程度に縛りをかけるのを、**弱情報事前分布 (weakly informative prior distribution)** と言います。
4512 たとえばこれが身長データであれば、 $\pm 3m$ の範囲は十分すぎ、なんなら負の数を取るはずがないので過
4513 剰な安全策を置いているとも言えます。それでも「なんらかの事前分布を持つのは、主観的な思い込みだ」と
4514 いうのであれば、 $\pm\infty$ の一様分布を考えるか、ベイズ法をやめるかということになります。

4515 次に σ_i です。これは標準偏差パラメータなので負の数を取ることはあり得ません。ここでは**半コーシー**
4516 **分布 (half-cauchy distribution)** を事前分布におくことにします。コーシー分布とは正規分布によくに

*2 正規分布は、平均 $\pm 1SD$ の範囲に全体の 68% が、 $\pm 2SD$ の範囲に 95% が、 $\pm 3SD$ の範囲に 99.7% が入ることが理論的にわかっているのです。わからない人は R で計算して確認すること！

4517 た形ですが、正規分布より裾が重く*3、標準偏差（分散）の事前分布に適していると言われている分布です
 4518 (Gelman et al., 2006)*4。このコーシー分布は正規分布のように左右対称ですが、0 を中心に折り
 4519 たたんで正の値しか取らないように考えたのが半コーシー分布です。今回はスケールパラメータに 5 を置い
 4520 ています。とくにこの数字に深い意味はありません。気になるようでしたら 10 でも 100 でも変えていただ
 4521 いて結構です*5。

4522 これらをおくことで、モデルの設計図ができあがりました (図 19.2)。あとはこれをもとに、**確率的プログラ
 ミング言語 (stochastic programming language)** に書き起こしていけば良いでしょう。

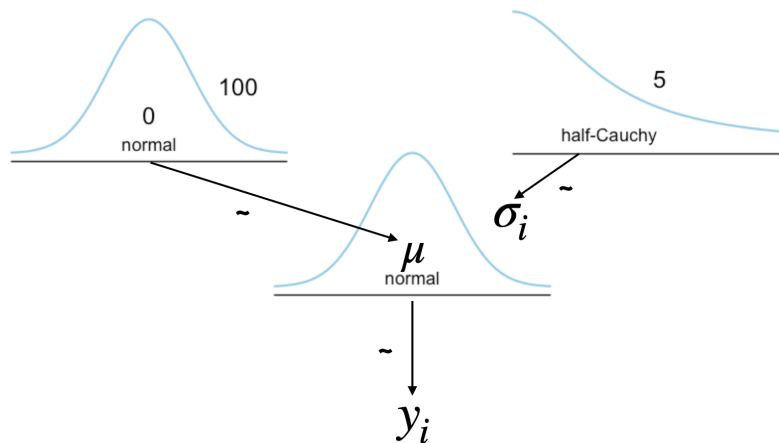


図 19.2 モデルの全体像

4523

19.2 Stan コードの書き方

4524

4525 では図 19.2 という設計図に基づいて、Stan のコードを書いていきましょう。

19.2.1 Stan コードの文法

4526

4527 Stan は事後分布からの乱数を生成するための言語です。その言語仕様は C や C++ にていているところ
 4528 があります。R で書くのと違うところとして「ブロック単位の記述」「セミコロンによる行の区切り」「変数の宣
 4529 言」の 3 つを押さえておきましょう。

4530 ■**ブロック単位の記述** Stan に書くべきことは**ブロック**の中に記述します。ブロックは中括弧 (`{ ... }`) で
 4531 括ったもので、次の 6 つの種類があります。

- 4532 1. data ブロック
- 4533 2. transformed data ブロック
- 4534 3. parameters ブロック

*3 より極端な値が出る可能性が多いという意味です

*4 この論文では他にも、student の t 分布などを提案していますが、パラメータ数がより少ないコーシー分布を選びました。

*5 半コーシー分布がどんな数字になるのか気になる、という人は R で `rcauchy` 関数を使って乱数を発生させ、その挙動を見てみれば良いでしょう。乱数発生によって大体の分布感 (分布の感覚?) を得ることができる、ということをお出ししてください。

- 4535 4. transformed parameters ブロック
- 4536 5. model ブロック
- 4537 6. generated quantities ブロック

4538 Stan のプログラムに含まれるブロックは、この順番でなければならぬと決まっているので注意してくださ
4539 い。この中で、1,3,5 番目、すなわち data, parameters, model ブロックが基本的に必要なもので、2.4.6
4540 番目のブロックはその派生や応用です。まずは基本ブロックから説明します。

4541 data ブロックは Stan と外部とのやりとりをするブロックです。外部からデータを取り込んで、それに基づ
4542 いて乱数を発生させます。Stan で計算せずに外部から与えるものはすべてこのブロックの中で書く必要が
4543 あります。parameters ブロックは、パラメータ、すなわち Stan で推定したいものを書くところです。Stan
4544 は事後分布からの乱数発生機であり、事後分布は**確率分布**ですからその形状はパラメータによって特徴づ
4545 けられます。事後分布の形が分からないというのは、事後分布のパラメータがどのような値になっているのか
4546 分からない、ということと同じ意味です。その「分からない」「知りたい」パラメータをリストアップしておくブロッ
4547 クです。model ブロックは、確率モデルすなわち尤度と事前分布を書くところです。ここで書かれた尤度と事
4548 前分布が、data ブロックのデータと組み合わせると、事後分布が Stan 内部で計算され、そこからの乱数が
4549 次々取り出されてくることになるのです。

4550 ■**セミコロンによる行の区切り** このブロックの中に、Stan 言語の文を書いていくことになるのですが、
4551 ここで 1 つ覚えておいて欲しいのが、**Stan の句点はセミコロン**ということです。句点がない文章は読みにく
4552 いどころか、機械にとっては意味不明な文字列になりかねませんので、必ず 1 つの命令・一文がおわった
4553 らセミコロンで閉じるようにしてください。これを忘れてエラーになるのが初心者にもっとも多いミスのひとつ
4554 です。

4555 ■**変数の宣言** たとえば R では数字でも分析結果でも、オブジェクトに代入でき、オブジェクト名は任意に
4556 つけることができました。またオブジェクトはいつどのタイミングで発生させても OK で、思いついたらすぐオ
4557 ブジェクトに代入、ということができました。Stan をはじめ、いくつかの高級言語ではこれに厳格な制限をか
4558 けるものがあります。すなわち、どんなオブジェクトを作るかを最初に**宣言 (declaration)** する必要がある
4559 のです。Stan では使われるオブジェクト (変数) が、どういう数字になるものなのか、どの範囲の数字になる
4560 のかをブロックの最初に宣言します。この数字の種類のことを**型 (type)** といい、その宣言された数字の型に
4561 適合しない数字はエラーとなって弾かれます。なんでそんな窮屈な、と思うかもしれませんが、逆に「ありえな
4562 い数字は許さない」という意味で、エラーを未然に防いでくれることにもなるのです。

4563 変数の形として、Stan は整数 (**int 型**)、実数 (**real 型**)、ベクトル (**vector 型**)、行列 (**matrix 型**) など
4564 があり、また範囲を lower や upper で指定します。たとえば標準偏差は負の数を取らないので、標準偏差
4565 の変数を宣言したいときは、real<lower=0> sig; のように書きます。

4566 19.2.2 Stan コードを書いてみよう

4567 されこれらを踏まえてコードを書いてみるわけですが、設計図を書いたときにデータを下に、それに紐づく
4568 確率分布を上の上に、と下から上に書いていったように、model ブロック → parameters ブロック → data
4569 ブロック、と書き進んでいったほうがわかりやすいでしょう。

4570 Stan のコードはメモ帳など、エディタで開くことのできるテキストファイルでいいのですが、拡張子を .stan
4571 にしておきましょう。RStudio で書く場合は、.stan にしておくことで Stan のファイルだと認識してくれるので、
4572 ブロックや型宣言の用語などを強調表示してくれるようになります。RStudio でコードを書き始めるときは、

4573 ファイル (File) >新しいファイル (New File) と進むと Stan というファイルがあります。これを選ぶとすでに
 4574 Stan のサンプルコードが書かれているファイルが出てきますが、そのファイルの中身は全部消して、次のコード
 4575 19.1 のように書き、ファイル名をつけて保存しましょう。ここでは `sevenScientist.stan` とファイル名を
 4576 つけたものとします。

code : 19.1 7人の科学者コード

```

4577 1 data{
4578 2   array[7] real Y;
4580 3 }
4581 4
4582 5 parameters{
4583 6   real mu;
4584 7   array[7] real<lower=0> sig;
4585 8 }
4586 9
4587 10 model{
4588 11   for(i in 1:7){
4589 12     //likelihood
4590 13     Y[i] ~ normal(mu, sig[i]);
4591 14     //prior
4592 15     sig[i] ~ cauchy(0,5);
4593 16   }
4594 17   //prior
4595 18   mu ~ normal(0,100);
4596 19 }
4597

```

4598 ■コード解説

4599 **model ブロック** 尤度と事前分布を書くところ。確率分布からのサンプリングは変数を左に、`~`で**確率分**
 4600 **布**につなぎます。`normal`とか`cauchy`が確率分布の名前なので RStudio ではハイライト表示され
 4601 ているはず。一行の終わりはセミコロンで。複数のデータや変数、 Y_i や σ_i は、 i ごとに変わるという
 4602 ことをあらわしていますが、プログラ的には `Y[i]` のように書いて、大括弧で変数を括って添字であ
 4603 ることを表現します。データは 7 件あって、それぞれについて i さんの測定値 Y_i と測定誤差 σ_i があ
 4604 るので、`for` 文で i を繰り返しています。変数 `sig[i]` についての事前分布も同様に繰り返していま
 4605 す。ちなみに変数名 `mu`, `sig`, `Y` などは任意で、予約語でなければ好きな名前にしてもらって構いま
 4606 せん。

4607 **parameters ブロック** `model` ブロックで、変数 `mu` や `sig` という任意に命名して使っていたわけですが、
 4608 それは何なのかを宣言してやらねばなりません。これらは求めたいパラメータなのですから、このブ
 4609 ロックに型と共に宣言するわけです。変数 `mu` は実数型、変数 `sig` も実数型ですが 7 つの要素を持
 4610 つ配列、しかも負の数を取らないので下限はゼロですよ、という宣言をしているところになります。範
 4611 囲の設定に `<>` の記号を使っていること注意してください。配列については次の `data` ブロックで説明
 4612 します。

4613 **data ブロック** 最後にデータブロックです。モデルの変数でもパラメータでもない、外部から与えられる
 4614 もなので、変数名 `Y` で宣言します。宣言のときに 1 つの変数に複数の要素がある (**配列**といいま
 4615 す) ことを明示するため、`array[size]` と書きます。この大括弧 (`[]`) の中身は整数で、今回は 7
 4616 件のデータなのでサイズは 7 としてあります。そしてその配列の数字が整数 (Integer) なのか、実数

4617 (Real) なのか、といった型宣言をしたうえで、変数名を書きます*6。

4618 いかがでしょうか。書き方 (文法) がわかると比較的単純、あるいは率直な書き方をしていることがお分か
4619 りいただけるかと思います。尤度と事前分布とデータを書くだけで、Stan は事後分布からの乱数を生成して
4620 くれます。確率密度関数を書いたり解いたりしなくても、結果の近似計算ができるというのはとても便利なこ
4621 とではないでしょうか。

4622 19.3 Stan を使った MCMC の実践

4623 さて Stan ファイルが準備できましたから、これはいったん保存しておいて、今度はこれを R の方から利用
4624 してやることになります。R からは `rstan` パッケージあるいは `cmdstanr` パッケージを使って、Stan を呼
4625 び出して結果を返してもらう、という流れです。Stan の言語は一度機械語に翻訳 (コンパイル) されるので、
4626 その間はしばらく待つ必要がありますが、それが終わると Stan から事後分布の乱数が次々生成され、R に
4627 返ってきます。R はそれら乱数からのサンプルを事後分布の代表値からなるデータセットとして扱って、記述
4628 統計や可視化を通じて確率分布の解釈に使うのですね。

4629 それではどのように R から呼び出すのか、実際のコードを見てみましょう。呼び出し方は `rstan` パッケー
4630 ジと `cmdstanr` パッケージのどちらを使うかによって少し異なります。共通部分のあと、両方の呼び出し方に
4631 ついて記述しますので、自分の環境にあった方を試してください。

4632 ■共通部分 どのパッケージを使うにしても、環境をクリアしたり、必要なパッケージを読み込んだりする部
4633 分は共通です。今回は次の 5 行を共通部分として読み込んでおいてください (code19.2)。

code : 19.2 環境の準備と共通部分

```
4634 1 rm(list = ls())
4635 2 library(tidyverse)
4636 3 library(bayesplot)
4637 4 # データ
4638 5 x <- c(-27.020, 3.570, 8.191, 9.898, 9.603, 9.945, 10.056)
4639
4640
```

4641 ■`cmdstanr` パッケージによる呼び出し その上でさっそく `cmdstanr` パッケージを使って、Stan を呼び
4642 出し事後乱数生成機を作ってみましょう。これは次のようにコードを書きます (code19.3)。

code : 19.3 `cmdstanr` パッケージによる呼び出し

```
4643 1 library(cmdstanr)
4644 2 model <- cmdstan_model("sevenScientist.stan")
4645 3 fit1cmdstan <- model$sample(
4646 4   data = list(Y = x),
4647 5   chains = 4,
4648 6   parallel_chains = 4
4649 7 )
4650 8 sample <- fit1cmdstan$draws() %>%
4651 9   posterior::as_draws_df() %>%
4652
```

*6 Stan の以前のバージョンでは `Y[i]` や `sig[i]` のように、大括弧 `[]` を後ろにつけて書いていました。しかしバージョン 2.32.0 はその表記は誤りとされ、ここで解説しているような `array` 表記がスタンダードになります。古い記法のままで警告が出たり、RStudio も古いバージョンだと新しい記法に変更要請 (赤い波線がコードについて) が出たりしますが、今すぐ問題になるというものではありません。とはいえ、なるべく新しいバージョンを使ったほうがよいので、表記法も新しいものに心がけましょう。

```
4653 10 as_tibble()
4654
```

4655 ■コード解説

4656 1 行目 パッケージの呼び出し

4657 2 行目 モデルのコンパイルを実行して model オブジェクトに入れます。コンパイルには少し時間がかかります。

4658 4-8 行目 乱数発生。ここで R のデータを Stan に渡したり、並列化の指定をしたりします。実行すると画面に色々表示されるかと思えます。

4659 9-11 行目 得られた乱数をデータセットに整形。まず posterior::as_draws_df で MCMC サンプルだけを取り出し、as_tibble で tibble 型に変換します。

4663 ■rstan パッケージによる呼び出し rstan パッケージを使って Stan ファイルを呼び出す場合は、少し書き方を変えて次のように書きます (code19.4)。

code : 19.4 rstan パッケージによる呼び出し

```
4665 1 library(rstan)
4666 2 options(mc.cores = parallel::detectCores())
4667 3 rstan_options(auto_write = TRUE)
4668 4 model <- stan_model("sevenScientist.stan")
4669 5 fitlrstan <- sampling(model, data = list(Y = x))
4670 6 sample <- fitlrstan %>%
4671 7   rstan::extract() %>%
4672 8   as.data.frame() %>%
4673 9   as_tibble()
4674
4675
```

4676 ■コード解説

4677 1 行目 パッケージの呼び出し

4678 2-3 行目 CPU の並列化や上書き保存など、パッケージの設定

4679 4 行目 モデルのコンパイルを実行して model オブジェクトに入れます。コンパイルには少し時間がかかります。

4681 5 行目 乱数発生。ここで R のデータを Stan に渡しています。実行すると画面に色々表示されるかと思えます。

4683 6-9 行目 得られた乱数をデータセットに整形。まず rstan::extract で MCMC サンプルだけを取り出し、as.data.frame でデータフレーム型に、as_tibble で tibble 型に変換します。

4685 どちらのパッケージで実行してもらっても結構ですが、コンパイル → サンプリングという流れは同じです。

4686 またサンプリングの際に、R からデータを渡すことになりませんが、データは R 環境では x というオブジェクト名だったものを、Stan の data ブロックで宣言した Y という変数にして渡すのだ、ということが書かれています。データは複数のものを与えることがあるので、list 関数で並べて渡すことにします。そのほかの設定については、順次説明します。

4690 最終的に tibble 型でえられた MCMC サンプルのデータセットができあがります。これはパッケージによって少し変数名が異なりますが、内容は同じで、今回の設定だと 4000 行のデータになっているはず。つまり事後分布からの乱数を 4000 個発生させたことになります。

4693 そしてこの各行が 1 回のサンプリングで得られた**事後分布からの代表値**です。一部を見てみましょ
 4694 う (出力 19.1)。1 回目のサンプリングで、9.98, 23.4, 4.63, 1.13, 0.123, 5.17, 0.0264, 1.24, -9.10 とい
 4695 う数字のセットがあります。変数名にあるように、前から $\mu, \sigma_1, \sigma_2, \dots, \sigma_7$ の代表値と、lp__ の値が
 4696 得られています。最後の lp__ は**対数尤度 (log posterior)** の略で、ここでは正確には事後対数尤
 4697 度と呼ばれるものです。今回の代表値から計算された尤度の対数ということで、パラメータの推定値
 4698 とは違いますのでここでは気にしないことにしましょう。ということで改めて、1 回目のサンプリングで
 4699 9.98, 23.4, 4.63, 1.13, 0.123, 5.17, 0.0264, 1.24 という数字のセットが得られました。このように、今回求めた
 4700 かったパラメータは 8 つあったわけですが、その 8 次元からなる空間の 1 つの座標点を取り出した 1 回目、
 4701 ということです。2 回目は 9.90, 27.2, 13.1, 2.86, 1.23, 0.779, 0.0449, 0.446 というセットでした。以下 3,4,..
 4702 と続き、4000 サンプルを得たわけですが。このように結果は 8 次元からなるパラメータの**同時確率空間 (joint**
 4703 **probability space)** から得られており、一行一行が 1 回 1 回のサンプリングステップなのです。

R の出力 19.1: MCMC サンプルの結果 (ごく一部)

```
# A tibble: 4,000 × 9
  mu sig.1 sig.2 sig.3 sig.4 sig.5 sig.6 sig.7 lp__
<dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1  9.98  23.4  4.63  1.13 0.123 5.17  0.0264 1.24 -9.10
2  9.90  27.2 13.1  2.86 1.23  0.779  0.0449 0.446 -7.72
3  9.99  20.9  5.12  1.11 0.132 4.48  0.0397 0.937 -8.97
```

4704

19.4 MCMC 結果の診断

4705

4706 さて、MCMC が実行できましたが、その中身を少しみてみましょう。サンプリング結果を代入したオブジェ
 4707 クト名をそのまま入力すると、要約された結果が示されます。パッケージによって出力例は少し異なりますが、
 4708 ほぼ同じ情報が含まれています。

cmdstanr の出力 2: MCMC の結果出力 (cmdstanr パッケージの場合)

variable	mean	median	sd	mad	q5	q95	rhat	ess_bulk	ess_tail
lp__	-9.70	-9.35	1.91	1.69	-13.28	-7.30	1.00	914	1694
mu	9.64	9.81	0.67	0.36	8.28	10.36	1.01	879	1101
sig[1]	54.19	31.91	105.83	16.95	15.41	142.60	1.01	663	215
sig[2]	10.48	6.71	20.30	3.96	2.88	25.82	1.00	1272	935
sig[3]	4.30	2.76	6.19	2.01	0.67	12.16	1.01	1281	1104
sig[4]	2.52	1.31	3.88	1.47	0.11	8.58	1.01	603	252
sig[5]	2.68	1.37	5.65	1.39	0.19	9.02	1.00	1291	1535
sig[6]	2.66	1.26	5.11	1.47	0.09	9.46	1.01	492	312
sig[7]	2.87	1.42	5.64	1.56	0.15	9.49	1.01	904	1049

4709

rstan の出力 2: MCMC の結果出力 (rstan パッケージの場合)

```
Inference for Stan model: sevenScientist.
4 chains, each with iter=2000; warmup=1000; thin=1;
post-warmup draws per chain=1000, total post-warmup draws=4000.
```

	mean	se_mean	sd	2.5%	25%	50%	75%	97.5%	n_eff	Rhat
mu	9.74	0.02	0.60	8.11	9.62	9.89	9.98	10.52	607	1.02
sig[1]	38.61	1.37	43.68	14.36	23.33	27.16	41.93	123.28	1016	1.01
sig[2]	9.44	0.57	12.88	2.86	4.63	6.22	10.57	33.54	515	1.00
sig[3]	4.64	0.40	14.66	0.63	1.33	2.49	4.98	17.98	1354	1.01
sig[4]	1.96	0.18	4.14	0.04	0.17	0.79	2.18	10.41	511	1.01
sig[5]	3.35	0.23	7.54	0.11	0.73	1.87	4.56	13.23	1044	1.01
sig[6]	2.00	0.25	3.65	0.03	0.22	0.77	2.44	10.55	216	1.02
sig[7]	2.50	0.17	5.71	0.09	0.44	1.24	2.49	12.28	1133	1.01
lp__	-9.45	0.07	1.67	-13.80	-10.19	-9.10	-8.39	-7.07	570	1.01

4710

4711 ここで確認して欲しいのは、**Rhat** です。これは「複数チェーンの収束の程度」を表しています。まず
 4712 MCMC によってサンプルを得る方法は乱数だったことを思い出してください。乱数ですから毎回同じ数字か
 4713 らスタートするわけではありません。そしてある乱数は次の数字を生むステップになっており、一步一步、事後
 4714 分布の形に近づいていくものです。しかし違う数字からスタートしても、目標とする事後分布は同じはずで
 4715 ずから、最終的にはこのステップの鎖はよく絡まったもの、つまりどこからスタートしても同じあたりをウロウロす
 4716 るものになるはずでず。ですから、MCMC がきちんとできているかどうか、事後分布を正しく代表した数字
 4717 になっているかどうかを確認する必要があります。この **Rhat** という数字は、1.0 であればピッタリ一致とい
 4718 う数字で、1.1 よりも小さければ十分「絡まっている」と判断できます。今回はいずれも十分な収束が得られ
 4719 たと言えるでしょう。

4720 このことを視覚的に確認することもできます。**トレースプロット (trace plot)** というのがそれで、MCMC
 4721 のステップを可視化し、十分に絡まっていることを見とくのです (図 19.3)。

4722 ちなみにトレースプロットの描画のコードは、code 19.5 の通りです。パッケージによって少し作法が違うの
 4723 で注意してください。cmdstanr で推定したものを使う場合は、bayesplot パッケージの力を借りて描画し
 4724 ます。

code : 19.5 トレースプロットの描画

4725

```
4726 1 # cmdstanr パッケージのオブジェクトからトレースプロットを描く場合
4727 2 fit1cmdstan$draws("mu") %>% bayesplot::mcmc_trace()
4728 3 # rstan パッケージのオブジェクトからトレースプロットを描く場合
4729 4 traceplot(fit1rstan)
4730
```

4731 これについては悪い例を見た方が早いかもしれません。図 19.4 に、チェーンが絡まない場合のトレースプ
 4732 ロットを示しました。このように、複数のチェーンがそれぞれバラバラ、迷子になっている場合は、手元に生成さ
 4733 れた乱数が事後分布の適切な代表になっていないと判断することになります。

4734 次に、**有効サンプルサイズ (effective sample size)** を見ておきましょう。rstan パッケージを使っ
 4735 ている例の場合は n_eff、cmdstanr パッケージを使っている例の場合は ess_tail を見ましょう。これは
 4736 MCMC サンプルがうまく取れたかどうかを診断した結果で、ここが少なくとも 3 桁ほどなければ不十分な
 4737 だ、と考えます。いくらあったらいいかについては、どれだけ MCMC サンプルを発生させたかに依存するの
 4738 で一概には言えませんが、少なくとも効果的に得られているのが 100 サンプルぐらいいないとダメなんだ、とい

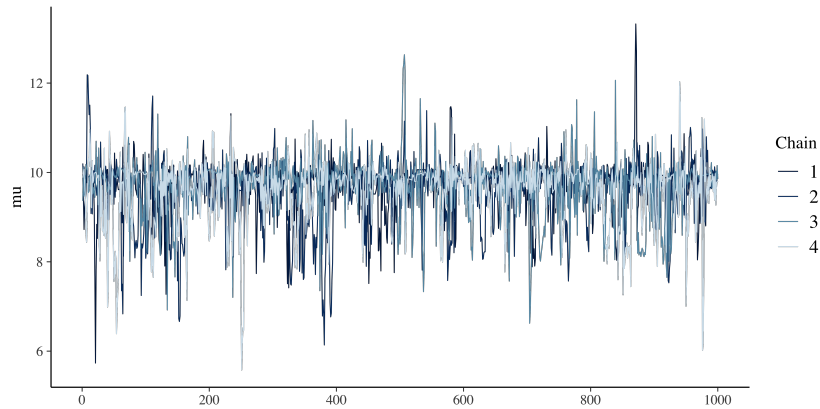


図 19.3 トレースプロットの例

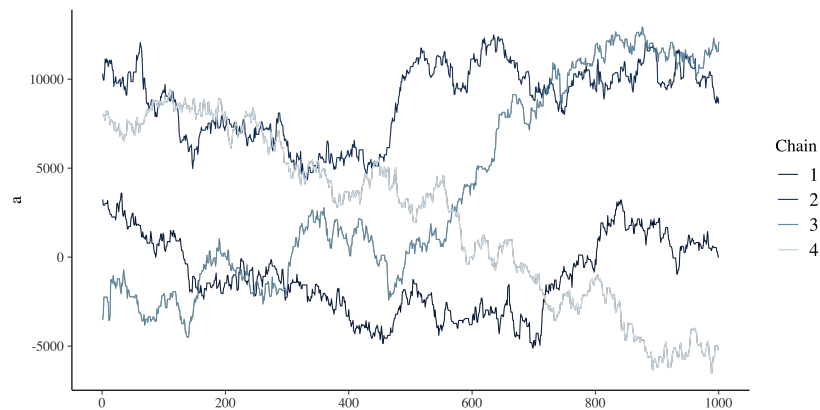


図 19.4 チェインが絡まない場合

4739 うことです。

4740 ちなみに、MCMC サンプルのチェーン (発生させた系列数) や、MCMC サンプルの数は、関数のオプションで指定します。今回は指定なしだったので、デフォルトの値で 4 チェイン、各々 1000 サンプルずつ作るようにしていました。これらを明示的に設定するには次のようにします (code19.6,code19.7.)。

code : 19.6 cmdstanr パッケージによるオプション指定

```
4743 1 fit1cmdstan <- model$sample(
4744 2   data = list(Y = x),
4745 3   chains = 4,
4746 4   parallel_chains = 4,
4747 5   iter_warmup = 1000,
4748 6   iter_sampling = 4000
4749 7 )
4750
4751
```

code : 19.7 rstan パッケージによるオプション指定

```
4752 1 fit1rstan <- sampling(model,
4753 2   data = list(Y = x),
4754 3   chains = 4, iter = 5000, warmup = 1000
4755 4 )
4756
```


4757

4758 ここにある `chains` というのが生成する MCMC サンプルの連鎖系列の数です。「どこからスタートし
4759 ても同じ事後分布にたどり着く」ことを検証するために、少なくとも複数のサンプルが必要で、デフォルト
4760 では 4 本のチェーンを作っていることになります。cmdstanr パッケージの方は、並列化する本数も別途
4761 `parallel_chains` 関数で指定します。

4762 また生成する乱数の数ですが、ここにウォームアップ (`wamup`) という単語が出てきています。これは
4763 MCMC サンプルの際、最初の数ステップは調整のための準備運動期間中なので、この期間に得られた
4764 MCMC サンプルは代表値として使わない、という意味です。cmdstanr パッケージの場合、捨てるステップ
4765 数とサンプルするステップ数を明示的に指定していますからわかりやすいですね。最終的に得られる MCMC
4766 サンプルの数は $4000 \times 4 = 16000$ で計算できます。rstan パッケージの場合、全体の反復 (`iter`) 回数か
4767 らウォームアップの分を引き算しなければなりません。引き算した残りの期間をチェーン数発生させることにな
4768 ります。今回の例では $(5000 - 1000) \times 4 = 16000$ サンプルできあがることになります。

4769 19.5 MCMC の結果の解釈

4770 さあ長々と話をしてきましたが、どうやら Rhat や有効サンプルサイズにも問題がなかったので、MCMC
4771 は事後分布からの適切な代表値を産んでいたと考えられますから、その結果を考えてみましょうか。

4772 ここで、前の章 (セクション 18.5.3, Pp.197) で準備した関数を使ってみましょう*7。これらをコード 19.8 の
4773 ように使うと結果が得られます。私の結果では出力 1 のようになりました。

code : 19.8 stan の結果を準備し関数に与える

4774
4775
4776

```
1 fit %>% MCMCtoDF() %>% MCMCsummary()
```

MCMC の結果 1

```
# A tibble: 8 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu          9.689    9.863    9.929    0.647    7.982    10.602
2 sig[1]     45.886   31.238   26.266   70.414   14.011   158.906
3 sig[2]     10.037    6.781    5.206   13.746    2.635    36.035
4 sig[3]      4.609    2.872    2.067    8.394    0.494    18.808
5 sig[4]      2.442    1.060    0.208    6.055    0.029    12.608
6 sig[5]      2.856    1.421    1.085    9.126    0.109    12.860
7 sig[6]      2.688    1.175    0.443    6.511    0.058    13.277
8 sig[7]      2.715    1.251    0.234    6.232    0.037    13.120
```

4777

4778 ここでは事後確率最大値 (Maximum A Posterior) を見ることにしましょう。これによると、 μ すなわ
4779 ち真の測定値は 9.596 ぐらいじゃないか、という推測がなされています。10 より少し小さいくらいではない
4780 か、と思っていたようですが、こうして計算してみると確かにそれぐらいのようですね。そして 7 人それぞれの
4781 測定精度が σ_i として算出されています。これをみると、最初の測定者は 23.101, 2 番目の人は 6.849 になっ
4782 ています。これが大きいということは、誤差の幅広さを表しているのですから、測定精度が悪いということ

*7 伴走サイト https://kosugitti.github.io/psychometrtics_syllabus/ からは、これらの R コードをダウンロードで
きます。Stan ファイルは自分で準備し、読み込み先やファイル名などは適宜変更してください。

4783 味します。確かに最初の二人は非常に精度が悪く、他のベテラン測定者が 0.3 とか 0.9 ぐらいの誤差である
4784 ことに比べれば、圧倒的に悪いということが言えますね。

4785 もちろんこれは点推定値で、真の測定値が 9.596 に違いないとか、最初の人々の測定精度が 23.101 だと断
4786 言すると、おそらくほぼ確実に外れた予測ということになるでしょう。そこで幅を持った予測、すなわち**区間推
4787 定 (interval estimation)**を行えば良いのです。区間を考えれば、真の測定値は 10.547 から 7.753 の
4788 範囲にあるだろうとか、最初の測定者の精度も最悪 151.573、最善なら 14.259 ぐらいである可能性がある
4789 わけです。

4790 バイズ推定のポイントとして、これらは事後分布、すなわち**確率分布による予測**であるということが挙げら
4791 れます。真の測定値 μ のありそうな範囲が 10.547 から 7.753 にある確率が 95% だ、という言い方ができる
4792 ところです。**モーメント法**による**信頼区間 (confidential intervals)**とは違い、バイズ法の予測は**確信区
4793 間 (Credible Intervals)**という言い方になります*8。

4794 いかがでしたでしょうか。今回はモデリングとバイズ推測の実際、Stan コードの書き方を実践してみました
4795 た。この「7 人の科学者」がおもしろい点は、少ないサンプル (たった 7 件のデータ!) で推測ができること、簡
4796 単なコードで検証できること、というのがありますが、なにより標準偏差 (分散) を検証対象とすることにある
4797 と私は思います。心理学の実験のほとんどが、平均の比較、操作の効果の話になっているのですが、データの
4798 散らばりというのも測定精度のように心理学的に意味のあるものとして考えられ、検証の対象にしたっていい
4799 のです。もちろん正規分布でなくてもいいですし、さまざまな分布のさまざまなパラメータに意味があるなら、
4800 それをわからないものとして検証しようと言え、というのはとても可能性が広がる話ではないでしょうか。

4801 19.6 課題

4802 次の計算をする R/Stan コードを記述し、提出してください。なお提出されたコード単体でバグがなく動く
4803 ことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別
4804 TA か小杉までメールで連絡し、指導を受けてください。

- 4805 1. 8 人目の科学者が現れて、測定値 18.25 だったと報告してきました。この人のデータも含めた分析モ
4806 デルにするために、R および Stan のコードを書き換え、MCMC 推定してください。また 7 人の時と
4807 くらべて結果がどう変わったでしょうか? 気づいたことを報告してください。
- 4808 2. MCMC サンプルの精度を上げるため、5 つのチェーンで、warmup 期間を 10000、最終的な
4809 MCMC サンプルを 10 万点得られるように R のコードを書き換えてください。
- 4810 3. MCMC サンプルを使って、 μ と σ の分布の形を可視化するコードを書いてください。

*8 信頼区間の場合は、ここにあると予測したとして、その予測が当たる確率を表しています。すなわち当たるか外れるか、0/1 の話
をしているのであって、この「範囲に存在する」という大きさ・幅の話をしているわけではないことに注意です。

4811 第 20 章

4812 モデリングの目から見た検定 1 ; 二群の 4813 平均値の差

4814 前はデータ生成メカニズムという観点で考えることで、標準偏差 (分散) を「測定精度」と考えた検討が
4815 できることを示しました。今回は、心理学でもよく使われている平均値の差を検討対象とするモデルを考える
4816 ことにします。

4817 皆さんは帰無仮説検定のことを覚えているでしょうか。正規分布を仮定した母集団からの標本統計量は、
4818 正規分布に従うことを利用して、母平均の区間推定を行い、群間に差があるかないかといった判断を確率的
4819 に行うのが帰無仮説検定でした。心理学では要因計画と帰無仮説検定が合体し、群間の差の形でデータが
4820 得られるようにすることで結論を導き、考察を重ねてきました。また帰無仮説検定という手法は、誤用や誤解
4821 が多く、ひどい時には悪用もされるということについてもみてきたと思います。ところで、帰無仮説検定の考え
4822 方はデータが得られたところから考え始めるデータ駆動型の分析モデルでした。ではデータがどのように出て
4823 きているのかを考える、データ生成モデリングの考え方をつかおうと、この平均値の差についての検討はどの
4824 ように姿を変えるのでしょうか。

4825 ここでは **t 検定** を例に話を進めていきたいと思います。

4826 20.1 t 検定の過程と実際

4827 t 検定は、二群の平均値の差を比較し、結論を下すという最も典型的な帰無仮説検定の例の 1 つです。と
4828 くに独立した二群の検定は、最も単純な要因計画 (一要因二水準 Between デザイン) ということができるで
4829 しょう。これについて、よく思い出せない人はデータ解析基礎の資料や、山田・村井 (2004)、清水 (2021) な
4830 どを読んで、分析の流れや仮定をもう一度確認しておいて欲しいと思います。

4831 非常に駆け足ながら概略を説明してみますと、次のようになります。

- 4832 1. 標本の母集団が正規分布していると考える。母集団から無作為に選ばれた標本が二群あるとする。
- 4833 2. 二群の一方に何らかの処置を加え、他方には何もしない (あるいは検討したい処置と同等の効果のな
4834 い処置を施す^{*1})。前者を実験群、後者を統制群という。データ (標本) は正規分布するはずだから、
4835 その平均値を見ることで誤差や個人差は相殺される。つまり群間の平均値の差が効果の大きさであ
4836 るということができる。

*1 たとえば「単語リストは声に出して覚えたほうが記憶の定着度が高い」ということを検証したい場合、声に出す群と出さない群で比較することになりますが、声に出さないことが (頭の中で反芻するなど) 従属変数に与える別の効果を持っていることがあるので、音のない映像を見せるなどして効果のない同等の処置をした群を比較対象にする、ということを考えたりするわけです。

- 4837 3. ところで標本平均値などの**標本統計量**も正規分布に従うことがわかっている。母集団が $N(\mu, \sigma)$ で
 4838 あれば、標本平均は $N(\mu, \frac{\sigma}{\sqrt{n}})$ に従う (ここで n はサンプルサイズ)。
- 4839 4. 標本平均が従う分布 (散らばりの位置と幅) が分かれば、母平均がどのあたりにあるのかは区間推定
 4840 することが可能。実験群・統制群ともに母平均の値を推定する。これが異なっていれば標本を超えて
 4841 母集団で効果があった、と結果を一般化できる。もちろん点推定値は母平均の値とピッタリ同じとは思
 4842 えないし、標本平均もピッタリ同じになるはずがないから、点推定値ではなく区間推定で考えたい。
- 4843 5. ところで標本平均の従う分布の中には、 σ つまり母 SD がパラメータとして入っている。母数がわから
 4844 ないという前提のもとでは、どの程度の幅で散らばるのがわからないことと同じ。そこで σ の代わり
 4845 に $\hat{\sigma}$ を用いて推定することを考える。この場合、標本平均は正規分布ではなく t 分布に従うことがわ
 4846 かっている。
- 4847 6. t 分布を使った区間推定をしても、差があるのかないのか判断するために何らかの基準が必要であ
 4848 る。そこで「差がない」という主張と「差がないとは言えない」という主張を戦わせて、判定するという形
 4849 式を取る。前者の主張を**帰無仮説**、後者の主張を**対立仮説**という。また判断も確率的になるので、勝
 4850 敗を決める基準を 5% とする。この確率は**有意水準**とか**危険率**と呼ばれる。
- 4851 7. データから t 分布に従う統計量を計算し、その値が出てくる確率を算出する。この確率は p 値と呼ば
 4852 れる。これが有意水準よりも小さいようであれば、帰無仮説の仮定の下で算出した数字が滅多に生じ
 4853 得ないことを意味するから、帰無仮説が間違っていたのだと判断してこれを**棄却**し、対立仮説を**採択**
 4854 する。

4855 さて、この二群の平均値を検定することの流れですが、とくに「帰無仮説のもとでの t 値を算出して判定す
 4856 る」というあたりがややこしかったかもしれません。数式で書くと嫌われそうですが、帰無仮説というのは二群
 4857 の平均値に差がない、という仮定だったので、実験群から考えられる母平均 μ_A と統制群から考えられる母
 4858 平均 μ_B に差がない、つまり $\mu_A = \mu_B$ 、からの $\mu_A - \mu_B = 0$ という位置母数が 0 の理論分布のを考えて
 4859 いる、ということがポイントです。検定統計量である t 値は次のように計算できるのです。

$$t = \frac{\bar{X}_A - \bar{X}_B}{\sqrt{\frac{(n_1-1)\sigma_A^2 + (n_2-1)\sigma_B^2}{n_1+n_2-2} \left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

4860 とっても複雑な式に見えますが、複雑そうなのは分母で、分子は標本平均の差を表しているに過ぎません。
 4861 参照する t 分布は、標準正規分布が $N(0, 1)$ だったように、位置と幅のパラメータを持ち、 $t(0, df)$ で考えら
 4862 れる分布にこの統計量を照らし合わせると、差 0 で自由度に応じて変わる幅 df の分布から「どの程度極端
 4863 な値が出てきたのか」の確率を計算できるわけです。ちなみに分母は、母分散ではなく標本から計算される分
 4864 散 s^2 からバイアスを除いた不偏推定量 $\hat{\sigma}^2$ を使っています。この計算はサンプルサイズ n ではなく $n-1$ で
 4865 割ることによって計算されるのですが、2 つの群それぞれについて $n_j - 1$ で割ったものを足し合わせるため
 4866 に、いったん $n_j - 1$ 倍して二群のサンプルサイズ -2 で割り直す、という作業をしているため、分母がとくに
 4867 複雑に見えているだけです。

4868 ともあれ、こうして検定するんだったという流れを思い出したところで、データ生成モデリングの観点か
 4869 らこれを考え直してみましよう。仮定としておいてあるのは、両群ともに同一の正規分布から得られた標
 4870 本であり、操作によってその平均値が異なっているはず、ということだけです。つまり実験群のデータは
 4871 $X_{i,A} \sim N(\mu_A, \sigma)$ 、統制群のデータは $X_{i,B} \sim N(\mu_B, \sigma)$ という分布が前提とされているのです。

4872 これをそのまま設計図にし、コードにしてみましよう。設計図として図 20.1 が、これをもとにコードが
 4873 code:20.1 のようにかけていけばいいでしょう。コードが設計図をほぼそのまま文字起こしていることを確
 4874 認してください。

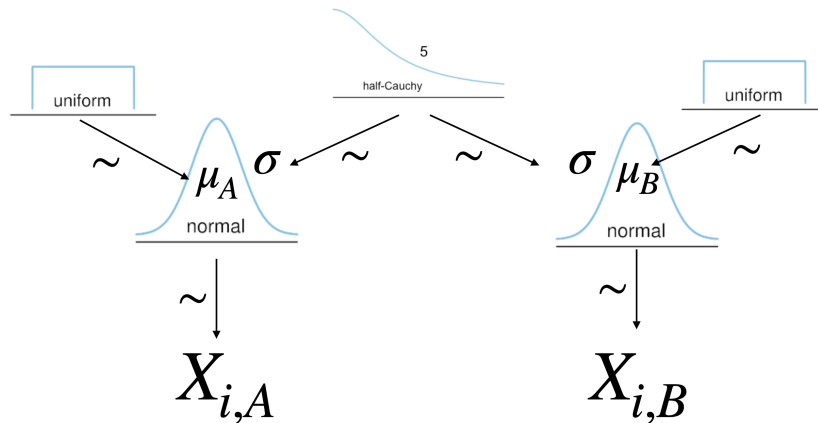


図 20.1 平均値の差を比べるときの設計図

code : 20.1 二群の平均値のコード

```

4875 1 data{
4876 2     int<lower=0> N1; // Number of Subjects in Group 1
4877 3     int<lower=0> N2; // Number of Subjects in Group 2
4878 4     array[N1] real X1; // Data in Group 1
4879 5     array[N2] real X2; // Data in Group 2
4880 6 }
4881 7
4882 8 parameters{
4883 9     real mu1;
4884 10    real mu2;
4885 11    real<lower=0> sig;
4886 12 }
4887 13
4888 14 model{
4889 15     // likelihood
4890 16     X1 ~ normal(mu1,sig);
4891 17     X2 ~ normal(mu2,sig);
4892 18     // prior
4893 19     mu1 ~ uniform(0,100);
4894 20     mu2 ~ uniform(0,100);
4895 21     sig ~ cauchy(0,5);
4896 22 }
4897
4898

```

4899 ここで Stan コードにちょっとした工夫を 2 点入れています。1 つ目は data ブロックにある、変数 $N1$, $N2$
4900 の存在です。二群のデータについて、サンプルサイズがとくに決まっていませんから、ここで外部から入力す
4901 ることにしています。2 つの群それぞれのサンプルサイズをデータとして取り込み、その数と同じだけデータ数
4902 を配列として宣言しているところがポイントです。もう 1 つのポイントは、尤度のところですが、丁寧に書くなら
4903 ば、次のようにするべきでしょう (コード 20.2)。

code : 20.2 丁寧な尤度の記載

```

4904 1     ...
4905

```

```

4906 2  model{
4907 3      for( i in 1:N1){
4908 4          X1[i] ~ normal(mu1 , sig);
4909 5      }
4910 6      for( i in 1:N2){
4911 7          X2[i] ~ normal(mu2 , sig);
4912 8      }
4913 9      ...
4914 10 }
4915

```

4916 このように、それぞれの群において for 文を回し、各データ点が正規分布から出てきているよ、とい
 4917 うことを明示するのです。しかしそうしなかったのは、Stan が「分かりきったことは書かなくていいよ」
 4918 という優しい設計になっているので、変数 X1 の要素すべてがある分布に従うのであれば、まとめて
 4919 $X1 \sim \text{normal}(\mu1, \text{sig})$; と書いても良い、つまり X1 の配列要素を逐一指定しなくても同じ尤度関数を
 4920 あてがってくれるというのを利用しています。

4921 ともかくこれでできるのです。では少し具体例をつかって、これがどういう結果をもたらしているのかを確
 4922 認しましょう。

4923 20.1.1 t 検定の具体例

4924 統計学の新しい指導法の教育効果を見るため、全国の大学の心理学科から学生を無作為に
 4925 16 名選び、従来の指導法グループ(統制群)と、新指導法のグループ(実験群)にランダムに 8
 4926 名ずつわりつけました。プログラム終了後、心理統計のテストを行いました。統制群のスコアは
 4927 20, 40, 60, 40, 40, 50, 40, 30, 実験群のスコアは 30, 50, 70, 90, 60, 50, 70, 60 となりました。新しい
 4928 指導法は、心理学科の学生に効果があると言えるでしょうか? 5% 水準で検定してください。

4929 これを t 検定するのは簡単ですね。計算式は複雑でも、機械がやってくれるから楽なのが NHST^{*2} のいい
 4930 ところです。

code : 20.3 帰無仮説検定のコード

```

4931 1  groupA <- c(30, 50, 70, 90, 60, 50, 70, 60)
4932 2  groupB <- c(20, 40, 60, 40, 40, 50, 40, 30)
4933 3  ## t 検定
4934 4  t.test(groupA, groupB, var.equal = TRUE)
4935
4936

```

4937 結果は出力 20.1 のようになります。検定統計量である t 値、検証のための自由度 df 、帰無仮説のもとで
 4938 の出現確率 p 値、が表示されています。5% より小さいので有意差あり、という判断ができます。もっとも、 t
 4939 値とか自由度とかは、データに関係のない数字なのでピンと来ないというところもあるかと思います。

*2 帰無仮説検定 (Null Hypothesis Significance Test) の略です

R の出力 20.1: 検定の結果

```
> t.test(groupA, groupB, var.equal = TRUE)

Two Sample t-test

data:  groupA and groupB
t = 2.6458, df = 14, p-value = 0.01919
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 3.786937 36.213063
sample estimates:
mean of x mean of y
    60      40
```

4940

20.1.2 モデルで推定してみる

つづいて、先ほどのコードを使ってモデリングによる推定をしてみたいと思います。私の環境下では、次のような数字になりました。

MCMC の結果 2

```
# A tibble: 3 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu1    60.006    59.976    59.700    5.608    48.869    71.420
2 mu2    39.990    39.997    39.650    5.679    28.654    51.128
3 sig    15.534    15.048    14.070    3.111    10.893    22.877
```

4944

これを見ると、未知のパラメータ μ_1, μ_2, σ の値がそれぞれ推定されていますが、「判定」のような結果は出てきていません。というのも、勝負するような形で考えているのではなく、パラメータがどこにありそうか、ということを示す**確率分布**を求めているからです。それでも、**EAP** をみると 60.006 と 39.990、95%CI をみると実験群は [48.869, 71.420]、統制群は [28.654, 51.128] とあり、実験群の平均値が 51 より小さい値になる確率はほとんどなく (2.5% 以下)、統制群の平均値が 49 以上より大きくなることもまたほとんどないわけですから、確実に差があると言ってもほぼ過言ではない、と言い切れるでしょう。

モデリングの利点は、 t 値や p 値のような直接のデータと関係ない値を出すのではなく、データに直接関係する数字で考えさせてくれるので、イメージしやすいところもあるかもしれませんね。

4952

20.1.3 分散の等質性

ところで、 t 検定の場合は「分散が同じであるかどうか」という条件によって、算出方法を補正するという考え方があるのを覚えていますでしょうか。**Welch の補正**というやつで、こちらの方が一般に条件が緩いものですから、 t 検定では普通こちらの方が使われます。

4956

R の出力 20.2: Welch の補正の例

```
> t.test(groupA, groupB, var.equal = FALSE)

Welch Two Sample t-test

data:  groupA and groupB
t = 2.6458, df = 12.274, p-value = 0.021
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 3.570406 36.429594
sample estimates:
mean of x mean of y
      60      40
```

4957

出力 20.2 にあるように、オプションとして指定するだけで瞬時に答えが出てきます。t 値、自由度、p 値が先ほどと異なって出ていますが、結果はほぼ変わりありません。

モデリング上ではこれをどう表現すれば良いのでしょうか。これは実は簡単で、分散が実験群と統制群で異なるというのですから、別々に推定してやれば良いのです。

code : 20.4 分散が異なる場合の推定モデル

```
4962 1      ...
4963 2 parameters{
4964 3     real mu1;
4965 4     real mu2;
4966 5     real<lower=0> sig1;
4967 6     real<lower=0> sig2;
4968 7 }
4969 8
4970 9 model{
4971 10    // likelihood
4972 11    X1 ~ normal(mu1, sig1);
4973 12    X2 ~ normal(mu2, sig2);
4974 13    // prior
4975 14    mu1 ~ uniform(0, 100);
4976 15    mu2 ~ uniform(0, 100);
4977 16    sig1 ~ cauchy(0, 5);
4978 17    sig2 ~ cauchy(0, 5);
4979 18 }
```

4980

4981

20.2 差の分布

4982

さて、二群のデータから考えられるそれぞれの母平均が推定されました。今回は実験群と統制群の平均値差が大きく離れており、一方の 95% 上限が他方の 95% 下限を上回っているという状況でしたので、「差がある」と判断できましたが、そうでないこともありそうです。つまり、微妙な差のとき、ですね。次のようなデータで試しに推定してみましょう。まずは 5% 水準で検定をしてみます。

code : 20.5 帰無仮説検定のコード

4987


```

4988 1 groupA <- c(30, 50, 70, 90, 60, 50, 70, 60)
4989 2 groupB <- c(30, 45, 60, 40, 60, 50, 40, 30)
4990 3 ## t検定
4991 4 t.test(groupA, groupB, var.equal = FALSE)
4992

```

4993 今度は $t(12.175) = 2.0761, p = 0.0597$ で有意とは言えなくなりました*3。このデータを、先ほどのモデル
 4994 で推定してみましょう。目にも分かりやすくするために、推定された平均値の事後分布をプロットしてみたいと
 4995 思います (図 20.2)。

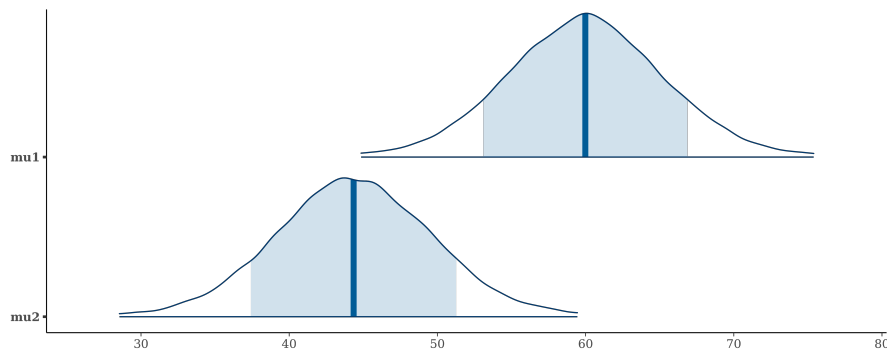


図 20.2 二つの平均値パラメータの事後分布

4996 今回は 80% 確信区間と、確率分布の両端を 99% のところまで描画してみました。2つの平均値パラメー
 4997 タは、分布の代表値で見ると異なっているのですが、重複している領域が多く一概に「一方が他方より大き
 4998 い」と言いにくい結果になっています。帰無仮説検定のようにズバっと「差がある」「ない」と言い切れればいい
 4999 のですが、このような場合はどうしたらいいのでしょうか。

5000 正解は「このままにしておく」です。急いで差があるとかないとか言い切るのではなく、これが区間推定の結
 5001 果そのものですから、この程度の重複があり得るもの・それほど明確な答えは出ないもの、として考え続ける
 5002 しかありません。

5003 しかしまあ、このままでは考えるヒントが少ないので、ここで少し工夫をしてみましょう。帰無仮説検定では、
 5004 $\mu_A = \mu_B$ という帰無仮説が反駁できるかどうかの問題なのでした。今回、 μ_A や μ_B を代表する数字の候補
 5005 が色々ありますので、これを使って「実際の程度差があったのか」を計算してみれば良いのです。MCMC
 5006 サンプルは、毎回推定したいパラメータの同時分布からの代表値を持ってきているという話をしました。たと
 5007 えば 1 回目のサンプルは $\mu_A = 62.4, \mu_B = 25.7$, 2 回目のサンプルは $\mu_A = 60.4, \mu_B = 36.8$ と言った具
 5008 合に、です。それらは事後分布からのサンプル、事後分布を代表した実現値の 1 つですから、これを使って
 5009 $\mu_A - \mu_B$ の計算をできます。すると差の分布の代表値が MCMC サンプルで計算できることになります。

5010 これを使って、 $\mu_A - \mu_B = 0$ になる確率はどれぐらいか、というのを考えれば良いのかもかもしれません。もっ
 5011 とも、連続変数におけるある一点の確率はゼロです*4、代表値といってもピッタリゼロになることはほぼあり
 5012 得ないでしょうから、「どの程度ならゼロと見做すか」ということを考える必要があります。この範囲は**実質的
 5013 に等価な範囲 (Region Of Practical. Equivalence; ROPE)** と呼ばれ、たとえばこのテストの例で
 5014 は ± 5 点ぐらいは誤差みたいなものだけど、5 点以上変わったら意味のある違いだ、ということであれば「5
 5015 点以上点差が出た確率」を計算できます！

5016 この計算は、MCMC サンプルを取り出した R のオブジェクトを使って計算することもできますが、実は

*3 5.9% だから惜しい！とか「有意傾向にある」なんて変なこと言い出したらダメですよ。勝負は 5% だと決めて始めたのですから。

*4 確率は相対的な面積で表される数字ですから、面積を持たない点については確率が定義できないのです。

5017 Stan の中で計算させてしまうこともできます。generated quantities ブロックが、この MCMC サンプル
5018 を取り出したあとの処理をするブロックに該当します。ここで計算される量のことを**生成量 (generated
5019 quantities)**と呼びます。実際にどのように使うのかをみてみましょう。

code : 20.6 生成量ブロックの利用

```
5020
5021 1     ...
5022 2 generated quantities{
5023 3     real diff;
5024 4     int<lower=0, upper=1> FLG;
5025 5     diff = mu1 - mu2;
5026 6     if(diff > 5){
5027 7         FLG = 1;
5028 8     }else{
5029 9         FLG = 0;
5030 10    }
5031 11 }
5032
```

5033 コード 20.6 には生成量ブロックを使う例を示しました。ここもブロックですので、変数の宣言が必要です。
5034 今回はまず diff という変数を実数型で宣言してみました。これはその下で $\mu_1 - \mu_2$ として定義されている
5035 ことからわかるように、差の分布 (の代表値) を生成す流のです。これを見ると、差の生じる確率を考慮すること
5036 ができます。次に整数型 (int 型) で FLG という変数を宣言しました。上限が 1 で下限が 0、という実質的
5037 には 0 か 1 の数字しか入らない、ある/なし、成立する/しないといった情報しか持たない変数です。いわゆる
5038 「フラグが立ったかどうか」の指標であり、その実態は下の if 文で表現されています。ここでは先ほど計算し
5039 た diff が 5 よりも大きいのであれば 1、そうでなければ 0 を代入する、というコードになっています。

5040 これを実行すると、パラメータのサンプルに加えて、そのパラメータから計算されるさまざまな量も同時に算
5041 出されます。出力例を見てみましょう。

MCMC の結果 3

```
# A tibble: 6 × 7
  name      EAP      MED      MAP      SD      L95      U95
<chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 diff    15.650    15.605    14.940     8.156    -0.516    31.706
2 FLG      0.911      1.000      1.000     0.285     0.000     1.000
3 mu1     60.035    60.035    60.112     6.737    46.602    73.575
4 mu2     44.385    44.374    44.291     4.594    35.249    53.696
5 sig1    18.509    17.472    15.960     5.369    11.414    31.792
6 sig2    12.496    11.752    10.761     3.683     7.629    21.594
```

5042
5043 ここにあるように、差の 95% 確信区間は [-0.516, 31.706] であることがわかります。差が大きければ
5044 31.70、小さければ逆転して-0.52 になる可能性すらあるのです。この確信区間は 0 を含んでいますから、差
5045 がない可能性は否定できません。そういう意味で、帰無仮説検定的表現を借りるなら、5% 水準で差がない
5046 とは言い切れない、ということができるでしょう。

5047 とはいえ、せっかく分布として情報が得られているのに、点推定にして勝負をするのはもったいないともい
5048 えるわけです。結論を急がずに、どの程度の差があるのかをじっくり考えてみるのがいいでしょう。また変数
5049 FLG を見ると、その平均が 0.908 となっています。これは変数 diff が 5 よりも大きい時は 1、そうでなけれ
5050 ば 0 という条件の数字で、その平均は 1 になった数を総 MCMC サンプル数で割った相対度数になっている

5051 のですから、確率を表す数字だと考えても良いことになります。つまり、差が5よりも大きくなる確率は90.8%
5052 である、と考えることもできるのです。これを応用すると、さまざまな仮説を考えることができそうですね。

5053 20.3 帰無仮説検定を省みる

5054 さてこのように考えてくると、帰無仮説検定について異なる視点で見ることができるようになったのではな
5055 いでしょうか。ここでは、帰無仮説検定の独特さを3つの点から考えてみたいと思います。

5056 20.3.1 不平等な対決

5057 帰無仮説検定は帰無仮説 vs 対立仮説という勝負に持ち込んで判断する方法である、ということを目
5058 頭に述べました。平均値の差の検定の文脈で言えば、帰無仮説は $H_0: \mu_1 = \mu_2$ であり、対立仮説は
5059 $H_1: \mu_1 \neq \mu_2$ です。このように、帰無仮説は「差がない」の一点張りであり、対立仮説は「それ以外ならなん
5060 でもあり」という状態を指しているのです。そもそも、帰無仮説というのは非常に限定的な仮説だったので
5061 す。無に帰してほしい仮説とは言え、あまりにも不利な勝負なのでした。その上で、勝負は「差がない」「ある」
5062 のどちらかにしかありません。ごく微妙な差であっても、判定によれば差が「あった」といえますし、ギリギリで
5063 もなければ「なかった」というしかないのです。このような過度に結果を際立たせる報告は、科学的な研究実
5064 践の文脈では拙速なことになりかねませんので、注意が必要なのでした。

5065 **データ生成モデリング**のやり方で同じデータを分析してみると、差も分布の形で描かれますから、「あった
5066 か、なかったか」という単純な問題にせずに色々考えてみたくなるのではないのでしょうか。

5067 20.3.2 量的な判断を

5068 もっとも、出力 20.2 の t 検定の結果をよくみてみると、t 値や p 値の下に 95% 信用区間が示されており、
5069 ここでは [3.570406, 36.429594] という数字が出ています。そうです、帰無仮説検定をやっているようで
5070 も、しっかり区間推定した結果も表示されていたのです。これは差の信用区間で、この区間が0を跨いでい
5071 ないということは、差がないとは言えない ($\mu_A - \mu_B = 0$)、ということを意味しています。この区間の幅をみ
5072 ると、3 から 36 と随分ひろくとられているように思えます。つまり、今回の検定結果は差があるとされたけれ
5073 ども、それほど確かなものではないかもしれないぞ、と慎重に判断することもできたはずなのです。

5074 この差についての量的な評価については、**効果量 (effect size)** という指標で表現されるのでした。効果
5075 量とは、標準化された差の大きさです。標準化するというのは、標準偏差で割って幅を整える、あるいは標準
5076 偏差何個分という単位で表現することでもあります。比べられるスコアはさまざまですから、その標準偏差で
5077 表現することで一般的に議論できるのです。平均値の差の検定においては、**Cohen の d (Cohen's d)** な
5078 どが指標としてもいいられますが、それは $\frac{\bar{X}_1 - \bar{X}_2}{\sigma}$ で計算されるのでした。

5079 これは生成量を使って簡単に計算できます。今回の場合、計算する場合は、生成量に次のようにすれば良
5080 いのです。

code : 20.7 生成量ブロックで効果量の算出

```
5081 1      ...
5082 2  generated quantities{
5083 3      real diff;
5084 4      real cohen_d1;
5085 5      real cohen_d2;
5086 6      diff = mu1 - mu2;
```

```

5088 7      cohen_d1 = diff / sig1;
5089 8      cohen_d2 = diff / sig2;
5090 9  }
5091

```

5092 どちらの群の標準偏差を基準にするか、ということを考える必要がありますが^{*5}、相対的な大きさの比較も
 5093 簡単にでき、しかもその結果も分布の形で得られる、すなわち効果量の確信区間 (効果が少なくともどの程
 5094 度ありそうかとか、大きければどの程度ありそうかとか) を考えることもできるのです。

5095 もちろんこれらの計算は、帰無仮説検定の文脈、言い換えれば伝統的な統計手法 (モーメント法による推
 5096 論) であってもできたことなのですが、データ生成モデルという観点からみるといっそう理解がしやすいので
 5097 はないかと思います。

5098 20.3.3 方向性を持った仮説も

5099 今回は最後に、「実験群が統制群よりも 5 点以上大きくなる確率」というのを考えました。この「一方が他
 5100 方よりも大きい」という仮説は、帰無仮説検定の文脈では**片側検定 (one-tailed test)** と呼ばれます。一
 5101 般的に「差がある」という対立仮説は、プラスであれマイナスであれ違いが生じている、という意味であり、統
 5102 計量は分布の右側でも左側でも、とにかく極端な値になりさえすれば良いという判断でした。そうではなく
 5103 「 $A > B$ のはずだ」という仮定であれば、小さくなる可能性を考えなくていいのですから、統計量のどちらか
 5104 一方の極について考えれば良いことになります。

5105 とはいえ検定ですから、ある有意水準で一方が他方より大きいと判断すると間違える確率が云々、という
 5106 判定に落とし込むことという点では同じでした。

5107 今回は生成量を使って、5 点以上大きくなる**確率**を求めることができました。これは検定統計量や p 値の
 5108 ように架空の値ではなく、もとのデータに直結した数字や確率になっていますから、解釈が比較的自然的にでき
 5109 るというのが利点です。プログラム上の表現も、たとえば今回のように `if` 文を使って表現するのは非常に直
 5110 感的で、「もしあれがこうなってそれがこうなったらどうなる？」と色々な仮説を考えていくこともできます
 5111 ね。帰無仮説検定のロジックは、結果判定のために仮説を限定的な型にはめ込んでしまいます。しかしその型
 5112 を飛び越えていろいろなことをやってもいいのだ、できるんだというのは、データがどうやって生まれてきてい
 5113 るかを考えている「創造主」としての特権かもしれません。

5114 ただし注意してもらいたいのは、これらの検証の仕方は、今回のデータと仮定されたモデルという前提の上
 5115 で成立する割合を確率とみなしているに過ぎない、ということです。データが変わればその数字も変わるで
 5116 しょうし、そもそもデータが正規分布から出てきていないのかもしれないかもしれません。本当は XX 分布から出てきてい
 5117 るのに、 YY 分布を仮定したモデルで計算したら、仮説が正しい可能性が 100% ! といっても虚しい (明らか
 5118 に間違っている) ことになりますね。心理学のデータの場合はえてして、どういう分布やメカニズムになるのか
 5119 がはつきりわからないものです。ただ、あまりにもわからなさ過ぎて、さまざまな影響が混ざり合っているの
 5120 で、色んな要素が相殺しあって結局ほとんど正規分布とみなしていいよ、ということだったりします。心という
 5121 目に見えないものに、不良設定問題を解くためのツールで立ち向かうという、無理に無理を重ねるような推論
 5122 の世界であるということに、自覚は持っていたいですね^{*6}。

^{*5} 二つの標準偏差の平均をとった、プールされた標準偏差で計算することもあります。

^{*6} いいんです、それでも。だってそんなよくわからない人間というのが好きなんですから。

20.4 今回のまとめ

我々がデータを分析する時は、そのデータの数字がどうであったか、ということを実先に考えたいはずで
 す。身長の違いは何センチあったのか、最終学歴によって生涯年収は何円差がつくのか、といった具体的な単
 位に基づく**実質的な差**が一番大事なはずで。心理学をやっていると、尺度をはじめとした「絶対的なスコア」
 が出てきませんから、尺度で何ポイント差があったのか、ということがあまり意味のある実態と対応しないこと
 も少なくありません。そう言った時のために、**標準化された差**を考えるようになったのです。その計算は非常
 に面倒ですが、機械がすぐに計算してくれることによって結論を急ぎすぎ、単位も実際のデータとも関係のな
 い統計量を参照して**有意差**を求めるようになってしまったのでは意味がありません。

データがどういうメカニズムで生まれてきたかを考え、母数をバイズ推定するというやり方は、いちいちコー
 ドを書かなければならないこともあって、非常に手間がかかるようにも思えます。しかし自分が何を仮定して
 いたのか、どういうメカニズムの元で考えていたのかを自覚し、またその推定値を使って自由に仮説を考える
 ことで、決まりきった分析手続きに落とし込むという単調な作業から解放され、クリエイティブにデータに向き
 合えることでもあります。

20.5 課題

次の計算をする R/Stan コードや回答を記述し、提出してください。なお提出されたコード単体でバグがな
 く動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜
 日別 TA か小杉までメールで連絡し、指導を受けてください。

■**フライドポテトの研究** あるコンビニで、ホットスナックのフレンチフライを 2 つ買ったところ、どうも一方
 より他方の方が長くて大きい気がした。そこでそれぞれの袋から取り出して長さを測定してみた (単位 cm)。
 測定結果が表 20.1 である^{*7}。このデータを用いて次の間に答えなさい。なおこの数値はシラバスのサイト上

表 20.1 二つの店舗のポテトの長さ

A	8.4	11.3	8.1	11.2	5.8	6.3	7.1	10.9	7.1	6.5	5.0	3.0	7.2	6.5	6.4	6.4	9.3	8.3
B	6.7	7.2	4.2	11.0	7.5	8.9	7.0	8.0	7.2	4.2	6.0	9.0	8.6	9.0	5.0			

のコードで取得可能です。

- ポテトの長さは正規分布に従うと仮定します。また両群の分布の幅は同じであると仮定します。t 検定
 で二群の平均値に差があるかどうか、判断してください。t 値や自由度、 p 値を参照しながら説明す
 ること。効果量も算出するとお良いです。
- ポテトの長さは正規分布に従うと仮定します。また両群の分布の幅は同じであると仮定します。その上
 で、群 A と群 B の平均値を推定するモデルを作り、バイズ推定してください。その上で、結果の平均
 値、中央値、MAP 推定値、90% 確信区間を報告してください。
- 両群のポテトの長さの平均が、3cm 以上違うようであればクレームをつけに行こうと思います。3cm
 以上の差がある確率はどれぐらいですか。推定してください。

^{*7} このデータは実際に筆者がコンビニの二つの店舗でポテトを購入し、長さを測定した結果に基づいています。

- 5152 4. よく考えてみれば, 両郡の分布の幅が同じであるという仮定はおかしいような気がしてきました。そこ
5153 で, それぞれの群毎に分布の幅を推定するモデルに作り替え, ベイズ推定してください。その上で, 結
5154 果の平均値, 中央値, MAP 推定値, 90% 確信区間を報告推定してください
- 5155 5. 異なる分散を指定したモデルで, 両群のポテトの長さの平均差が 5cm 以内であれば, クレームをつ
5156 げに行くのはやめておこうと思います。差が 5cm 以内である確率はどれぐらいですか。推定してくだ
5157 さい。

第 21 章

モデリングの目から見た検定 2 ; パラメータの世界とデータの世界

5161 前回は、二群の平均値差の検定を行うデータ生成モデルを考えました。平均値差の検定には正規分布が
 5162 仮定されますから、データ生成モデルも正規分布からデータが作られていると考えれば良かったわけです。
 5163 仮定に忠実に設計図を書き、それに対応した Stan コードを書くとも平均の推定値が出てきます。検定の時も
 5164 母平均の推定値を使って考えるのですが、検定統計量にしてしまうので実感が湧きにくいところがあったか
 5165 もかもしれません。データ生成モデルの場合は、直接パラメータを考えるのでイメージがしやすいという側面があ
 5166 りました。

5167 さらに、MCMC による推定ですから事後分布からの代表値が、実現値として手元のオブジェクトに代入
 5168 されています。これをつかった計算をすることで、たとえばパラメータの差を計算でき、その分布を考えること
 5169 ができるのでした。こうした計算は Stan の `generated quantities` ブロックを使うことで、MCMC の結
 5170 果と同時に考えることもできるのでした。

5171 さて今回も、この `generated quantities` ブロックをつかって、生成量からいろいろ考えてみましょう。

21.1 事後予測分布

5173 前回、二群のデータ $X_{i,A}$ および $X_{i,B}$ に対して、次のようなデータ生成モデルを考えました。

$$X_{i,A} \sim N(\mu_A, \sigma), X_{i,B} \sim N(\mu_B, \sigma)$$

5174 このパラメータ μ_A, μ_B, σ の推定値、 $\hat{\mu}_A, \hat{\mu}_B, \hat{\sigma}$ の分布からの代表値が MCMC によって得られるのでし
 5175 たね。MCMC サンプルは iteration ごとに $1, 2, 3, \dots, M$ 個得られたとすると、その期待値すなわち **EAP** は、
 5176 次の式で求めていることと同じです*1。

$$\hat{\mu}_{A\text{eap}} = \frac{1}{M} \sum \mu_A^i = \frac{1}{M} (\mu_A^1 + \mu_A^2 + \dots + \mu_A^M)$$

$$\hat{\mu}_{B\text{eap}} = \frac{1}{M} \sum \mu_B^i = \frac{1}{M} (\mu_B^1 + \mu_B^2 + \dots + \mu_B^M)$$

$$\hat{\sigma}_{\text{eap}} = \frac{1}{M} \sum \sigma^i = \frac{1}{M} (\sigma^1 + \sigma^2 + \dots + \sigma^M)$$

*1 ここで上つきの数字は MCMC のステップ番号を表しているもので、冪乗の数字ではありません。また M は Stan のデフォルト
 では 4000 です。

5177 さて、今回は得られたデータ $X_{i.}$ から μ, σ を推定したわけですが、これはいわばデータ生成メカニズムの
 5178 設定値を見つけたことと同じですね。例え話で考えるなら、こんな感じです。ある企業 A が製品 a を量産して
 5179 いるとします。ライバル会社 B が、この製品 a の類似品を作ろうと考えたとします。そこで B 社は A 社の製
 5180 品 a を作っている製造機械と同じ型番の機械を購入します。この機械には製品を生成するに当たっていくつ
 5181 かの設定をしなければならないとします。アナログな例えで恐縮ですが、ツマミが 2 つ 3 つ付いていて、それ
 5182 をひねれば何らかの製品ができるといった感じです。でも A 社が製品 a をつくるのに、どうい設定にしてい
 5183 るのかはわからない。そこで製品 a をいくつかサンプルとして入手します。入手したサンプル a_1, a_2, \dots と見比
 5184 べながら、この辺りかな？この辺りかな？とツマミを調節し、サンプル a_1, a_2, \dots と同じような製品 b_1, b_2, \dots が
 5185 できるのはこの辺の設定だな、というのを見つけたという感じ。

5186 この例え話を MCMC の用語を使っていうと、入手したサンプル a_1, a_2, \dots がデータです。このデータをつ
 5187 くる製造機械が**確率分布**です。機械には設定のツマミがついているんですが、そのツマミは**パラメータ**の
 5188 メタファーになっています。サンプルを参考に、この辺りかな、この辺りかな、とツマミの値を探るステップが
 5189 MCMC の各段階で、まずは $\mu_A^1, \mu_B^1, \sigma^1$ 、次にそれからちよっと動かして $\mu_A^2, \mu_B^2, \sigma^2$ 、またちよっと動かし
 5190 て・・・と 4000 回試すことで「だいたいこのデータを作るための設定はこの辺り」というのが決まってくるわけ
 5191 ですね。「この辺り」というのを点推定する場合は **EAP** や **MAP** などを使いますし、幅を持って推定したい
 5192 場合は**確信区間**で表現するのです。

5193 ということで、いくつかの手元のデータから、製造メカニズムの設定を手に入れました。このとき B 社はきつ
 5194 と、「じゃあこの推定値をつかって（模造品である）製品をジャンジャン作っていこう」となるでしょう。この推定
 5195 値から作られる新しい製品（模造品、新しいデータとも言えます）のことを MCMC ではとくに**事後予測分布**
 5196 (**posterior predictive distribution**) と言います。分布となっているのは、作られる製品も散らばりま
 5197 すので、その散らばり方を分布で考えるからです。

5198 事後予測分布は、点推定値を使って計算することもできます。つまり、データが $X_i \sim N(\mu, \sigma)$ というメカ
 5199 ニズムで作られていて、それから推定値 $\hat{\mu}, \hat{\sigma}$ を得たわけですから、これを使って $X_{new} \sim N(\hat{\mu}, \hat{\sigma})$ とすれ
 5200 ば良いのです。これはただの乱数生成でもありますから、たとえば R で `rnorm(mu, sigma)` とするとき `mu`
 5201 と `sigma` に推定値を入れてやればいいでしょう。でも点推定値では決め打ちが過ぎますので、区間推定した
 5202 上限と下限も入れますか？いや、それならいっそ、MCMC で得られた事後分布からの代表値を全部入れて
 5203 しまえばいいのではないのでしょうか。

5204 それを実現するために、`generated quantities` ブロックを使うといいでしょう。コード 21.1 にその例を
 5205 書いてみました。データやモデルは前回の二群の平均値差の検定と同じものです。

code : 21.1 事後予測分布を作るコード

```

5206 1 data{
5207 2     int<lower=0> N1; // Number of Subjects in Group 1
5208 3     int<lower=0> N2; // Number of Subjects in Group 2
5209 4     array[N1] X1; // Data in Group 1
5210 5     array[N2] X2; // Data in Group 2
5211 6 }
5212 7
5213 8 parameters{
5214 9     real mu1;
5215 10    real mu2;
5216 11    real<lower=0> sig1;
5217 12    real<lower=0> sig2;
5218 13 }
5219 14
5220

```

```

5221 15 model{
5222 16     // likelihood
5223 17     X1 ~ normal(mu1,sig1);
5224 18     X2 ~ normal(mu2,sig2);
5225 19     // prior
5226 20     mu1 ~ uniform(0,100);
5227 21     mu2 ~ uniform(0,100);
5228 22     sig1 ~ cauchy(0,5);
5229 23     sig2 ~ cauchy(0,5);
5230 24 }
5231 25
5232 26 generated quantities{
5233 27     real Xpred1[N1];
5234 28     real Xpred2[N2];
5235 29     for(i in 1:N1){
5236 30         Xpred1[i] = normal_rng(mu1,sig1);
5237 31     }
5238 32     for(i in 1:N2){
5239 33         Xpred2[i] = normal_rng(mu2,sig2);
5240 34     }
5241 35 }
5242

```

5243 ここで、generated quantities には、normal_rng という関数が入っています。この_rng は確率分布
5244 の後ろにつけて、乱数を発生させるという意味です。R でいう r****みたいなもんですね。このコードでは各
5245 群と同じサイズだけの事後予測データセットを作っています。まあ同じ μ と σ からの乱数ですから、データの
5246 数だけ作ったりしなくてもいいんですがイメージしやすくするためにそうしてみました。この μ と Xpred を図
 にしてみたのが図 21.1 です。この図の左側は μ_1 の事後分布が、右側には X_1 の事後予測分布が描かれ

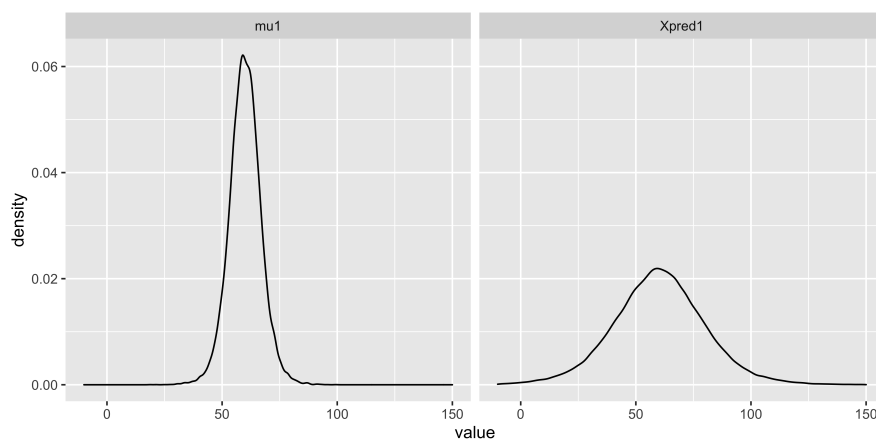


図 21.1 事後分布と事後予測分布

5247 ています。左側は**パラメータの世界**の話であり、右側は**データの世界**の話になっていることに注意してくだ
5248 さい。
5249

5250 パラメータ μ_1 は、EAP で 59.96,95%CI[48.81,71.22] となっています。そこからでてくるであろうデータ
5251 は、25 から 100 弱の幅に分布していますね。平均が 60 であっても、SD が 20 弱ありますから、データのレ
5252 ベルではその散らばる幅が大きくなるのも当然ですよ。でもこれはとても大事なことで、我々は「群間に差
5253 があった」となるとその効果について色々考えますが、それはパラメータの話、あるいは**平均因果効果**の話で

5254 あって、個々のデータではまた事情が違います。たとえば、男女差があると言っても平均的な男性、そうで
5255 ない女性などさまざまなケースがあるように。たとえば、平均的に 80% 成功する手術があると言われても、自
5256 分の身に置き換えると生きるか死ぬかは割合の問題ではないように。

5257 また、この事後予測分布の形は、もとのデータの分布の形と似ているかどうかを視覚的に確認するために
5258 使われます。図 21.2 に、事後予測分布ともとのデータのヒストグラムを合わせて表示してみました。事後予測
分布は MCMC サンプルの数だけありますが、ここでは一部だけ取り出しています。これを見ることで、もと

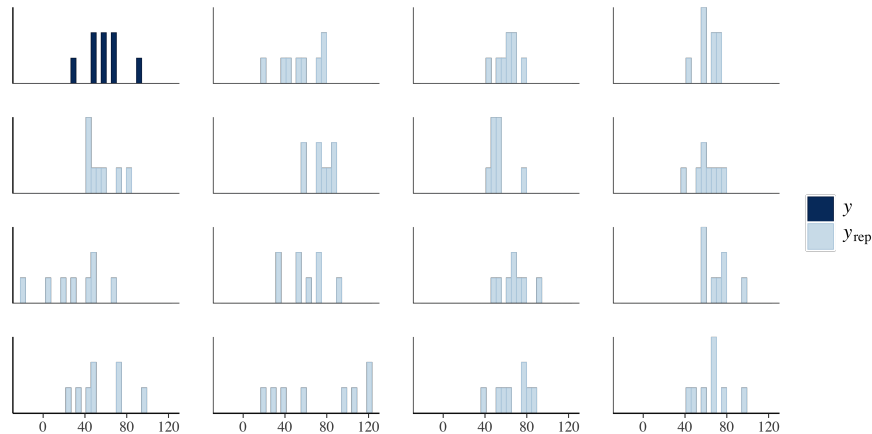


図 21.2 データの分布と事後予測分布の視覚的な確認

5259 のデータと事後予測分布の形が似ているかがわかります。もし模倣した機械の設定値が正しいのであ
5260 れば、手元にある商品サンプルと同じような分布をするデータ分布が得られるはずですね。視覚的な確認で
5261 すから、見て「うん似てるな、似てないな」というのを見て考えるだけではあります。今回はデータの数が 8
5262 件と少ないので何ともイメージしにくいところですが、データが多くなるとデータの分布の形状で比較すること
5263 もしやすくなるでしょう。もし分布の形が似ていなければ、購入した製造機械が違うものだったのかもしれな
5264 い、ということになります。というのも実際のデータ分析の場合は、製造機械（確率分布、モデル）が正しいか
5265 どうかはわからないので、それすらも仮定、検証すべき対象だからです。

21.2 データレベルの仮説

5268 **事後予測分布**の利用方法は他にもあります。事後予測分布は、確率モデルと推定されたパラメータを利用
5269 した「新しいデータの生成」だったわけですが、この新しく生まれるであろうデータについての、データの世界
5270 での仮説を考えることができます。

5271 帰無仮説検定で扱っていたのは、パラメータの世界の仮説でした。たとえば二群の平均値差の検定につい
5272 て言えば、帰無仮説は $\mu_1 = \mu_2$ であり、対立仮説は $\mu_1 \neq \mu_2$ です。これはギリシア文字 μ を使って書いて
5273 あることから明らかな通り、母数の仮説であり、 $\bar{X}_1 = \bar{X}_2$ のような、データについての仮説ではありません
5274 ね。ではデータについての仮説というのはどういうものがあり得るのでしょうか？たとえば次のようなものが考
5275 えられます。

- 5276 • 無作為に選んだ統制群のデータと実験群のデータを比べたとき、実験群が統制群を上回っている確
5277 率はどれくらいあるだろうか？
- 5278 • 無作為に選んだ統制群のデータと実験群のデータの差が、基準点 c より大きくなる可能性はどれくら
5279 いあるだろうか？

5280 第1の点は**優越率 (probability of dominance)**という指標です。これはパラメータの世界の比較で
 5281 はなく、実データの比較ですからイメージしやすいのではないのでしょうか。たとえば男女差が平均的にはある
 5282 とされているとしても、それはパラメータの違いであって、実際これから新しく出会う男性(女性)が自分より
 5283 優れている(劣っている)可能性は、と考えた方が実質的な意味がありそうです。性差の話の多くは、平均値
 5284 差だけでしか議論されませんが、効果量^{*2}で考えると小さいことが少なくありません。それでも有意差が検出
 5285 されてしまうのは、非常に大きなサンプルサイズをとっているから、ということもあつたりします。サンプルサイ
 5286 ズを大きくするのは研究者の工夫や努力なのですが、小さな効果を取り上げて針小棒大に騒ぎ立てるのは科
 5287 学的に良い態度とは言えません。そんな時に優越率で考えてみると、差があると言っても無作為に二群から
 5288 新しい情報を得た場合の優越率が50% ちよつとしかない、つまり超えるか超えないかが半々ぐらいですから
 5289 「何の情報ももたらさない」と同じこと、ということになつたりします。くどいようですが、実データのレベルで
 5290 の比較になるので、群間の差の表現としてよりわかりやすいということが言えるでしょう。第2の点は、優越率
 5291 に一定の違いの差 c を含めて考えるもので、**閾上率 (probability beyond threshold)** といいます。こ
 5292 れもデータレベルでの仮説の1つで、たとえば新しいトレーニング方法を使えばタイムが3秒小さくなる確率
 5293 は $X\%$ 、などということが分かれば具体的にイメージしやすい伝達方法になると思いませんか。

5294 ちなみにこの優越率、閾上率などの指標は古くから指摘されてきていたものであり^{*3}、ベイズ推定を使わな
 5295 い、**モーメント法**による推定でも計算可能な量でした。まあしかし、このような統計量が実際に使われている
 5296 ケースってほぼ見かけないんですけどね。しかし、これらの指標はベイズ推定を使うことによって、というより
 5297 乱数による近似計算をすることによって、遥かにイメージしやすくなっています。ここで示された数値はいずれ
 5298 も、generated quantities ブロックの中に if 文をつかって書いてやれば表現できることだからです！

5299 実際に考えてみましょう。これら「条件が成立する確率」については、成立すれば1、成立しなければ0とい
 5300 うフラグをたててMCMC サンプルを生成し、その平均値すなわち相対頻度で考えてやれば良いのでした。

code : 21.2 優越率や閾上率を計算するコード

```

5301 1 data{
5302 2     int<lower=0> N1; // Number of Subjects in Group 1
5303 3     int<lower=0> N2; // Number of Subjects in Group 2
5304 4     array[N1] real X1; // Data in Group 1
5305 5     array[N2] real X2; // Data in Group 2
5306 6     int<lower=0> C; //constant
5307 7 }
5308 8
5309 9 ...
5310 10
5311 11 generated quantities{
5312 12     real Xpred1;
5313 13     real Xpred2;
5314 14     int<lower=0, upper=1> FLG1;
5315 15     int<lower=0, upper=1> FLG2;
5316 16     Xpred1 = normal_rng(mu1, sig1);
5317 17     Xpred2 = normal_rng(mu2, sig2);
5318 18     //probability of dominance
5319 19     if(Xpred1 > Xpred2){
5320 20         FLG1 = 1;
5321 21     }else{
  
```

*2 標準化された平均値差のことでしたね。

*3 たとえば南風原・芝 (1987) による優越率の解説は1987年の論文です。

```

5323 22     FLG1 = 0;
5324 23     }
5325 24     //probability beyond threshold
5326 25     if(Xpred1 - Xpred2 > C ){
5327 26         FLG2 = 1;
5328 27     }else{
5329 28         FLG2 = 0;
5330 29     }
5331 30 }
5332

```

コード 21.2 は、generated quantities ブロックで 2 つの FLG 変数を用意しています。FLG1 が優越率、FLG2 が閾上率のために用意したものです。このブロックで、推定されたパラメータ $\mu_1, \mu_2, \text{sig}_1, \text{sig}_2$ からそれぞれ生成される二群の「新しいデータ」を作り (Xpred1, Xpred2), 成立するかどうかを検証します。FLG1 は if 文の中にあるように、 $X_{\text{pred1}} > X_{\text{pred2}}$ が成立すれば 1, しなければ 0 になる数字です。FLG2 は、同じくその差分が一定の値 C を超えていれば成立, 超えていなければ不成立というフラグです。ちなみにこの値 C は data ブロックで外部から読み込むことになっています。

これと前回のデータを使って推定させてみましょう。

code : 21.3 優越率, 閾上率を推定するコード

```

5340
5341 1 groupA <- c(30, 50, 70, 90, 60, 50, 70, 60)
5342 2 groupB <- c(20, 40, 60, 40, 40, 50, 40, 30)
5343 3 dataSet <- list(X1 = groupA, X2 = groupB, N1 = 8, N2 = 8, C = 3)
5344 4 ### rstan の場合
5345 5 sampling1 <- sampling(modelR, dataSet, warmup = 1000, iter = 4000, chains = 4)
5346 6
5347 7 ### cmdstanr の場合
5348 8 sampling2 <- modelC$sample(
5349 9     data = dataSet,
5350 10     chains = 4,
5351 11     iter_sampling = 5000,
5352 12     iter_warmup = 1000,
5353 13     parallel_chains = 4
5354 14 )
5355

```

MCMC の結果 4

```

# A tibble: 8 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 FLG1      0.800      1.000      0.999      0.400      0.000      1.000
2 FLG2      0.764      1.000      1.001      0.425      0.000      1.000
3 mu1      59.933     59.936     59.833     6.914     46.147     73.895
4 mu2      40.006     40.002     40.141     4.656     30.869     49.362
5 sig1     18.598     17.495     15.846     5.480     11.310     32.054
6 sig2     12.640     11.941     11.071     3.725     7.623     21.903
7 Xpred1    59.880     59.828     58.098     20.590     18.660     100.984
8 Xpred2    39.887     39.978     40.457     13.984     11.642     67.691

```

ここにあるように、優越率は 80.0%, 閾上率は 76.4% ということになりました。実験群 (GroupA) のほうが

5358 統制群 (GroupB) よりもスコアが高くなる確率が 80%, スコアが 3 点以上大きくなる確率は 76% もあるの
5359 ですから, 実験群の効果は結構おおいぞ, といったことがわかるようになります。

5360 こうしたデータに基づく仮説は, 単位がもとのデータに対応しているから想像しやすく, より実践的な意味
5361 を想像しやすいですね。

5362 まとめておきます。パラメータの世界に関する仮説は, 推定されたパラメータはどれぐらいの値なのか, どれ
5363 ぐらいの「差」「違い」があるのかを考えることでした。あるいは, 推定されたパラメータはどれぐらいの幅に分
5364 布するのか, という差の分布を考えることで, 効果の有無の確からしさを量的に評価できるのでした。デー
5365 タの世界に関する仮説は, データ生成機がどういうデータを作るのか, ひいてはそれが今のデータと同じ形を
5366 しているかを見ることで, モデルが正しいかどうかを検証できます。また未来のデータの分布はどんな形にな
5367 るのか, という意味では, データという具体的な数字でもって未来予測をできます。

5368 ただしいずれの世界においても注意しておいて欲しいのは, データ生成機すなわち確率モデルが正しい場
5369 合であり, これが間違っていると generated quantities などで作られる数字もまったく意味のないもの
5370 になってしまう, ということです。もちろん **t 検定**など帰無仮説検定を行う時も, データが正規分布していなけ
5371 れば意味のない検定結果になるので, 「そもそも統計的な仮定が間違ってた」というのでは意味がないのはい
5372 ずれも同じです。推測統計学は手元の少しのデータから, 未知の世界を予測検証するという不良設定問題に
5373 対応するためのものですから, 仮定や前提が正しいかどうかについては常に敏感でなければなりません。

5374 21.3 パラメータ・リカバリ

5375 そもそも仮定が間違っていた, というのを考えすぎて何もできなくなっても良くないですから, そこはいつ
5376 たん OK だとしましょう。そう考えると, 仮定・前提をつけた上での話になりますが, 未知の世界や未来を予測
5377 できるというのは, とんでもなく強力なツールであるような気がします。株価や競馬の予想ができれば, 大儲
5378 けできる薔薇色の人生が待っているかも。ということで, バイズ推定の勉強にも身が入るというものです。推
5379 定を可能にしてくれた MCMC に感謝!

5380 半ば冗談, 半ば本気なのですが, Stan や JAGS などの**確率的プログラミング言語 (stochastic**
5381 **programming language)** が登場したおかげで, 事後分布の計算が可能に, 簡単便利にできるよう
5382 になったというのは大変ありがたいことです。これのおかげでどんなモデルでも, 設計図が書ければ答えが出せ
5383 ちゃう。複雑な確率の式とか微分積分を考えなくても, なんとかなるような気がします。しかし簡単なツールが
5384 出てきたときの常で, ちょっと注意が必要なのですが, MCMC は優秀すぎて正しくない計算式からでも答え
5385 が出てしまう, という可能性があるのです。

5386 MCMC がうまく行ったかどうかは, **Rhat** や**有効サンプルサイズ**, **トレースプロット**などをチェックすれ
5387 ばよい, というのはすでにお話した通りです。しかしここで指摘したいのはそうではなく, これらの MCMC
5388 がうまく行っていたとしても, 推定結果が正しいとは限らない可能性です。MCMC は「あり得そうな答えの
5389 可能性」を大量に出力し, その分布で我々は考えるのですが, その答えの可能性群が全部的外れなところ
5390 に行ってるかもしれないのです。実際にデータを分析していると感じるのですが, MCMC はかなり複雑なモデ
5391 ルを考えても答えを出してくれます。とにかく答えてくれる機械が手に入ったとしても, その信憑性をチェック
5392 することを忘れてはいけません。

5393 ではこれをチェックするためにはどうすればよいのでしょうか。ひとつは, 答えが先にわかっている問題を作っ
5394 て, 機械が正解を当てられるかどうかを検証するということが考えられます。この「わかっている正しい値」を
5395 当てることができるのであれば, 実際のデータを用いた「わからない値」を推測する方法として有用だとい
5396 うことがわかります。正解を当てることができないようであれば, この計算機械は当てにならないので, 実際の
5397 データを使っても正しく推測できているとは言えないでしょう。

5398 この検証方法のことをパラメータリカバリ (parameter recovery) と言います。実際のデータ分析を
 5399 する場合は、いきなり未知の答えに推測機をあてがうのではなく、推測機が事前に設定したパラメータをリカ
 5400 バリ (回復) できることを確認してから利用するようにしましょう。実際の研究実践の場合、その分析手順は、
 5401 以下ようになります。

- 5402 1. データ生成メカニズムを考える (紙とペンによる設計図の作成)
- 5403 2. モデルを式で表現する
- 5404 3. モデルの式からデータ生成をシミュレーション
- 5405 4. Stan でコードを書く
- 5406 5. シミュレーションデータを Stan が再現していることを確認する
- 5407 6. 実際のデータでやってみる
- 5408 7. 結果の解釈 (図による確認も)

5409 ここで第 3 のステップ、仮想データの生成から第 5 のステップ再現の確認が、パラメータリカバリという作業
 5410 になります。推定したい平均値などのパラメータを事前に仮に入れてみて、その答えが復元できるかどうかを
 5411 見るのですね。

5412 なんだかまるでっかしい! という人もいるかもしれません。自分で設定したパラメータ通りになるのって当たり
 5413 前じゃないか、何がおもしろいのか? というわけですね。ですが意外と当たり前でないことがあるんです。間
 5414 違った答えであれば、実践では役に立ちません。あるいは、いきなりデータでいいじゃない、固いこと言うな
 5415 よ、と思う人もいるかもしれません。これが設定ミスの場合には答えを出さない・答えられないという機械であれ
 5416 ばそれでもいいのかもしれませんが、MCMC は大概なんでも答えてしまうので、いきなりの運用は怖いと
 5417 思いませんか。またあるいは、基礎実験とかそのほかの授業でも、今までそんなのやったことないよ、という意
 5418 見もあるかもしれません。大学の授業で行うような演習・実習は、そもそも効果をはっきりあること、再現する
 5419 現象であることがある程度確からしいものを選んでやっていますからいいのですが、これから皆さんは研究
 5420 者の卵として、新しい現象、十分に確認されていない現象を研究対象にすることになります。そのとき、いきな
 5421 り出てきた結果がどれほど信用できるでしょうか。実際、最近では心理学実験が再現できないケースも数多く
 5422 指摘されており、その理由のひとつとして統計モデルの乱用 (悪用・誤用) にあったことは何度もお伝えしてき
 5423 た通りです。さらに、ごく微小な効果を偶然検出してしまった、という可能性もあります。事前に推定機の精度
 5424 がわかっていたらこうした問題を回避できるのです。この話はベイズ推定をする場合だけでなく、帰無仮説検
 5425 定など従来のやり方にも通じ、効果量や例数設計の話に直結します。

5426 21.3.1 二群の平均値差の例

5427 ここまで扱ってきた、二群の平均値差を検証する例で考えてみましょう。もともとの設計図 (図 20.1) がここ
 5428 でも役に立ちます。データがどういうメカニズムで出てきたかを考えていたわけですから、これを使えばデー
 5429 タ生成を考えることができますね。乱数生成のアプローチはデータを取り始める前にも有用なのです!

5430 今回はパラメータを仮に設定してやることができます。二群の平均とその差を次のように決めてしまいま
 5431 しょう。

code : 21.4 平均値差の設定

```
5432
5433 1 mu1 <- 50
5434 2 diff <- 10
5435 3 mu2 <- mu1 + diff
5436 4 sig1 <- 5
5437 5 sig2 <- 8
```


5438

5439 第一の群の平均は `mu1 <- 50` としています。第二群の平均は `mu2 <- 60` というようにしてもいいので
 5440 すが、差を考えるという意味を強調するためにいったん変数 `diff` を作って差の量を定め、それを `mu1` に加
 5441 えるという方法をとっています。2つの群は正規分布に従いますから、それぞれの標準偏差を `sig1, sig2`
 5442 として設定しました。ここから仮想データを作ります。データの生成は乱数に従うので、`rnorm` 関数の出番
 5443 です。

code : 21.5 仮想データの生成

5444

5445

5446

5447

5448

5449

```
1 set.seed(12345)
2 N <- 10
3 X1 <- rnorm(N, mu1, sig1)
4 X2 <- rnorm(N, mu2, sig2)
```

5450 これで実際に生成された仮想データは次のようになっています。

R の出力 21.1: 検定の結果

```
> X1
[1] 52.92764 53.54733 49.45348 47.73251 53.02944 40.91022 53.15049
[8] 48.61908 48.57920 45.40339
> X2
[1] 59.07002 74.53850 62.96502 64.16173 53.99574 66.53520 52.90914
[8] 57.34738 68.96570 62.38979
```

5451

5452 ちなみにこの2つのデータセット、平均は $\bar{X}_1 = 49.33528$, $\bar{X}_2 = 62.28782$ となっています。理論的には
 5453 それぞれ 50, 60 であるはずなのですが、乱数による実現値なので少し誤差が出ていますね。これは実際の研
 5454 究状況でもある話です。理論的な値と、目の前の実現値というのは必ず少しずれがあるもの。その中で、きち
 5455 んと理論値を推定できるかどうかが問題なのです。さて、このデータを使って分析していきましょう。

code : 21.6 パラメータリカバリ・検証

5456

5457

5458

5459

5460

5461

5462

5463

5464

5465

5466

5467

5468

5469

```
1 ## t検定(モーメント法による推定と判定)
2 t.test(X1, X2)
3 ## ベイズ推定(MCMCによるベイズ推定と差の分布)
4 dataSet <- list(X1 = X1, X2 = X2, N1 = N, N2 = N, C = 3)
5 modelC <- cmdstan::cmdstan_model("ttest03.stan")
6 sampling2 <- modelC$sample(
7   data = dataSet,
8   chains = 4,
9   iter_sampling = 5000,
10  iter_warmup = 1000,
11  parallel_chains = 4
12 )
```

5470 まずは `t` 検定の結果ですが、ここでは $t(14.797) = 5.2059$, $p = 0.0001113$ で有意差ありと判定されまし
 5471 た。実際 `diff <- 10` で差がある設定なのですから、正しく判定できて良かったな、というところです。ベ
 5472 イズ推定の方はどうでしょうか。結果は次のようになりました(ここでは前回のコード 20.6 を再利用しました)。

MCMC の結果 5

```
# A tibble: 6 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 diff   -12.936  -12.957  -13.250   2.762  -18.458  -7.445
2 FLG     0.000   0.000   -0.001   0.000   0.000   0.000
3 mu1    49.349  49.350  49.499   1.477  46.381  52.302
4 mu2    62.285  62.285  62.256   2.355  57.600  66.977
5 sig1    4.482   4.257   3.868   1.194   2.842   7.406
6 sig2    7.189   6.870   6.363   1.821   4.615  11.579
```

5473

5474 これを見ると, $\mu_1 = 50$ と設定したときの EAP 推定値が 49.349, 95%CI で [46.381, 52.302] です
 5475 から, ほぼ正しい点推定値, また確信区間の中に真値を含んでいます。推定はうまく行っている, というこ
 5476 事です。 μ_2 についても同様で, 真値が 60 であり, EAP 推定値が 62.285, 95%CI で [57.600, 66.977] で
 5477 す。点推定値は真値から 2 点ほどずれていますが, 95% の区間推定にすれば真値を含んでおり, 推定はう
 5478 まく行っているようです。標準偏差の推定値についても, それぞれうまく行っており, 差についても EAP 推定
 5479 値が -12.936, 95%CI で [-18.458, -7.445] です*4。確かに確信区間に真値は含んでいるのですが, 17.45
 5480 は真値 10 から少し外れすぎかな... というようなことがわかりますね。

5481 このように, 今回の結果をみると概ね正しく推定できていると言えるでしょう。その精度についても「どの程
 5482 度誤差が生じるものなのか」が具体的にわかったことと思います。しかし次のような設定の場合はどうで
 5483 しょう。

code : 21.7 仮想データの生成その 2

5484

```
5485 1 mu1 <- 50
5486 2 diff <- 18
5487 3 mu2 <- mu1 + diff
5488 4 sig1 <- 10
5489 5 sig2 <- 15
5490 6 N <- 3
5491
```

5492 この設定でデータを作って, 推定した結果は次のようなものです (出力 6)。

MCMC の結果 6

```
# A tibble: 6 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 diff   -5.434  -5.630  -5.104  12.855  -31.789  21.966
2 FLG     0.171   0.000   0.000   0.376   0.000   1.000
3 mu1    53.962  53.972  54.117   4.147  45.756  62.294
4 mu2    59.396  59.671  59.608  12.124  32.811  84.278
5 sig1    6.194   5.034   3.778   4.332   2.356  17.443
6 sig2    20.651  17.148  13.108  12.966   8.682  53.683
```

5493

5494 MCMC の収束は問題なく, 推定値は出力されていますが, 本来 $\mu_1 = 50$ であるところが 53.962, 95%
 5495 確信区間は [45.756, 62.294] であり, $\mu_2 = 68$ であるべきところに至っては EAP 推定値が 59.396, 95%

*4 生成量は $\text{diff} = \mu_1 - \mu_2$ で作っているのので, 正負が逆転しているのはご容赦ください。

5496 確信区間は [32.811, 84.278] です。確信区間はなんとか平均値を挟み込みましたが、38 から 84 までと
 5497 46 点も幅がある推定はあまり有効な予測とは言えませんね。差は 18 あるという設定なのですが、EAP
 5498 推定値はわずか 5.434 であり、95% 確信区間は [-31.789, 21.966] です。確信区間にゼロをふくんでい
 5499 ますから、「差がないかも」ということになります。実際、t 検定の結果は帰無仮説を棄却できませんでした
 5500 ($t(2.2337) = 0.52828, p = 0.6451$)。タイプ 2 エラー (Type II Error) を犯してしまっています。

5501 どうしてこんなにうまくいかなかったのでしょうか。今回は正規分布という正しい確率モデルを使っていたの
 5502 に、です。その答えは明らかに、 $N < 3$ という点にあります。つまり、データが 3 件しかないので予測の精度
 5503 が悪くなったのです。このように、モデルが正しくても推定精度が十分ではない、あるいは判定をするときに間
 5504 違った判断をする可能性がある、ということがわかります。推定精度や効果の判定については、サンプルサイ
 5505 ズ、データの散らばり (群内分散)、データ間の差の大きさ (群間分散) などの要因が影響してきます。ここで
 5506 設定値を色々変えてみて、どういう関係にあるのかをみておくのも良いでしょう。

5507 このように仮想データを使ったパラメタリカバリを通じて、推定や検定の精度を設計したりチェックしたり
 5508 できることがわかりました。これを踏まえて考えると、前回の例では $N = 8$ のデータで分析をしましたが、
 5509 $N = 8$ だと平均値差を 8~12 点ほど外した予測をしてしまうかもしれない、ということがわかります。逆に
 5510 「5 点差までの推定誤差に収めたい」というような希望があった場合、データは何件ぐらい取れば良いでしょ
 5511 うか? このような問いについても、 N の設定値をさまざまに変えてみることで、実験や調査を実際に始める
 5512 前に予想を立てておくことができます。これが例数設計です。例数設計の重要性を確認するとともに、こうした
 5513 乱数を作るアプローチで簡単に実践できますので、実際の研究の前にはぜひ一度試してみてください。

5514 21.4 今回のまとめ

5515 生成量をつかうことで、次のような検証をできるのです。

- 5516 • 事後予測分布を作ることで、想定した確率モデルが正しかったかどうか検証する。
- 5517 • 事後予測分布を利用して、優越率や閾上率などデータに基づく仮説を検証する。

5518 さらに、データ生成モデルを応用して、そもそも擬似データをつくることで推定モデルの精度や、どれぐらいの
 5519 サンプルサイズが必要かといったことを、実践の前に検証できます。

5520 データ生成アプローチと乱数発生技術の組み合わせで、より慎重かつ確実な研究アプローチができること
 5521 を理解してください。

5522 21.5 課題

5523 次の計算をする R/Stan コードや回答を記述し、提出してください。なお提出されたコード単体でバグがな
 5524 く動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜
 5525 日別 TA か小杉までメールで連絡し、指導を受けてください。

5526 ■フライドポテトの研究その 2 あるコンビニチェーンの 2 つのお店 A,B で、それぞれホットスナックのフ
 5527 レンチフライを買ったところ、どうも一方より他方の方が長くて大きい気がした。そこでそれぞれの袋から取り
 5528 出して長さを測定してみた (単位 cm)。測定結果が表 21.1 である^{*5}。このデータを用いて次の間に答えなさい。
 5529 なおポテトの長さは正規分布に従い、また両群の分布の幅は同じであると仮定します。ちなみに、この数
 5530 値はシラバスのサイト上のコードで取得可能です。

*5 このデータは実際に筆者がコンビニの二つの店舗でポテトを購入し、長さを測定した結果に基づいています。

表 21.1 二つの店舗のポテトの長さ

A	8.4	11.3	8.1	11.2	5.8	6.3	7.1	10.9	7.1	6.5	5.0	3.0	7.2	6.5	6.4	6.4	9.3	8.3
B	6.7	7.2	4.2	11.0	7.5	8.9	7.0	8.0	7.2	4.2	6.0	9.0	8.6	9.0	5.0			

- 5531 1. 店舗 A のポテトの平均が、店舗 B のポテトの平均より長い確率はどれくらいあるか。
- 5532 2. 今後、店舗 A で購入するポテトが店舗 B で購入するポテトよりも長い確率はどれくらいあるか。
- 5533 3. 今後、店舗 A で購入するポテトが、店舗 B で購入するポテトよりも 2cm 以上ながい確率が 30% あ
- 5534 れば店舗 B では買わないようにしようと思う。店舗 B を利用すべきかどうか、推定に基づいて判定し
- 5535 てください。
- 5536 4. 今後、店舗 A で購入するポテトと店舗 B で購入するポテトの長さの違いが 5cm 未満である確率が
- 5537 80% であれば、家により近い店舗 B を利用してもいいように思う。店舗 B を利用すべきかどうか、推
- 5538 定に基づいて判定してください。
- 5539 5. このコンビニチェーンのポテトの長さ平均を、できるだけ正確に推定したいと思う。標準偏差は今回
- 5540 の MAP 推定値を使うこととして、何本くらいポテトのサンプルを集めればその 95% 確信区間を
- 5541 0.05cm 以内に収めることができるだろうか。シミュレーション結果に基づき、おおよその数字を答えて
- 5542 ください。

第 22 章

モデリングの目から見た検定 3 ; 多群の 平均値差モデル

22.1 要因計画モデル

さてここまでは二水準の平均値差を考えるモデルでしたが、ここからそれを三水準に増やしてみましょう。水準数が多くなると**要因計画 (Factorial Design)** や**実験計画 (Experimental Design)** といった話になってきます。**分散分析 (ANalysis Of VAriance; ANOVA)** は要因計画と帰無仮説検定の合わせ技で、効果がある/ないの判定を分散の比をつかって行うことでした*1。

今回は要因計画をデータ生成モデルから考え、**ベイズ推定 (Bayesian Inference)** をすることになります。**帰無仮説検定**のような Yes/No 判断はせず、平均値や差の大きさを直接見積もることになるのは、二群の時と同じです。

まずは簡単な**群間計画 (Between Design)** の話から進めましょう。対応のない二群の時のように、それが三群になるモデルですので、モデルの設計図は非常に単純です (図 22.1)。設計図を見れば、コードはそ

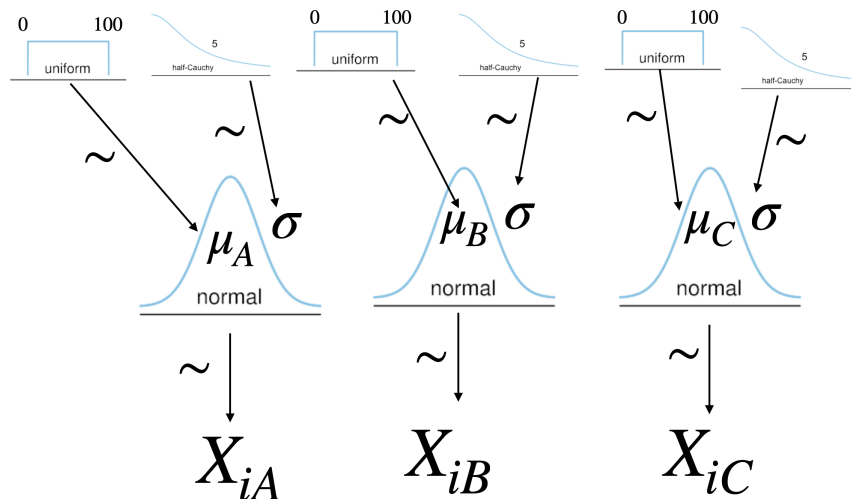


図 22.1 多群の平均値の差を比べるときの設計図

*1 忘れた人は、データ解析基礎のテキストに戻って確認してください。

5556 のまま置き換える形でかけますね。設計図の通りに書いてみたコードが code:22.1 にあります。

code : 22.1 三群の平均値のコード

```

5557
5558 1 data{
5559 2     int<lower=0> N1;
5560 3     int<lower=0> N2;
5561 4     int<lower=0> N3;
5562 5     array[N1] real X1;
5563 6     array[N2] real X2;
5564 7     array[N3] real X3;
5565 8 }
5566 9
5567 10 parameters{
5568 11     real mu1;
5569 12     real mu2;
5570 13     real mu3;
5571 14     real<lower=0> sig;
5572 15 }
5573 16
5574 17 model{
5575 18     // likelihood
5576 19     X1 ~ normal(mu1,sig);
5577 20     X2 ~ normal(mu2,sig);
5578 21     X3 ~ normal(mu3,sig);
5579 22     // prior
5580 23     mu1 ~ uniform(0,100);
5581 24     mu2 ~ uniform(0,100);
5582 25     mu3 ~ uniform(0,100);
5583 26     sig ~ cauchy(0,5);
5584 27 }
5585 28
5586 29 generated quantities{
5587 30     real diff12;
5588 31     real diff13;
5589 32     real diff23;
5590 33
5591 34     diff12 = mu1 - mu2;
5592 35     diff13 = mu1 - mu3;
5593 36     diff23 = mu2 - mu3;
5594 37 }
5595

```

5596 最後の generated quantities ブロックにあるように、各群の差の分布も計算して出せるようにしてあ
5597 ります。これをみながら、群 A と B, A と C, B と C などどこに差があるのか、その大きさはどれぐらいかと
5598 いうことを考えることができます。左の大きさを標準偏差 sig で割ることで効果量を出すこともできます。

5599 二群の平均値差の時もそうでしたが、検定と違って差の大きさを直接検証できているところがポイ
5600 ントです。また、検定の場合は分析によって有意差が確認できた場合、そのあと**多重比較 (multiple
5601 comparison)**へと進むのです。これは帰無仮説が $H_0 : \mu_1 = \mu_2 = \mu_3$ となっていて、これを否定する
5602 ことから「どこに差があったか」をさらに詳しくみる必要があったからです。この時、統計的検定の繰り返しに
5603 よって**タイプ 1 エラー (Type I Error)**が増えること、すなわち $\alpha = 0.05$ と 5% 水準で検定していても、
5604 検定を繰り返すと 5% を大きく上回る問題が生じるため、有意水準を調節するなど工夫が必要だという話で

5605 した*2。

5606 ところが、ベイズ推定の場合はこうした問題が生じません。なぜなら、確率を伴う判断をしていないからで
5607 す。帰無仮説検定の場合は、「帰無仮説が正しいという条件のもとで」計算された検定統計量 (F 値や t 値)
5608 が、帰無仮説の条件下ではどれほど生じやすいかということを考え、その確率 p 値を計算していましたが、こ
5609 の p 値はパラメータがどこにあるかといった確率ではありません。判断のための指針としての確率に過ぎない
5610 のです。一方ベイズ推定する場合、推定されたパラメータの**事後分布**は、パラメータがこの辺りにあるだろう、
5611 という確率であり、生成量である差の分布もきの大きさがこれぐらいだろうという範囲、大きさに直接関わる
5612 確率です。ですから何 % 区間で考えてもいいですし、その判断が間違うかどうかの確率というのは問題に含
5613 まれないのです*3。検定の多重性の問題が生じない、というのは複雑な要因計画を考える場合には非常に助
5614 かりますね。

5615 22.2 パラメータの変形と制約

5616 22.2.1 パラメータの変形

5617 ところで、先ほどは差の分布を generated quantities ブロックで表現しましたが、違う表現も可能で
5618 す。差分も未知なるパラメータと考えると推定するのです。

5619 順を追って説明しましょう。まず、最初の群の平均値を μ_1 とします。第二の群平均 μ_2 は、 μ_1 と δ_1 だけ違
5620 う、つまり $\mu_2 = \mu_1 + \delta_1$ と考えるのです*4。同様に、第三の群平均 μ_3 は $\mu_3 = \mu_1 + \delta_2$ と考えます。こうす
5621 ると、推定すべきパラメータは $\mu_1, \delta_1, \delta_2$ であり、generated quantities ブロックを使わなくても、直接
5622 差分パラメータを推定できるようになります。これを書いてみたのが、code:22.2 になります。

code : 22.2 差分を直接推定するコード

```
5623 1 ...
5624 2 parameters{
5625 3   real mu1;
5626 4   real delta1;
5627 5   real delta2;
5628 6   real<lower=0> sig;
5629 7 }
5630 8
5631 9 model{
5632 10  // likelihood
5633 11  X1 ~ normal(mu1,sig);
5634 12  X2 ~ normal(mu1+delta1,sig);
5635 13  X3 ~ normal(mu1+delta2,sig);
5636 14  // prior
5637 15  mu1 ~ uniform(0,100);
5638 16  delta1 ~ normal(0,100);
```

*2 お前は何をいつてるんだ、と思った人は、データ解析基礎のテキストに戻って確認しておいてください。大雑把に解説しますと、5% 水準というのは一回の判断で間違いが含まれる可能性を 5% にしましょう、という判定基準のことでしたが、二回判定を繰り返すと $1 - (0.95^2) = 0.0975$ と 9.75% 水準になってしまうということでした。3 群の平均値を比べるときは、3 箇所比較する必要があり、**t 検定**を 3 回も繰り返すと 5% 水準ではなくなる、という問題です。

*3 もっとも、すべて研究者の仮定したモデルのもとで、という話ではありますから、その仮定が間違っていたら全部意味のない数字である、ということにはなるかと思えます。とはいえ、帰無仮説検定もデータが正規分布に従うという仮定で行うのですから、これについてはどっこいどっこい、というところでしょうか。

*4 ここで δ はデルタと読む、ギリシア文字の d のことです。先ほどは生成量として `diff` と書いていましたが、ここでは未知のパラメータなのでギリシア文字に書き換えました。


```

5640 17   delta2 ~ normal(0,100);
5641 18   sig ~ cauchy(0,5);
5642 19   }
5643

```

5644 data ブロックに違いはありませんので省略してあります。変更点として、まず parameters ブロックで推
 5645 定すべきパラメータが変わっていることを確認してください。次に model ブロックで、3つのデータがそれ
 5646 ぞれ $\text{normal}(\mu_1, \text{sig})$, $\text{normal}(\mu_1 + \delta_1, \text{sig})$, $\text{normal}(\mu_1 + \delta_2, \text{sig})$ から出てきているよ
 5647 うに書き変わっています。正規分布の位置パラメータをずらすことで書き換えているのです。最後にその下、
 5648 事前分布の設定ですが、差分 δ_1, δ_2 は正負どちらにも出てくる可能性があり、その大きさがわかりませんの
 5649 で、平均を 0 にした正規分布としてあります。これで推定した結果は、先ほどの生成量を使った結果と変わり
 5650 はありません。数式の展開で書き方が変わっただけで、モデルの形は同じだからです。

5651 ここで transformed parameters ブロックの紹介をしておきましょう。code:22.2 は悪くはないのです
 5652 が、モデルのところ、 $\text{normal}(\mu_1 + \delta_1, \text{sigma})$ としているのがちょっと美しくないですね。元のコード
 5653 にあった、 $\text{normal}(\mu_2, \text{sigma})$ のほうがわかりやすいのに、と思った人もいると思います。まあ今の所、そ
 5654 こまで複雑ではないのでそれほど問題にはならないでしょうが、これからもしパラメータが複雑な構造をする
 5655 ようになったとき、尤度のところはもう少しスッキリ書きたいということも生じてくるでしょう。そういうときに
 5656 transformed parameters ブロックを使います。このブロックは、parameters ブロックの次に書くもの
 5657 で、言葉の通りパラメータを変形 (トランスフォーム) するためのものです。

5658 これを使って書いたのが code:22.3 です。

code : 22.3 パラメータ変換を入れたコード

```

5659 1   ...
5660 2   parameters{
5661 3     real mu1;
5662 4     real delta1;
5663 5     real delta2;
5664 6     real<lower=0> sig;
5665 7   }
5666 8
5667 9   transformed parameters{
5668 10    real mu2;
5669 11    real mu3;
5670 12    mu2 = mu1 + delta1;
5671 13    mu3 = mu1 + delta2;
5672 14  }
5673 15
5674 16  model{
5675 17    // likelihood
5676 18    X1 ~ normal(mu1, sig);
5677 19    X2 ~ normal(mu2, sig);
5678 20    X3 ~ normal(mu3, sig);
5679 21    // prior
5680 22    mu1 ~ uniform(0,100);
5681 23    delta1 ~ normal(0,100);
5682 24    delta2 ~ normal(0,100);
5683 25    sig ~ cauchy(0,5);
5684 26  }
5685
5686

```

5687 ここで注目すべきは、まず transformed parameters ブロックで mu2, mu3 という変数名が宣言され、そ
 5688 れが式によって表現されていることです。ここでパラメータ mu1, delta1 が mu2 に変形されていることがわ
 5689 かります。このようにして作られたパラメータであれば、model ブロックの中で使ってやることができます。現
 5690 に、尤度のところが $X2 \sim \text{normal}(\mu2, \text{sig})$ のようにスッキリした形になっていますね。ただし注意すべき
 5691 は、事前分布のところが変わっていない、ということです。事前分布はパラメータに対しておかれますので、こ
 5692 こで $\mu2 \sim \text{uniform}(0, 100)$; などと置くべきではありません。

5693 このように変換することで、コードが多少なりともスッキリしたと思います。面倒だな、と思うかもしれませんが
 5694 が、コードがスッキリすることは自分にとっても誰にとっても「読みやすい」ことにつながります。そして読みや
 5695 すいコードは、間違い (バグ) があつたときに修正しやすいのです。あまりごちゃごちゃしたコードを書くと、問
 5696 題が生じたときに対応が難しくなりますから、なるべくこうした作法を活用して、コードは単純明快なものにし
 5697 ておくと良いでしょう。

5698 22.2.2 パラメータの制約

5699 さてここでこのモデルをもう一段階、発展させてみましょう。今のコードは第一の群、 μ_1 を基準に、 μ_2, μ_3
 5700 を構成していました。別に第二、第三の群を基準にしても良かったですね。つまり基準点は恣意的になっ
 5701 ています。それで悪いということではないのですが、 μ_1 だけ特別扱いです。そうではない作法として、全体平均
 5702 μ を考えて、次のように定式化してみましょう。

$$\mu_1 = \mu + \delta_1 \quad (22.1)$$

$$\mu_2 = \mu + \delta_2 \quad (22.2)$$

$$\mu_3 = \mu + \delta_3 \quad (22.3)$$

5703 こうすることで、全体平均からの差分としてモデルを表現できます。ここで全体平均 μ はどういう意味がある
 5704 でしょうか。群の違いを独立変数、結果のスコアを従属変数とした線形モデルの一環として考えると、全体平
 5705 均は群を通じて変化のない水平線 (図 22.2) という関係を表していることになります。横一線で群ごとに違い
 5706 がないのですから、いわばこれは帰無仮説のモデル、すなわち帰無モデル (Null Model) とでもいうべき
 5707 ものです。

5708 あるいは、群ごとの平均値の違いはこの全体平均からの差分で表現されます。この群ごとの違いは、その
 5709 群に入ったからこそ生じた効果であるとも言えます。つまり、 δ は効果の大きさなのです。翻ってかんがえ
 5710 と、全体平均は「なんの処置もしなければ理論的にこの値になるはず」という操作なしの状態、天然の状態と
 5711 でもいえるでしょう。ランダム化比較実験 (Randomized Control Trial) は、被験者をランダムに割り
 5712 付けることによって、平均的に見れば個別の効果が相殺し合うことを用いています。平均化することで誤差
 5713 が相殺するので群ごとの比較ができますし、群の平均値が変化したということは、その群に属したことの効果
 5714 が現れたと考えるわけです。帰無モデルの状態は、この個別の効果を相殺させた状態であることを表してい
 5715 ます。

5716 さてこのように考えて、これを数式の形、モデルの形に書き起こしていくことを考えましょう。ところがこのま
 5717 まだと、未知のパラメータは $\mu, \delta_1, \delta_2, \delta_3$ と 4 つあることになります。今までは μ_1, μ_2, μ_3 , あるいは $\mu_1, \delta_1, \delta_2$
 5718 の 3 つだったのに、です。同じ数式を書き直したただけなのに、パラメータの数が変わるのはおかしいですね。
 5719 実はこれには 1 つ、制約が欠けているのです。今回は全体平均 μ からの相対比較になっていますから、
 5720 $\delta_1 + \delta_2 + \delta_3 = 0$ という関係が隠れていることになります*5。そこで、この制約をモデルの中に書き込まな

*5 もしこれが $= 0$ にならなければ、全体平均の計算に矛盾が生じますよね。全体平均 μ は、 $\mu = \frac{1}{3}(\mu_1 + \mu_2 + \mu_3)$ ですから、

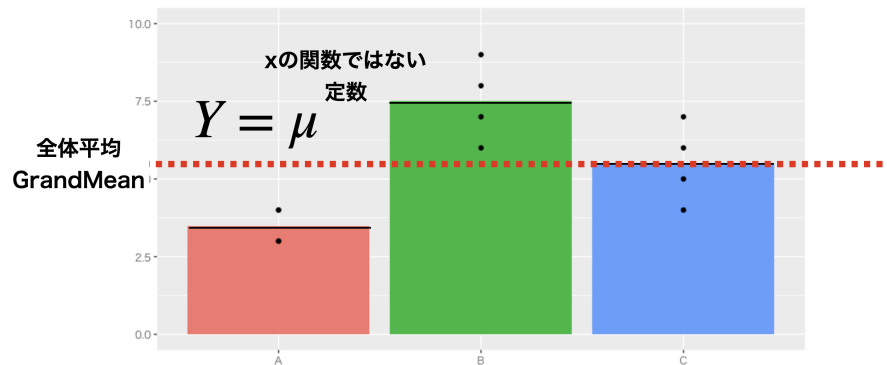


図 22.2 全体平均は水平線

5721 ければなりません。 $\delta_1 + \delta_2 + \delta_3 = 0$ という式から、変数を移項することでこれを表現します。すなわち

5722 $\delta_3 = 0 - (\delta_1 + \delta_2)$ のようにします。

5723 これを使って書いたのが code:22.4 です。

code : 22.4 制約を入れたコード

```

5724
5725 1 ...
5726 2 parameters{
5727 3   real gm;
5728 4   real delta1;
5729 5   real delta2;
5730 6   real<lower=0> sig;
5731 7 }
5732 8
5733 9 transformed parameters{
5734 10  real delta3;
5735 11  real mu1;
5736 12  real mu2;
5737 13  real mu3;
5738 14  delta3 = 0 - (delta1 + delta2);
5739 15  mu1 = gm + delta1;
5740 16  mu2 = gm + delta2;
5741 17  mu3 = gm + delta3;
5742 18 }
5743 19 ...
5744

```

5745 このようにすることで、先ほどと等価なモデルを別の形で書き表すことができるようになりました。

5746 ここまでをまとめておきます。三群の平均値差を求めるモデルを、3つの方法で書いてみました。

- 5747 1. 3つの群平均 μ_1, μ_2, μ_3 を個別に推定するモデル
- 5748 2. ある群とそこからの差分, $\mu_1, \delta_1, \delta_2$ を推定するモデル
- 5749 3. 全体平均とそこからの差分, $\mu, \delta_1, \delta_2, \delta_3$ を推定するモデル。ただし $\sum \delta = 0$ という制約をつける

$$\mu = \frac{1}{3}(\mu + \delta_1 + \mu + \delta_2 + \mu + \delta_3) = \frac{1}{3}(3\mu + \sum \delta_j) = \frac{1}{3}3\mu + \frac{1}{3}\sum \delta_j$$
となり、最初の項が $\frac{1}{3}3\mu = \mu$ ですから、第二の項は $\frac{1}{3}\sum \delta_j = 0$ でないと計算があいまいしません!

5750 この方法はいずれも同じ推定結果を返します。第一のモデルの生成量から差分を計算すれば、 δ_1, δ_2 を計算
5751 したことと同じです。

5752 同じことを異なる方法で実装するのは、生産的でないように思うかもしれません。しかし第3の形式にして
5753 おくと、モデルの一般化が容易になります。すなわち、多群の平均値差の比較をするときに、群の数が変わる
5754 たびにコードを書き換える手間がなくなるのです！

5755 22.3 モデルの洗練

5756 ここまで、二群、三群の平均値差を求めるモデルを書いてきました。書き方はわかっているので、比較する
5757 群が四群、五群と増えても書き足していくだけでなんとかなりそうです。

5758 しかし、たとえば四群に増えたときに $\delta_1, \delta_2, \delta_3$ 、五群に増えたときに $\dots, \delta_3, \delta_4$ と手作業で増やしていると、
5759 N 群までふえたときは $\dots, \delta_{N-2}, \delta_{N-1}$ と何行も書き足していかなければなりません。しかもその度にコンパイル
5760 しているようでは、とても面倒なことはすぐに想像がつかます。

5761 そこで、群の数が変わっても同じコードで対応できるように、一般化を目指してコードを改良しましょう。そ
5762 のためには、何群の比較なのかを示す群の数を、data ブロックで外部から取り込むようにすれば良いでしょ
5763 う。その変数を使って差分の数、効果の数を計算してやれば良いのです。たとえば比較する群の数を G とす
5764 ると、全体平均 μ と、 $G - 1$ 個の差分をパラメータとして推定すれば良いことになります。

5765 その考え方を使って書いたコードは、たとえば code:22.5 のようになります。

code : 22.5 群の数を一般化したコード

```
5766 1 data{
5767 2   int<lower=0> Lv;           // 水準数
5768 3   int<lower=0> N;           // 各群のサンプルサイズ
5769 4   array[Lv,N] real X;      // 変数の値
5770 5 }
5771 6
5772 7 parameters{
5773 8   real gm;                  // 全体平均
5774 9   array[Lv-1] real raw_delta; // 全体からの差。水準数マイナス1個
5775 10  real<lower=0> sig;        // 誤差の分散
5776 11 }
5777 12
5778 13 transformed parameters{
5779 14   array[Lv] real delta;     // 差の大きさを作り直す
5780 15   array[Lv] real mu;       // 再構成される群ごとの平均
5781 16
5782 17   for(i in 1:(Lv-1)){
5783 18     delta[i] = raw_delta[i]; // ほとんどコピー
5784 19   }
5785 20
5786 21   delta[Lv] = 0 - sum(raw_delta); // 総和が0になるように最後だけ書き換える
5787 22
5788 23   for(i in 1:Lv){
5789 24     mu[i] = gm + delta[i];   // 群ごとに再構成
5790 25   }
5791 26
5792 27 }
5793 28
```

```

5795 29 model{
5796 30   // Likelihood
5797 31   for(l in 1:Lv){
5798 32     for(i in 1:N){
5799 33       X[l,i] ~ normal(mu[l],sig);
5800 34
5801 35     }
5802 36   }
5803 37
5804 38   // Prior
5805 39   gm ~ uniform(0,100);
5806 40   raw_delta ~ normal(0,100);
5807 41   sig ~ cauchy(0,5);
5808 42 }
5809

```

5810 ■コード解説

5811 **data ブロック** 群の数 L_v , 各群のサンプルサイズ N , 各データ点 X を入力します。ここで X が**二元配列**に
 5812 なっていることに注目してください。 $X[1,1]$ は第一群の第 1 番目のデータ, $X[1,2]$ は第一群の
 5813 第 2 番目のデータ, $X[2,1]$ は第二群の第 1 番目のデータ, 同様に $X[i,j]$ は第 i 群の第 j 番目の
 5814 データを表すこととなります。この形式は二元配列 (2 次元の変数セット) になっているといえます*⁶。
 5815 宣言の時に, 上で取り込んだ変数 L_v や N を使うことができます。最大 L_v 個の群, 各群 N 人のデー
 5816 タがあるという前提です。

5817 **parameters ブロック** 全体平均 gm と, 差分, 誤差の SD をパラメータとしています。差分はここでは
 5818 `raw_delta` とし, 上で宣言した L_v をつかって L_v-1 個の要素を持つ配列にしています。 `raw_delta`
 5819 って変な名前, と思うかもしれませんが, この後これを変換して δ を作りますので, 作る前の生 (raw)
 5820 の δ という名前にしてみました*⁷。

5821 **transformed parameters ブロック** 後の話をわかりやすくするために, 各群の平均 μ を構成する
 5822 ことにします。 μ の数は水準数と同じ L_v です。そして第 i 群の平均 $\mu[i]$ を作るために,
 5823 $\mu[i]=gm+\delta[i]$ という書き方をしたいので, 水準数と同じ数だけ $\delta[i]$ を作りたいと思
 5824 います。もとの `raw_delta` が $1,2,\dots,L_v-1$ 個ありますので, δ も L_v-1 番目まではそれと同じも
 5825 のをコピーします (for 文で繰り返し代入して行ってます)。そして最後に L_v 番目の要素を, 0 から
 5826 これまでの要素をすべて足した (`sum(raw_delta)`) ものを引くことで作っています。ちなみに `sum` は
 5827 `stan` の持っている総和の関数で, 配列要素のすべてを足し合わせるというものです。こうして δ が
 5828 L_v 個できあがりましてので, 各水準の平均値を $\mu[i]=gm+\delta[i]$ という数式で一般的に書
 5829 くことができるようになりました。

5830 **model ブロック** 尤度のところに, 群それぞれについて巡回する for 文, その中で $1,2,3\dots N$ 人まで巡回す
 5831 る for 文を組み込んで二重にぐるぐる回しています。事前分布は設定の通りです。

5832 このコードを使って, 実際に推定してみましょう。データは表 22.1 のものを使います。

5833 推定してみましょう (R コードは 22.6, 結果は 7)。各群の平均値, 全体平均からの偏差などが, 事後分布

*⁶ 行列でもいいじゃないか, と思われるかもしれませんが。それでも結構です。Stan は行列の時, `real` ではなく `matrix` で宣言することができます。しかしこのように `real` 型にしておいて, 後ろのカッコで配列の次元を表現すると, 何次元でも拡張することができるので今回はそのようにしました。

*⁷ 変数名は任意ですので, 著者の命名法がわかりにくいなあと思ったら, 自分で好きな変数名にもらって構いませんよ。その場合は, 以後すべての変数を読み替えてくださいね。

表 22.1 独立した三群から得られたスコア

groupA	6	6	5	5
groupB	7	4	7	6
groupC	4	3	4	6

5834 として出力されているのがわかると思います。これをみて、たとえば差の事後分布の 50% 確信区間でみると
 5835 δ_2, δ_3 はその区間に 0 を挟んでないので、一定の効果はあるだろうなどと判断できるわけです。もちろん生成
 5836 量をつかって、ある程度の差がある確率がどのくらいとか、優越率などデータレベルの仮説を立てることも可
 5837 能です。

code : 22.6 三群の平均値の比較

```
5838 1 Example <- matrix(c(6, 6, 5, 5, 7, 4, 7, 6, 4, 3, 4, 6),
5839 2                       ncol = 4, byrow = T)
5840 3 dataSet <- list(Lv = 3, N = 4, X = Example)
5841
5842
```

MCMC の結果 7

```
# A tibble: 10 × 7
  name          EAP      MED      MAP      SD      L95      U95
  <chr>      <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 delta[1]    0.257    0.256    0.227    0.569   -0.880    1.394
2 delta[2]    0.749    0.754    0.804    0.566   -0.387    1.878
3 delta[3]   -1.005   -1.009   -1.114    0.576   -2.158    0.162
4 gm          5.251    5.251    5.242    0.400    4.452    6.039
5 mu[1]       5.507    5.507    5.562    0.699    4.114    6.902
6 mu[2]       5.999    6.001    5.986    0.697    4.603    7.386
7 mu[3]       4.246    4.238    4.182    0.695    2.860    5.647
8 raw_delta[1] 0.257    0.256    0.227    0.569   -0.880    1.394
9 raw_delta[2] 0.749    0.754    0.804    0.566   -0.387    1.878
10 sig        1.339    1.261    1.134    0.391    0.823    2.301
```

5843

22.3.1 データサイズの一般化

5844

5845 さてこのモデルは、比較する群の数が増えても外部からデータとして群の数を与えることができますので、
 5846 毎回コンパイルする必要がありません。やったね！これで完璧、と言いたいところですが、1 つ気になるのは
 5847 サンプルサイズです。各群のサンプルサイズ N をデータとして与えるようになっていますが、たとえば A 群は
 5848 10 人、B 群は 12 人、C 群は 14 人と言ったように、群ごとにサンプルサイズが異なるとこのコードでは対応
 5849 できません。実際に実験をやる場合、各群のサイズを整えて人を集めたいところですが、調査研究などでは群
 5850 ごとに同じ人数を与えるような調整ができないこともあり、そういう意味ではこのコードではうまく対応できそ
 5851 うにありません。

5852 そこで、群の人数が異なる場合でも対応できるように、さらにこのコードを拡張していきたいと思います。そ
 5853 のためにまず、データを**整然データ (tidy data)** の形に整形しましょう。

5854 先ほどのサンプルデータ (表 22.1) は、 3×4 の長方形をしていました。これはデータのサイズが整ってい
 5855 るからできることで、サイズが変わってしまうと表の中に欠損ができてしまうことになります。また、たとえば B

5856 群の 3 番目の人の値を見る時、われわれは 2 行 3 列目のデータにさっと目が行きますが、機械的には群の
 5857 情報が行名に、人の群内整理番号が列名にあるので、1 つのデータを見る時に行・列それぞれを参照するこ
 5858 とになります (図 22.3 左)。ここで「どの群か」「群の中の何番目の人か」という情報はいずれも変数で、人
 5859 よって変わるものですから、データの持つ情報の 1 つなのです。データの持つ情報、変数なのに参照先が
 5860 変わっていることになります。

5861 このデータの情報をまったく損なうことなく、縦長に並べ直したのが図 22.3 の右側です。これは 1 つの行
 5862 がそれぞれの観測に対応しています。このような形式にすると、何行目かということ指定するだけで、ど
 5863 の群の何番目の人なのかという情報が手に入ります。またデータの中に欠損があっても、整然データの形にす
 る時にこれは観測に含まれないので、データの一部が欠けることがないという利点があります。

群内整理番号				
群	1	2	3	4
A	6	6	5	5
B	7	4	7	6
C	4	3	NA	6

群	番号	値
A	1	6
A	2	6
A	3	5
A	4	5
B	1	7
B	2	4
B	3	7
B	4	6
C	1	4
C	2	3
C	4	6

図 22.3 一般的なデータと整然データ

5864
 5865 このようなデータ形式にして、これを与える形でデータ分析をするとデータサイズに依存しないコードが書
 5866 けます。具体的には code:22.7 のようにします。

code : 22.7 整然データに対応したコード

```

5867 1 data{
5868 2   int<lower=0> Lv;           // 水準数
5869 3   int<lower=0> L;           // データ数
5870 4   array[L] int<lower=0,upper=Lv> idx; // データの ID
5871 5   array[L] real X;         // 変数の値
5872 6 }
5873 7 ...(中略)...
5874 8 model{
5875 9   // Likelihood
5876 10  for(l in 1:L){
5877 11    X[l] ~ normal(mu[idx[l]],sig);
5878 12  }
5879 13
5880 14  // Prior
5881 15  gm ~ uniform(0,100);
5882 16  raw_effect ~ uniform(-100,100);
5883 17  sig ~ uniform(0,100000);
5884 18 }
5885
5886

```


parameters ブロックと transformed parameters ブロックは code:22.5 と同じなので省略しました。まず data ブロックをみてください。群の数を Lv で入力するのはそのままですが、次にデータの総数 L を撮るようにしています。図 22.3 の例で言えば、C 群 3 列目のデータがありませんから、データ長は全部で $L = 11$ になります。次に L 行のデータがどの群に属するのかを指示するインデックス変数 idx を L 個用意しています。L 行目のデータが第何群に属しているのかの数字が入ります。たとえば 1 行目は A 群なので 1, 5 行目 ($L=5$) は B 群なので 2, とした数字が入るようになっていきます。さいごにデータの値そのものである X も、データ長 L と同じだけの 1 次元配列を用意してやります。

技術的に注目すべきは model ブロックの尤度のところで、for 文がデータ長 L の反復をしているだけです。つまり 1,2,3,...,L 行目のデータを順に参照していくのです。そして l 行目のデータがどの群に属するかは、変数 idx[l] が指し示してくれますから、mu[idx[l]] としてやることで指定できていることになります。idx[l] は、l 行目のデータが属している群の数字が代入されていますから、たとえば 1 行目 ($l=1$) の場合は mu[idx[1]] = mu[1] としていることと同じ、5 行目の場合は mu[idx[5]] = mu[2] としていることと同じ、ということになります。ここでは入れ子になった参照が行われています。ちょっとテクニカルですが、この技術を使うと表現力も広がりますので、仕組みをしっかりと理解しておいてください。

22.4 パラメータリカバリ

さあ、最後にパラメータリカバリをして、このコードが正しく推定できるか、あるいはどの程度のサンプルサイズがあればどの程度の精度で推定できるのかを検証しておきましょう。

私たちはデータがどのように造られるか、というそのメカニズムの方からアプローチしているわけです。これはリバースエンジニアリングと呼ばれる考え方ですね。そして今から仮想データを作ろうという時も、そのデータ生成モデルをそのまま利用すればいいのです。仮想データの生成は、データ生成メカニズムをリバースエンジニアリングで明らかにしていくアプローチのちょうど正反対、リバース・リバース・エンジニアリングです。

具体的には、たとえば次のようなコード (code:22.8) でデータを作ることができるでしょう。

code : 22.8 リバースエンジニアリング

```

5909 1 N <- 100
5910 2 Lv <- 5
5911 3 gm <- 50
5912 4 sig <- 3
5913 5
5914 6 raw_effect <- runif(Lv - 1, -10, 10)
5915 7 effect <- c(raw_effect, 0 - sum(raw_effect))
5916 8 mu <- gm + effect
5917 9
5918 10 X <- rnorm(N * Lv, mu, sig)
5919 11 dat <- data.frame(
5920 12   Idx = rep(1:Lv, N * Lv),
5921 13   value = X
5922 14 )
5923 15
5924 16 modelC <- cmdstanr::cmdstan_model("BetweenAnova2.stan")
5925 17 dataSet <- list(Lv = Lv, L = NROW(dat), idx = dat$Idx, X = dat$value)
5926
5927

```

■コード解説

- 5929 1 行目 各群のサンプルサイズ N を設定します。各群共通のサイズになります。
- 5930 2 行目 水準数です。ここでは 5 群の平均値の比較をすることにしました。
- 5931 3 行目 全体平均の設定です。事前分布に $[0,100]$ の一様分布を置いてますから、真ん中ぐらいにしてあげ
5932 ました。
- 5933 4 行目 誤差の散らばりを設定します。
- 5934 6 行目 水準数 -1 の効果を設定します。手入力でもいいのですが面倒なので、 -10 から 10 までの範囲で
5935 一様乱数によって生成しました*8。
- 5936 7 行目 水準数 L_v の効果にするため、先ほど作った水準数 -1 の `raw_effect` に `0-sum(raw_effect)`
5937 の計算結果をつけくわえています。関数 `c` は結合させる `combine` という意味です。
- 5938 8 行目 各群の平均値です。全体平均に効果を足しています。全体平均は 1 つの数字、効果は水準数の
5939 要素を持つベクトルですから、計算ができないように思えますが、このような場合 R は自動的にサ
5940 イズ合わせのため値の再利用を行います。つまり、 $50 + (\delta_1, \delta_2, \delta_3, \delta_4, \delta_5)$ を $(50, 50, 50, 50, 50) +$
5941 $(\delta_1, \delta_2, \delta_3, \delta_4, \delta_5)$ と解釈してくれるのです。
- 5942 10 行目 水準数 \times 各群のサイズの乱数を発生させています。乱数は正規分布によるもので、平
5943 均や SD は上で指定した通りです。ここでも値の再利用が行われていて、ベクトル `mu` は
5944 $\mu_1, \mu_2, \mu_3, \mu_4, \mu_5, \mu_1, \mu_2, \dots$ とくり返して当てはめられていきます。
- 5945 11-14 行目 数値を `data.frame` 型にくみ上げていってます。`rep` 関数は、指定した回数だけ繰り返すも
5946 ので、ここでは `1:L_v` すなわち `1,2,3,4,5` という数字をデータサイズ分繰り返していることになります。
- 5947 16 行目 ここでは `cmdstanr` パッケージを使ってコンパイルしています。
- 5948 17 行目 推定に使うデータセットです。データ長は `data.frame` の長さを返す関数 `NROW` を使って与えて
5949 います。インデックス変数は `dat` オブジェクトの `Idx` 変数、データは同じオブジェクトの `value` 変数
5950 を指定しています。

5951 いかがでしょうか。作ったデータ、オブジェクト、ベクトルなどの意味がわからない場合は、その都度 R でオ
5952 ブジェクト名を入力し、何が格納されているかを確認しながら進めると良いでしょう。ここで作った仮想データ
5953 を、先ほどの Stan オブジェクトに与えて乱数を生成し、パラメータリカバリがどの程度の精度でできるかを考
5954 えてみてください。サンプルサイズや誤差の大きさなどを変えながら確認してみると良いでしょう。

5955 22.5 課題

5956 今回は、基本課題と発展課題の 2 つを用意します。基本課題は必須ですが、応用課題は提出者に加点さ
5957 れるだけで必須ではありません (提出しなかったからと言って減点されることもありません)。

5958 ■基本; 整然データに対応した多群比較のコードを書く 最後に紹介した、整然データに対応した多群
5959 比較のコードを書き、パラメータリカバリのコードとデータを使って正しく動くかどうか確認してください。作っ
5960 た Stan ファイルや R コードを提出してください。

5961 ■発展; 自分のデータを使って分析する 基礎実験 1 や基礎実験 2 など、他の授業でとったいくつかの
5962 群の平均値を比較するようなデータを用意し、今回のモデルを適用して平均値の比較を行なってください。と
5963 くにもそのようなデータセットを持っていない場合は、こちらから提供するサンプルデータを使ってください。

*8 一様乱数の関数は、R では `runif` です。

第 23 章

モデリングの目から見た検定 4 ; 対応のある群の比較

23.1 対応のある群

23.1.1 多次元正規分布

ここまでは群間計画 (Between Design) のモデリングをしてきましたが、今回は群内計画 (Within Design) のモデリング的アプローチになります。

Within モデルは別名として、反復測定 (Repeated Measure) と呼ばれることがあり、とくに二群の場合は対応のある t 検定 (Paired t-test) と言ったりします。同じ人から 2 回データを取って、その変化をみる事前・事後デザイン (Pre-post Design) がその典型例です。2 回以上とる場合もあるので、その場合は反復測定とよばれるわけです。

この場合のデータ生成モデルは、何が違うのでしょうか。群間計画の t 検定は別名独立した二群の t 検定と呼ばれるように、2 つの群のデータがバラバラに出てきているのですが、今回は対応のある二群ですから、独立してないということになります。独立していない、つまり関係がある。データ分析上は、2 つのデータに相関がある場合を考える必要があります。同じ人の事前・事後のデータであれば、当然変化のもとになる「その人」という共通要因があるわけで、同じ人からデータを得ていますから当然相関していることを前提に考えなければなりません。

統計モデル上、独立した二群であれば別々の正規分布からデータが出てきていると考えますが、今回は 1 つの分布から 2 つのデータが出てきていることになります。ここで使うのは二次元正規分布 (Two-Dimensional Normal Distribution)、さらに群の数が多い場合に一般化した多次元正規分布 (Multidimensional Normal Distribution) あるいは多変量正規分布 (Multivariate Normal Distribution) と呼ばれる分布を使います。この正規分布の特徴は、それぞれの変数 (次元) に注目すれば正規分布なのですが、各次元に相関関係があることを組み込まなければなりません。

これまでの 1 次元の正規分布は、平均 (位置) μ と標準偏差 (幅) σ で特徴付けられました。多次元の正規分布も、各次元について平均と標準偏差があります。平均も複数の次元のセットになっているので、平均ベクトル $\boldsymbol{\mu}$ として表現することになります。標準偏差もベクトルで考える必要がありますし、さらにすべての組み合わせにおける相関関係がありますので、分散共分散行列 $\boldsymbol{\Sigma}$ を考えることになります。丁寧に数式で表現すると、1 次元正規分布に従う確率変数は次のように書きます。

$$X \sim Normal(\boldsymbol{\mu}, \boldsymbol{\sigma})$$

5992 これに対して多次元正規分布に従う確率変数ベクトル \mathbf{X} は次のように書きます。

$$\mathbf{X} \sim \text{MultiNormal}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

5993 この分散共分散行列の中身を見ますと、次のようになっています。

$$\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1m} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \cdots & \sigma_m^2 \end{pmatrix}$$

5994 この行列の**対角 (diagonal)** 要素には分散 σ_j^2 があり、非対角要素には共分散 σ_{jk} が入っている**正方対**
 5995 **称行列**になっていますね*¹。またこうしてみると、変数 j の分散 σ_j^2 は j と j の共分散、すなわち σ_{jj} であるこ
 5996 ともわかりやすいですね。ところでこれを考えることでどこに相関の要素が入っているのでしょうか？改めて、分
 5997 散や共分散、相関係数の定義式から考えてみたいと思います。

$$\text{分散} : \sigma_j^2 = \frac{1}{n} \sum (x_{ij} - \bar{x}_j)^2 \quad (23.1)$$

$$\text{共分散} : \sigma_{jk} = \frac{1}{n} \sum (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \quad (23.2)$$

$$\text{相関} : \rho_{jk} = \frac{\sigma_{jk}}{\sigma_j \sigma_k} \quad (23.3)$$

5998 この関係から逆に、共分散 σ_{jk} は次のように表すことができます。

$$\sigma_{jk} = \sigma_j \sigma_k \rho_{kj}$$

5999 ということで、共分散の式は分布の幅を表す標準偏差と、相関係数からできあがっていることがわかりまし
 6000 たね。

6001 23.1.2 対応のあるデータの生成モデル

6002 ではこれを使ってモデルの設計図を書いてみましょう。設計図は図 23.1 のようになるでしょう。

6003 あとはこれを使ってコードを書くだけです。この時のポイントは、データの渡し方にあります。

6004 ベクトル型とマトリックス型

6005 2次元正規分布に従うデータは、必ず2つセットになっています。1回の観測で2つの数字を持つもの
 6006 です。これを Stan で表現するときには、変数を `vector` 型で表現しておく必要があります。今回の場合は
 6007 `vector[2] mu;` などとします。これで `mu` という変数が2つセットで準備されることになります。これまで同
 6008 じ変数名で複数のデータを持つ場合、たとえば `real mu[2]` のようにしていましたが、なぜ今回もこの形式、
 6009 つまり**配列**にしないのか、と思われるかもしれませんが。配列でも2つの数字を扱ったりすることは可能なので
 6010 すが、ベクトル型にしておく利点が別にあります。ベクトルにはベクトルを使った演算がありますが、Stan には
 6011 ベクトルや行列専用の関数がありますので、変数をベクトルで宣言しておいてやるとこれらの関数を使って高
 6012 速化できるのです。

6013 また Stan には、`matrix` 型という行列の型もあります。`matrix[N,M] X;` のように宣言してやると、サイ
 6014 ズ $N \times M$ の行列変数 X が1つ用意されます。`matrix[N,M] X[L];` のように宣言すると、サイズ $N \times M$

*¹ これら行列の要素や名称については、第 6 講でやったことを思い出してください。

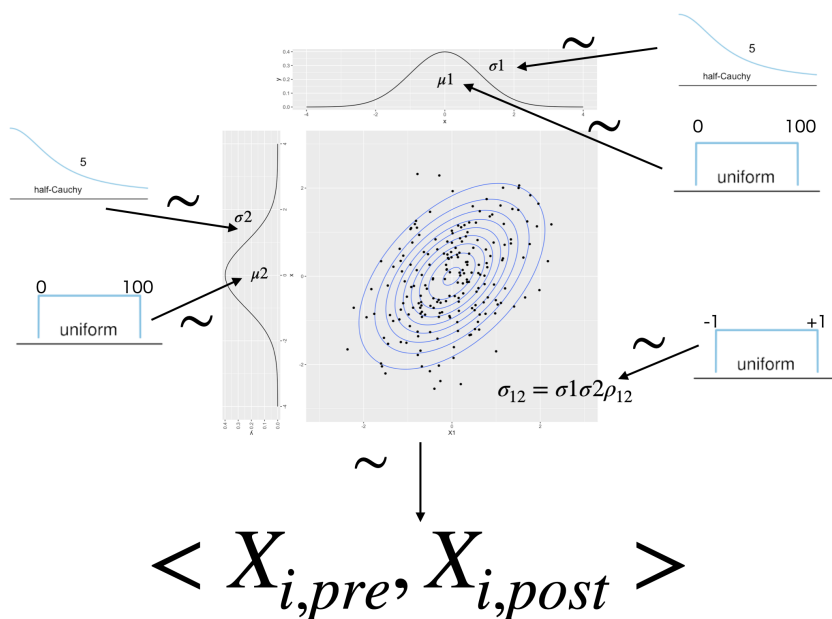


図 23.1 対応ある二群のデータ生成モデル設計図

6015 の行列変数 X を L 個用意することになります。他にも特殊な型として、`cov_matrix` 型とか `corr_matrix`
 6016 型などがあります。これらは正方行列ですから、サイズの指定は 1 つでよく、たとえば `cov_matrix[2]` とす
 6017 れば 2×2 の正方行列を用意したことになります。Stan における変数や配列の型について、松浦 (2016) を
 6018 参考にまとめたものを表 23.1 に、イメージ図を図 23.2, 23.3, 23.4 に示します。

表 23.1 Stan の配列と型

型の例	宣言例	解説
整数	<code>int X</code>	整数ひとつ
実数	<code>real X</code>	実数ひとつ
範囲つき実数	<code>real<lower=0> X</code>	0 より大きい数字の入る X
整数の配列	<code>int X[N]</code>	N 個の整数
実数の配列	<code>real X[N,M]</code>	$N \times M$ 個の実数からなる 2 次元配列
実数の配列	<code>real X[N,M,L]</code>	$N \times M \times L$ 個の実数からなる 3 次元配列
ベクトル	<code>vector[K] X</code>	長さ K のベクトルひとつ
ベクトル	<code>vector[K] X[N]</code>	長さ K のベクトルが N 個
ベクトル	<code>vector[K] X[N,M]</code>	長さ K のベクトルが $N \times M$ 個
行列	<code>matrix[J,K] X</code>	サイズ $J \times K$ の行列がひとつ
行列	<code>matrix[J,K] X[N]</code>	サイズ $J \times K$ の行列が N 個
特別な行列	<code>cov_matrix[J] X</code>	サイズ $J \times J$ の正方対称行列がひとつ
特別な行列	<code>corr_matrix[J] X</code>	対角が 1 のサイズ $J \times J$ の正方対称行列がひとつ

6019 これを踏まえて、今回の対応のある二群のコード例を見てみましょう (Code::23.1)。

code : 23.1 対応ある二群のデータ生成モデル

6020




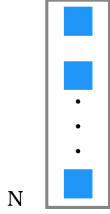
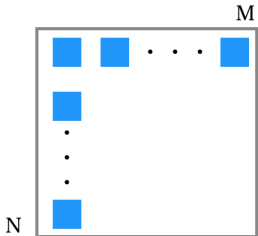
型	宣言例	サイズ
整数	<code>int X</code>	
実数	<code>real X</code>	
範囲つき実数	<code>real<lower=0> X</code>	
整数の配列	<code>int X[N]</code>	
実数の配列	<code>real X[N,M]</code>	

図 23.2 変数の型とサイズ (1)

```

6021 1 data{
6022 2   int N;
6023 3   array[N] vector[2] X;
6024 4 }
6025 5
6026 6 parameters{
6027 7   vector[2] mu;
6028 8   real<lower=0> sd1;
6029 9   real<lower=0> sd2;
6030 10  real<lower=-1,upper=1> rho;
6031 11 }
6032 12
6033 13 transformed parameters{
6034 14   cov_matrix[2] SIG;
6035 15   SIG[1,1] = sd1 * sd1;
6036 16   SIG[1,2] = sd1 * sd2 * rho;
6037 17   SIG[2,1] = sd2 * sd1 * rho;
6038 18   SIG[2,2] = sd2 * sd2;
6039 19 }
6040 20
6041 21 model{
6042 22   X ~ multi_normal(mu, SIG);
6043 23   //prior
6044 24   mu[1] ~ uniform(0,100);
6045 25   mu[2] ~ uniform(0,100);

```


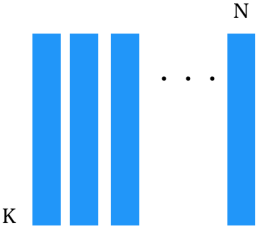
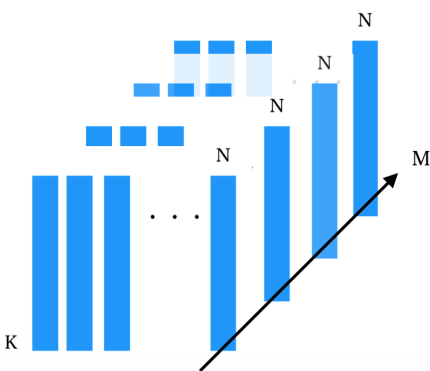
型	宣言例	サイズ
ベクトル	<code>vector[K] X</code>	
ベクトル	<code>vector[K] X[N]</code>	
ベクトル	<code>vector[K] X[N,M]</code>	

図 23.3 変数の型とサイズ (2)

```

6046 26 rho ~ uniform(-1,1);
6047 27 sd1 ~ cauchy(0,5);
6048 28 sd2 ~ cauchy(0,5);
6049 29 }
6050

```

6051 ■コード解説

6052 **data ブロック** サンプルサイズ N のデータですが、ペアになっている数字なので `vector[2]` で宣言した
6053 変数 X を N 個用意しています。

6054 **parameters ブロック** パラメータとして、平均をベクトル、標準偏差や相関係数はスカラで宣言しています。
6055 標準偏差や相関係数の範囲指定にも注意してください。

6056 **transformed parameters ブロック** 各標準偏差、相関係数を使って分散共分散行列 SIG を構成していま
6057 す。わかりやすくするために、各要素について書き下して書いています。

6058 **model ブロック** 尤度はデータベクトルが多次元正規分布関数 `multi_normal` から出てきていることを
6059 示しています。多次元正規分布関数の引数は、ベクトルと行列です。

6060 この書き方ですと、三群に拡張した時に大変になるのは目に見えていますね。こんな時のために、ベクトル
6061 と行列用関数をつかって書いておく方法があります。Code::23.2 には、多次元にまで拡張したコードを示し


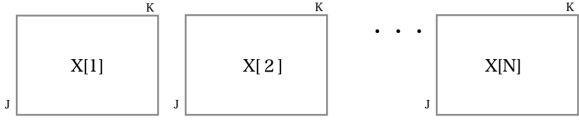
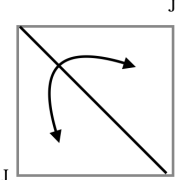
型	宣言例	サイズ
行列	<code>matrix[J,K] X</code>	
行列	<code>matrix[J,K] X[N]</code>	
特別な行列 特別な行列	<code>cov_matrix[J] X</code> <code>corr_matrix[J] X</code>	 <p>正方対称 $J \times J$ $X_{ij} = X_{ji}$</p>

図 23.4 変数の型とサイズ (3)

6062 ています。

code : 23.2 コードの一般化

```

6063 1 data{
6064 2   int N;
6065 3   int K;
6066 4   array[N] vector[K] X;
6067 5 }
6068 6
6069 7 parameters{
6070 8   vector[K] mu;
6071 9   vector<lower=0>[K] sds;
6072 10  corr_matrix[K] rho;
6073 11 }
6074 12
6075 13 transformed parameters{
6076 14  cov_matrix[K] SIG;
6077 15  SIG = quad_form_diag(rho,sds);
6078 16 }
6079 17
6080 18 model{
6081 19  X ~ multi_normal(mu,SIG);
6082 20  //prior
6083 21  mu[1] ~ uniform(0,100);
6084 22  mu[2] ~ uniform(0,100);
6085 23  rho ~ lkj_corr(1);
6086 24  sds ~ cauchy(0,5);
6087 25 }
6088
6089

```

6090 ■コード解説

6091 data ブロック データのサイズ K も外侮から指定できる変数とし、それを使ってサイズ K のベクトル X を
6092 N 個用意しています。

6093 parameters ブロック パラメータとして、平均と標準偏差をベクトルで、相関係数も行列として宣言しまし
6094 た。標準偏差ベクトルの範囲指定の方法に注意してください。

6095 transformed parameters ブロック Stan の関数 quad_form_diag を使って分散共分散行列を作成しま
6096 した。この関数は、正方行列 X とベクトル v から、 $\text{diag}(v)X\text{diag}(v)$ という計算をするものです。
6097 ここで diag は k 個の要素を対角項に持つ $k \times k$ の対角行列を作ることを意味します。今回の例で
6098 具体的に説明しましょう。ここでの正方行列は相関行列で、ベクトルは標準偏差を要素に持つ sds で
6099 すから、

$$\begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix} \begin{pmatrix} 1 & \rho_{12} \\ \rho_{12} & 1 \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix} = \begin{pmatrix} \sigma_1 & \sigma_1\rho_{12} \\ \sigma_2\rho_{12} & \sigma_2 \end{pmatrix} \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix} = \begin{pmatrix} \sigma_1^2 & \sigma_1\sigma_2\rho_{12} \\ \sigma_1\sigma_2\rho_{12} & \sigma_2^2 \end{pmatrix}$$

6100 となり、分散共分散行列が構成されていることがわかりますね。

6101 model ブロック 注目すべきポイントは相関行列の事前分布です。相関行列の事前分布には、Stan では
6102 lkj_corr という関数を使うのが一般的で、この引数が 1 になっているのは無情報分布であることを
6103 意味しています。

6104 このコードで推定される平均値、mu ベクトルの 2 つの要素が、事前・事後の平均値に該当します。ですから
6105 この差分を計算すれば「平均的にどれぐらい変化があったか」という平均因果効果を見積もったことになり
6106 ます。この差分については、生成量のブロックで算出しても構いません。もちろんその差分から標準化した効
6107 果の大きさを算出することもできますし、パラメータについての仮説、事後予測などデータについての仮説な
6108 ど、さまざまに検証できることはこれまでの通りです。

6109 23.2 ID をもったデータ構造

6110 さてでは今回の乱数生成機の精度を検証するために、パラメータリカバリをしてみましょう。

6111 仮想データの作り方は、データ生成モデルの裏返しです。2 群の平均値差を検証するためのデータでした
6112 ら、次のようにして作ることができます。

code : 23.3 対応ある二群の仮想データ

```
6113 1 mu <- c(50,50)
6114 2 sd1 <- 10
6115 3 sd2 <- 5
6116 4 rho <- 0.7
6117 5 SIG <- matrix(ncol=2,nrow=2)
6118 6 SIG[1,1] <- sd1 * sd1
6119 7 SIG[2,2] <- sd2 * sd2
6120 8 SIG[1,2] <- sd1 * sd2 * rho
6121 9 SIG[2,1] <- sd1 * sd2 * rho
6122 10 N <- 100
6123 11 library(MASS)
6124 12 X <- mvrnorm(N,mu,SIG)
6125 13 # Stan に 与 える データ セット
6126 14 dataSet <- list(N=N, X=X)
6127
6128
```

6129 ここでは二群の平均をいずれも 50, 標準偏差を 10 と 5 にし, 相関係数を 0.7 としたサンプルを 100 個作る
 6130 ことにしました。発生させる乱数は MASS パッケージの `mvrnorm` 関数を使います。この関数は, 発生させるサ
 6131 ンプルサイズと, 平均ベクトル, 分散共分散行列を引数に取ります。これを使って推定させた結果の例が出力
 6132 8 になります。

MCMC の結果 8

```
# A tibble: 9 × 7
  name          EAP      MED      MAP      SD      L95      U95
<chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu[1]      47.590  47.594  47.608  0.913  45.808  49.356
2 mu[2]      49.506  49.508  49.570  0.500  48.523  50.493
3 rho         0.693   0.697   0.703   0.052   0.582   0.784
4 sd1         9.087   9.059   9.071   0.640   7.919  10.434
5 sd2         4.980   4.960   4.930   0.352   4.345   5.724
6 SIG[1,1]    82.975  82.063  81.989  11.761  62.715  108.872
7 SIG[1,2]    31.597  31.189  30.790   5.515  22.086  43.574
8 SIG[2,1]    31.597  31.189  30.790   5.515  22.086  43.574
9 SIG[2,2]    24.923  24.599  24.207   3.550  18.883  32.761
```

6133

6134 中央値による点推定値 (MED 推定値) を見ると, ほぼ真の値の通りになっていますね。
 6135 ところで, 今回与えたデータ X は行列の形をしています。一行が 2 つの数字を持つペアになっているので
 6136 すが, これは**整然データ**の形にはなっていませんね。整然データであれば, 次のようなデータセットになってい
 6137 るはずです (R の出力 23.1)。

R の出力 23.1: 対応ある二群の仮想データ・整然版

```
# A tibble: 200 × 3
  ID name value
<int> <chr> <dbl>
1     1 pre  45.1
2     1 post 45.8
3     2 pre  45.1
4     2 post 51.2
5     3 pre  52.4
6     3 post 53.7
7     4 pre  46.8
8     4 post 48.6
9     5 pre  45.5
...

```

6138

6139 このようなデータ形式に対応できるように, Stan のコードを書き直してみましよう。ポイントは, 個体識別イ
 6140 ンデックスと, 事前・事後を表す変数のインデックスの両方を準備する, というところにあります。

6141 まず R から与えるデータの方から考えましよう。整然データの形に並び替えるのに加えて, 事前・事後と
 6142 いった水準を表す変数も数字で与えなければなりません。先の例 Code::23.4 では `pre, post` となっていた
 6143 ところを, `pre=1, post=2` と置き換えた別の変数を用意しました。

code : 23.4 対応ある二群の仮想データ・整然版

6144

```

6145 1 ...
6146 2 X <- mvrnorm(N, mu, SIG)
6147 3 tidy_data <- X %>%
6148 4   as.data.frame() %>% as_tibble() %>%
6149 5   rename(pre = V1, post = V2) %>%
6150 6   rowid_to_column("ID") %>%
6151 7   pivot_longer(-ID) %>%
6152 8   mutate(cond = if_else(name == "pre", 1, 2))
6153

```

6154 ■コード解説

6155 2 行目 真値の設定をして、行列 X を作っています。

6156 3 行目 X を加工して、tidy_data という変数に作り替えます。加工のプロセスは下記の通りです。

6157 4 行目 data.frame 型を経て tibble 型に変形します*2。

6158 5 行目 この段階で、変数名が V1, V2 となっていますが、わかりやすくするために pre, post に書き換えま
6159 した。

6160 6 行目 行番号を ID という変数名を持つものに変えています。

6161 7 行目 ID をキーに、横長だったデータを縦長に (tidy に) 変換する関数です。

6162 8 行目 縦長になった変数 (Code::23.4 はこの段階の出力です) から、name 変数の値が pre であるもの
6163 を 1 に、そうでないものを 2 にした新しい変数 cond を作りました。

6164 このコードでどういうデータセットができたのか、R のコンソールで tidy_data として確認すると良いでしょ
6165 う。1 つ 1 つのステップを確認したい場合は、途中で止めて一歩ずつ進むのも手です。

6166 さて Stan にデータセットを与えるときは、このデータセットから 1. データ全体の長さ、2. 被験者の数、3.
6167 被験者 ID、4. そのデータがどの水準なのかを示す識別子、5. 実際の値を取り出して渡すことになります。
6168 続いて Stan コードの例を見てみましょう。

code : 23.5 整然データ対応番

```

6169 1 data{
6170 2   int L;
6171 3   int N;
6172 4   array[L] int<lower=0,upper=N> IDindex;
6173 5   array[L] int<lower=1,upper=2> Condition;
6174 6   array[L] real val;
6175 7 }
6176 8
6177 9 transformed data{
6178 10  array[N] vector[2] pairX;
6179 11  for(1 in 1:L){
6180 12    pairX[IDindex[1],Condition[1]] = val[1];
6181 13  }
6182 14 }
6183 15 ... (中略) ...
6184 16 model{
6185 17  pairX ~ multi_normal(mu, SIG);
6186

```

*2 tibble 型にする必要は別にありません。著者の趣味です。tibble 型は data.frame 型の拡張版で表示した時に変数の型などがわかりやすいという利点があります。matrix 型からいきなり tibble 型に変更すると警告が出るので、いったん data.frame 型を経由しました。

```

6187 18 //prior
6188 19 mu[1] ~ uniform(0,100);
6189 20 mu[2] ~ uniform(0,100);
6190 21 rho ~ uniform(-1,1);
6191 22 sd1 ~ cauchy(0,5);
6192 23 sd2 ~ cauchy(0,5);
6193 24 }
6194

```

6195 ■コード解説

6196 data ブロック 1. データ全体の長さ L, 2. 被験者の数 N, 3. 被験者 IDIndex, 4. そのデータがどの水
6197 準なのかを示す識別子 Condition, 5. 実際の値 val を与えています。最後の 3 つはデータ長と同じ
6198 サイズです。

6199 transformed data ブロック データを変形する transformed data ブロックの登場です。ここではペア
6200 としての情報がある vector 型にした変数を作っています。その変数 pairX は、N 人分の情報が
6201 あり、それぞれ要素番号 1,2 がベクトルの要素に対応しています。ここでも入れ子になった識別子を使
6202 っていて、pairX[IDIndex[l],Condition[l]] は l 行目の個人 IDIndex[l], l 行目の条件
6203 Condition[l] から、たとえば IDIndex[l] = 15, Condition[l]=2 であれば 15 番目の人の事
6204 後の値 pairX[15,2] を値として代入する、ということをしています。

6205 model ブロック 作った pairedX というベクトルに対して多次元正規分布からのデータを与えています。

6206 少しややこしくなったように思えますが、データの構造が複雑化すると整然データの方がスッキリすることの
6207 方が一般的です。コードをよく読んで、どのような動きをするのかイメージをしっかりと掴むようにしましょう。

6208 23.3 個人差と変化量のモデルへ

6209 さて、コードは随分と複雑になってきましたが、やっていることは「事前・事後で変化はあるか」ということだ
6210 けです。それなのに多次元正規分布が出てくるなんて、面倒ですねえ！もう少し簡単な表現はできないもの
6211 でしょうか。たとえば事前 + 効果 = 事後のように。

6212 もちろん可能です。むしろその方が自然なモデリングと言えるかもしれないですね。ただその時、事前・事後
6213 が同一人物であることを忘れてはいけません。小杉の事前の状態 + 効果 = 山田の事後の状態、では意味が
6214 わからないからです。記号を使って書くと、 $X_i^{pre} + effect = X_i^{post}$ のように、添字 i で同じ個人であるこ
6215 とを明記しておく必要があります。

6216 これを踏まえて、少し話を拡張した 3 水準モデルを考えてみましょう。たとえば次のようなカバーストーリー
6217 はいかがでしょうか。

6218 臨床心理学者がある介入法をつかってメンタルケアに取り組んでいます。4 人のクライアントにこの手
6219 法を適用し、抑うつ度が改善されていくかどうかチェックし、この介入法に効果があるのかどうかを検
6220 証したいと思っています。

6221 チェックの結果が表 23.2 のようになっていて、ここの数字が抑うつ度スコア (低ければ低いほど健康) だと
6222 思ってください。記号で書く時のノーテーション (表記法) は、値 X_{ij} が i さんの第 j 期のスコア、ということ
6223 にしたいと思います。

6224 これを使って毎回の変化の大きさを見積もりたいと思います。前回やった、多水準モデルのことを参考に見
6225 ていきましょう。ポイントは、前は全体平均からの差分で効果を考えていたところが、今回は個人差がある

表 23.2 介入時期とスコアの推移

ID	1 期	2 期	3 期
1	10	5	9
2	9	4	5
3	4	2	3
4	7	3	5

6226 ので個人ごとの平均 μ_i を基準におくところですよ。

6227 最初の人のスコア, $X_{11} = 10, X_{12} = 5, X_{13} = 9$ はそれぞれ, その人の平均 μ_i から考えて,
6228 $\mu_1 + \delta_1, \mu_1 + \delta_2, \mu_1 + \delta_3$ となるはずだと考えましょう。理論通りにならない誤差を含めて考えれば, 先ほどの表 23.2 は表 23.3 のようになっているはずなのです。

表 23.3 介入時期とスコアのモデル式

ID	1 期	2 期	3 期
1	$X_{11} \sim N(\mu_1 + \delta_1, \sigma)$	$X_{12} \sim N(\mu_1 + \delta_2, \sigma)$	$X_{13} \sim N(\mu_1 + \delta_3, \sigma)$
2	$X_{21} \sim N(\mu_2 + \delta_1, \sigma)$	$X_{22} \sim N(\mu_2 + \delta_2, \sigma)$	$X_{23} \sim N(\mu_2 + \delta_3, \sigma)$
3	$X_{31} \sim N(\mu_3 + \delta_1, \sigma)$	$X_{32} \sim N(\mu_3 + \delta_2, \sigma)$	$X_{33} \sim N(\mu_3 + \delta_3, \sigma)$
4	$X_{41} \sim N(\mu_4 + \delta_1, \sigma)$	$X_{42} \sim N(\mu_4 + \delta_2, \sigma)$	$X_{43} \sim N(\mu_4 + \delta_3, \sigma)$

6229

6230 つまり, $X_{ij} \sim N(\mu_i + \delta_j, \sigma)$ ですね。個人のベースライン μ_i に, 各時期の効果 δ_j が加わったものを中心に, 正規分布に従う誤差を纏ってデータになる, という考え方です。ただしここでも注意が必要で, パラメータがいろいろありますが, 効果は相対的なもの, すなわち $\delta_1 + \delta_2 + \delta_3 = 0$ という総和ゼロの制約を考える
6232 必要があります。実質的には自由度 2, すなわち $\delta_3 = 0 - (\delta_1 + \delta_2)$ となることを忘れないでください。

6234 そしてもう 1 つ大事なことが, μ_i は人によって違う個人差成分ですが, この散らばりも正規分布に従うものとして考えることができますね。すなわち平均と散らばりをもった集合体が従う分布としての正規分布です。この平均を ψ , 標準偏差を τ として*3,

6235

$$\mu_i \sim N(\psi, \tau)$$

6237 と考えることにしましょう。これを設計図に書き込んだのが, 図 23.5 です。

6238 このように, 相関があることを個人 i に紐づけられた変数として表現し, 個人差からの差分として対応
6239 のあるモデルを考えることができます。ここで, 個人差に関係なく影響を与えている効果のことを**固定効果 (Fixed Effect)**と呼び, 個人差のように交換可能な要素のことを**変量効果 (Random Effect)**と呼び
6240 ます。また図 23.5 に示されているように, データ生成分布の中に個人差の分布が入れ子になって含まれて
6241 いますね。このようなデータは階層モデルといい, とくに今回のように線形な階層モデルは**階層線形モデル (Hierarchical Linear Model)**とよばれます。この展開については, 次回以降にお話しすることになるので
6242 しょう。

6245 ではこのモデルを実装していきましょう。設計図に沿えば Stan コードは書けるはずですよ。ここには Stan に
6246 与える R のコード例を用意しましたので, 自分なりのコードを書いて検証してみてください*4。

*3 ψ はギリシア文字でプサイと読みます。 τ も同じくギリシア文字でタウと読みます。

*4 これに限らず, コードに正しい書き方はありません。自分なりの書き方で結構です。判断すべきは正しく機能するかどうかであって, 書き方はあくまでも一例にすぎません。

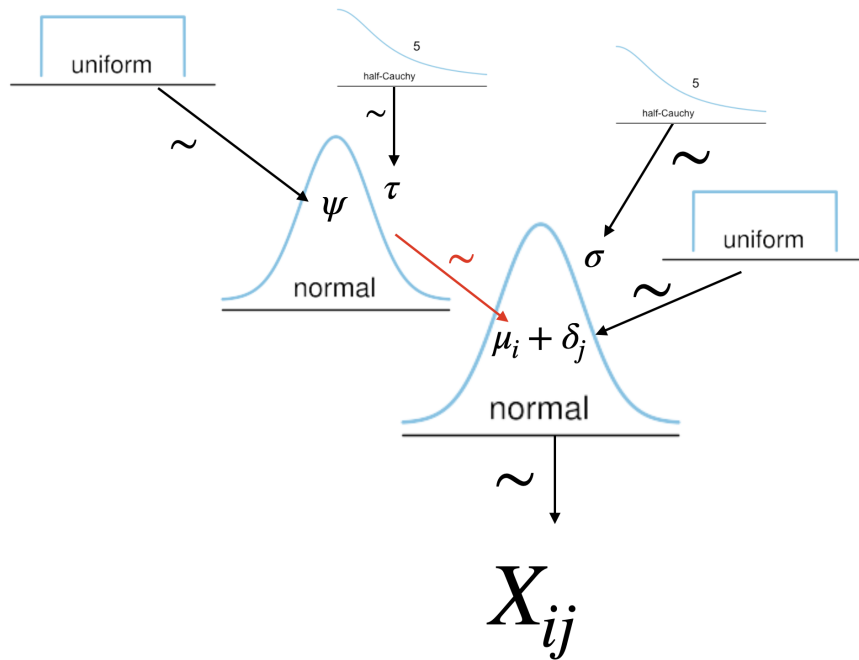


図 23.5 階層的データ生成モデル

code : 23.6 整然データ対応番

```

6247 1 dat_raw <- data.frame(ID = 1:4,
6248 2                       period1 = c(10,9,4,7),
6249 3                       period2 = c(5,4,2,3),
6250 4                       period3 = c(9,5,3,5))
6251 5 # 整然データにする
6252 6 tidy_dat <- dat_raw %>%
6253 7   pivot_longer(-ID) %>%
6254 8   mutate(name = as.factor(name) %>% fct_relevel("period1", "period2")) %>%
6255 9   mutate(cond = as.numeric(name))
6256
6257

```

23.4 課題

Code::23.6 を参考にしながら、表 23.2 の事例を検証するための Stan コードを書いてください。作った Stan ファイルや R コードを提出してください。

第 24 章

モデリングの目から見た検定 5 ; カテゴリカル分布をつかって

ここまでさまざまなモデルを使って、パラメータ間あるいはデータの間の差を検討するということを考えてきました。ここで差分を計算できているということは、数値が連続的であったということもできます。いいかならば、間隔尺度水準以上の尺度水準で得られたデータに対する検証だったわけです。

ところがデータの種類によっては順序尺度水準や名義尺度水準でしか得られないこともあります。試験に合格したのか、失敗したのか。男性か、女性か。47 都道府県のどこ出身なのか。あるいはたとえば、正解率や成功率、相対的な比率などでデータが現されたり考えたりすることもあります。比率は連続的な数字ではありますが、 $\frac{X}{Y}$ のように数え上げによる度数の比較ですから、元は離散的な、質や種類を区別するための数字による (度数に該当するかどうか) 情報だと言えるでしょう。たとえば比率のデータに対して t 検定のような分析をするのは適切ではなく^{*1}、カテゴリカルな分布を考えて検証する必要があります。

また、複数のカテゴリカル変数の組み合わせによってクロス集計表 (cross-tabulation table) あるいは分割表 (contingency table) を利用することもあります。男性の X 割が A 党に、女性の Y 割が B 党に投票した、といった表は、性別 × 支持政党のクロス集計表ですが、これを見ることで性別による支持政党の偏りがあるのかないのか、といったことがわかります。クロス集計表は社会科学的なデータにも多く見られ、カテゴリカルな性質の組み合わせについての情報を提供してくれます。この組み合わせに偏りがあるのかないのか、と言ったことを考える際にも、確率モデルとしてはカテゴリカルな分布を必要とします。

24.1 離散的な分布

それでは代表的な離散変数についての確率分布を見ていきましょう。離散分布の場合、確率関数は密度 density ではなく、質量 mass、つまり確率質量関数 (Probability Mass Function) と呼ばれることに注意してください。カテゴリの中に含まれるかどうか、どれぐらいの量がそこに含まれているかということを直接表しているからです。

■ベルヌーイ分布 離散的な分布の基本はベルヌーイ分布 (Bernoulli Distribution) です。この分布から出てくる数字は 0 か 1 であり、確率 θ で 1 が、 $1 - \theta$ で 0 が出る分布です。たった 2 つの状態しか区分しませんが、正答と誤答、賛成と反対、男性と女性、生と死などさまざまなもののメタファーとして利用できるため、応用範囲はむしろ広いと言えるでしょう。パラメータは θ ひとつですので、わかりやすいですね。Stan

^{*1} t 検定は正規分布するデータを前提にしている検定方法・確率分布であり、比率のデータはとうてい正規分布とは考えられないからです。正規分布の定義域は $-\infty$ から ∞ なのに対し、比率は 0 から 1 でしかないことから明らかです。

6288 での関数は `bernoulli` と書き*2, $y \sim \text{bernoulli}(\theta)$; のように使います。

6289 今後出てくることとなりますが、ロジスティック関数と組み合わせて使うことが多い分布です。ロジス
6290 ティック関数は、連続的な数字を 0 から 1 の範囲に変換する関数で、変換したものをベルヌーイ分布
6291 のパラメータにする合わせ技が非常に便利だからです。たとえば学力などは正規分布すると考えられ、
6292 $-\infty$ から ∞ の値を取りうるようになりますが、その能力を反映してテストに正答するか誤答するか、
6293 ということ考えると学力 θ をロジスティック関数に入れた ($\text{logistic}(\theta)$) ものをベルヌーイ分布にい
6294 れる ($X \sim \text{bernoulli}(\text{logistic}(\theta))$) からです。この合わせ技はよく使われるので、Stan では
6295 `bernoulli_logit` という関数があるほどです。

6296 ■二項分布 パラメータ θ をもつベルヌーイ分布に従うコイントスを N 回やったときに、何回表が出る
6297 か、というときに使うのが二項分布 (Binomial Distribution) です。10 問中何問正解するかとか、
6298 N 個提示した刺激のうち覚えていたのはいくつとか、 N 回シュートして入ったのは何回かといった、割
6299 合や比率に関係する分布です。パラメータは N と θ の 2 つあり、Stan での関数は `binomial` です。
6300 $y \sim \text{binomial}(N, \theta)$; のように書きます。

6301 ■カテゴリカル分布 ベルヌーイ分布や二項分布はコイントスの「表」のことしか考えていませんでした
6302 (裏は「表ではない」という考え方)。これに対してカテゴリカル分布はサイコロの出目のように、1,2,3,4,5,6
6303 のどれが出るか、という多段階の乱数発生を考える分布です。こん関数のパラメータはサイズ K のベクト
6304 ルで、たとえばサイコロの場合は $\theta = (1/6, 1/6, 1/6, 1/6, 1/6, 1/6)$ のようになります。Stan での関数は
6305 `categorical` と書き、使う時は $y \sim \text{categorical}(\theta)$; のようにします。このとき θ はベクトル
6306 で宣言されていなければなりません。

6307 ■多項分布 カテゴリカル分布と似ているようですが、この分布は N 回サイコロを降った時の各出目の回数
6308 であり、1 が k_1 回、2 が k_2 回、 \dots 、6 が k_6 回出た、というようにカウントしたものを出力します。ABO 型
6309 の血液分類で、サンプルの中に A 型が何人、B 型が何人 \dots というときに、この母比率を推定する時は多項
6310 分布のパラメータを推定する、ということになります。二項分布の多変量版だと考えてもいいかもしれません
6311 (松浦, 2016)。この分布に与えるパラメータはベクトル θ であり、Stan では $y \sim \text{multinomial}(\theta)$;
6312 のように使います。総数 N はデータベクトル y の要素の総和 ($N = \sum_{i=1}^K y_i$) ですから指定しなくても構いま
6313 せん。ここで注意すべきは、 θ はサイズ K のベクトルであり、すべての要素を足し合わせると 1 になる
6314 必要があります。Stan では合計 1 のベクトル専用の型があり、それは `simplex` と呼ばれます。

6315 24.2 χ^2 検定

6316 カテゴリカルなデータ、度数分布などをどうやって研究の時に使うのか、具体的な例とともに見ていきましょ
6317 う。ベイズ推定の話に入る前に、帰無仮説検定の話から先に考えてみたいと思います。

6318 表 24.1 には、3 つの携帯電話キャリアのユーザ数をとある大学で調査した例です。合計 123 名から回答
6319 を得て、それぞれ 51 名、45 名、27 名ということが判明しました。これを見て「A 社は多いね、学生人気だ
6320 ね」などと解釈してもいいのですが、これでは推測統計学の域を出ていませんね。違うサンプルを対象にすれ
6321 ば違う度数になる可能性があるからです。そこで推測統計学的な発想をします。

6322 ここで検証したいのは「どこかのカテゴリが多く、どこかのカテゴリは少ない」、つまり「散らばりに偏りがあ
6323 る」ということです。逆に考えると帰無仮説が出てきます。つまり「散らばりに偏りはない」が帰無仮説であり、

*2 対数尤度で書く場合は `bernoulli_lpmf` となります。lpmf は log probability mass function の略です。

表 24.1 携帯キャリア調査

キャリア	A 社	D 社	S 社	合計
観測度数	51	45	27	123
期待度数	41	41	41	

6324 対立仮説は「偏りがある」ということになります。散らばりに偏りが無い、という帰無仮説の世界の元では、期
 6325 待度数は総数を均等に割った値になるはずで、表 24.1 には期待度数 (Expected frequency) という
 6326 行を設けました。ちなみに得られたデータそのものは観測度数 (Observed frequency) と言います。

6327 今回のデータが帰無仮説の元ではどれぐらいあり得ない数字だったのでしょうか。これを考えるために、
 6328 χ^2 値を計算します。 χ はギリシア文字のカイであり、 x (エックス) とは違うので注意してください。 χ^2 値は
 6329 カイ二乗値と読みます。この値は次のように計算します。

$$\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$$

6330 ここで O_i は観測度数、 E_i は期待度数です。今回の場合は次のような数字になります。

$$\chi^2 = \frac{(51 - 41)^2}{41} + \frac{(45 - 41)^2}{41} + \frac{(27 - 41)^2}{41} = 2.439 + 0.390 + 4.780 = 7.609$$

6331 計算式から明らかなように、期待度数からのずれの大きさを表現していますから、期待通り = 帰無仮説通
 6332 りであればこの値は 0 になります。今回は $\chi^2 = 7.609$ ですが、この統計量がどれぐらい出てきやすいかとい
 6333 うことは、自由度 2 の χ^2 分布をつかって検証できます。ここでの自由度は選択カテゴリ数 - 1 です*3。自由
 6334 度 2 の χ^2 分布から 7.609 という実現値が出てくる確率は $p = 0.022$ で 5% より小さいですから*4、統計的
 6335 に有意な偏りがある、と判断できます。

6336 χ^2 は偏りを検証するための数字です。今回は均等である、という帰無仮説から算出しましたが、母集団
 6337 比率が事前にわかっているのであればそれを用いて、手元のサンプルがどれぐらい偏っているかを検証す
 6338 るということも可能です。たとえば日本人の血液型はおよそ A 型が 40%、O 型 30%、B 型 20%、AB 型
 6339 10% といわれていますから、期待度数をこの比率で割り振るとサンプルが偏っていたかどうかの検証になり
 6340 ます*5。他にも (比率によらず) 理論的な度数がわかっていることがあれば、どれぐらい合致しているかを検
 6341 証できます*6。

6342 また同様のロジックで、クロス集計表の検定を行うこともできます。クロス集計表の場合は**独立性の検定**と
 6343 呼ばれることもあります。偏りがあれば関係がある、独立ではないということですね。

6344 例として、表 24.2 には、先ほどの携帯キャリア調査の結果に男女の区別もつけてみました。これを見て「男
 6345 性は D 社が好きで、女性は A 社が好き」といった偏りがあると判断できるかどうかを検証できます。

6346 各セルの期待度数は行の周辺度数 × 列の周辺度数 ÷ 総度数で計算でき、各観測度数から期待度数を
 6347 引いて二乗したものを、さらに期待度数で割ったものを足し合わせることで χ^2 にするのは同じです。これが
 6348 (行の数 - 1) × (列の数 - 1) の自由度をもつ χ^2 分布に従いますので、今回は $\chi^2 = 6.4326$, $df = 2$, $p =$

*3 合計数が 123 と決まっていますから、最初の二つのカテゴリに入る数字が決まれば残りの一つは自動的に決まります。自由に値
 が変えられるのは二つまでなのです。

*4 この値は R で `1-pchisq(7.609, df=2)` とすることで得られます。もっと直接的にするなら、`chisq.test(c(51,45,27))` と
 すると χ^2 検定が実行されます。

*5 母比率の検定と言います。

*6 **適合度**の検定と言います。モデルフィットの指標としても使われ、たとえば第 13 講で説明した**構造方程式モデリング**の指標など
 でも用いられています。

表 24.2 携帯キャリア調査と性別の関係

キャリア	A 社	D 社	S 社	合計
男性	21	30	13	49
女性	30	15	14	74
合計	51	45	27	123

6349 0.0401 であることから有意 (に偏っている), と判断できます*7。これは分散分析の時と同様に, どこかに偏り
6350 があるという全体的な傾向を示しただけであることに注意が必要です。

6351 このようにして度数の検定を行うことができます。これに対して, データ生成メカニズムを考える場合は, あ
6352 る比率に沿って度数が出てきていると考えますから, もっと直接的に母比率の検定を考えることができます
6353 し, MCMC によるアプローチを使うと「どこが, どの程度大きいと言えるか」といったことについても検証で
6354 きます。次にその方法を見ていきましょう。

6355 24.3 カテゴリカル分布のモデリング

6356 ではまず, 表 24.1 にあげた携帯キャリア調査の結果を, データ生成モデルから考えてみましょう。設計図に
6357 するまでもない感じで, 51:45:27 というデータの比率から考えられる, 母比率ベクトル $\pi = \{\pi_1, \pi_2, \pi_3\}$
6358 を考え, そこからデータベクトル y が多項分布に従って出てくる, と考えるわけです。

6359 Stan のコードは Code:24.1 のようになります。ポイントは特殊なベクトル simplex で π を宣言している
6360 ところでしょうか。またデータベクトルは `vector[K] X`; としたいところですが, モデルからいってここは `int`
6361 型である必要があり, 型指定をしたベクトルというのがないので, 配列にしました。

code : 24.1 多項分布からのデータ生成モデル

```
6362
6363 1 data{
6364 2     int K;
6365 3     array[K] int X;
6366 4 }
6367 5
6368 6 parameters{
6369 7     simplex[K] pi;
6370 8 }
6371 9
6372 10 model{
6373 11     X ~ multinomial(pi);
6374 12 }
6375
```

6376 これを次のような R コードで実行してみましょう (Code:24.2)

code : 24.2 多項分布のパラメータ推定

```
6377
6378 1 # rstan の場合
6379 2 model <- rstan::stan_model("categorical1.stan")
6380 3 fit <- sampling(model,
6381 4   data = list(K = 3, X = c(51, 45, 27))
6382 5 )
```

*7 この計算は `chisq.test(matrix(c(21,30,13,30,15,14),ncol=3,byrow=T))` で算出しました。

```

6383 6 # cmdstanr の場合
6384 7 model <- cmdstanr::cmdstan_model("categorical1.stan")
6385 8 fit <- model$sample(
6386 9   data = list(K = 3, X = c(51, 45, 27)),
6387 10  chains = 4,
6388 11  parallel_chains = 4,
6389 12  refresh = 500
6390 13 )
6391

```

6392 構造も何もないモデルですから、推定はすぐに終わると思います。51 : 45 : 27 = 0.4146 : 0.3658 : 0.2195
 6393 ですから、ほぼ標本分布と同じような形で事後分布が推定されました (図 24.1)。これで終わり、といった

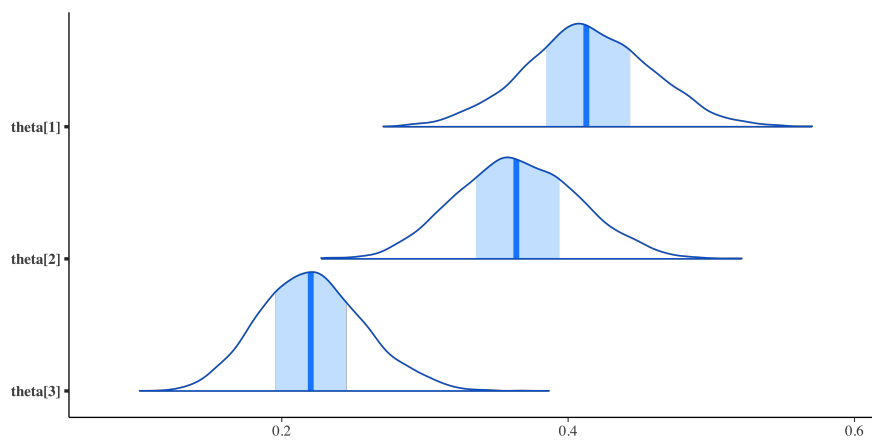


図 24.1 比率ベクトルの推定結果

6393
 6394 らそれまでなのですが、仮説検定のように考えてみましょう。三社それぞれが均等に選ばれている、つまり
 6395 $\pi_1 = \pi_2 = \pi_3$ というのはあり得そうにないですが (ピッタリ一致するはずはない)、たとえば

- 6396 • A 社よりも D 社の方が好まれている
- 6397 • S 社よりも D 社の方が好まれている
- 6398 • D 社が最も好まれている

6399 といった仮説は生成量を使って検証できそうです。1 つ目は $\pi_1 > \pi_2$, 2 つ目は $\pi_2 > \pi_3$ が成立する割合
 6400 から考えれば良いからです。3 つ目の仮説はこの 2 つと同じじゃないか、と思うかもしれませんが、正確には
 6401 $\pi_2 > \pi_1$ かつ $\pi_2 > \pi_3$ という両方の可能性が成立していること、と考える必要があります。こうした「A かつ
 6402 B」のような命題は**連言命題**といいます。

6403 これを検証する生成量の作り方は、次のようになります。他にもいろいろ考えられそうですね。

code : 24.3 多項分布からのデータ生成モデル

```

6404 1 ... 前略...
6405 2 generated quantities{
6406 3   int<lower=0,upper=1> FLG1;
6407 4   int<lower=0,upper=1> FLG2;
6408 5   int<lower=0,upper=1> FLG3;
6409 6
6410 7   if(pi[1] > pi[2]){ FLG1 = 1; } else { FLG1 = 0; }
6411 8   if(pi[2] > pi[3]){ FLG2 = 1; } else { FLG2 = 0; }

```



```

6413 9      if(pi[2] > pi[1] && pi[2] > pi[3]){ FLG3 = 1; } else { FLG3 = 0; }
6414 10 }
6415

```

24.3.1 対応のない二変数の場合

6417 続いてクロス表の検定に行きましょう。表 24.2 にあるように、携帯キャリアの選び方と性別の関係をみたい
6418 と思います。男性と女性は独立したカテゴリですから、先ほどの多項分布モデルが 2 つ同時にあることと同じ
6419 です。

6420 男性は 21:30:13 という観測度数が得られていますが、これは $\pi_M = c(\pi_1^M, \pi_2^M, \pi_3^M)$ という母比率を
6421 持っており、また女性は 30:15:14 という観測度数から、 $\pi_F = c(\pi_1^F, \pi_2^F, \pi_3^F)$ という母比率だったのではな
6422 いか、と推測することになります。

6423 「男性と女性とで携帯キャリアの選択率に違いがあるか」ということが検証したい命題になるかと思いき
6424 ますが、何を持って違いがあるとするか、ということを厳密に考えて生成量を作らなければなりません。たとえば、

- 6425 • 男性は女性よりも A 社を好む $\rightarrow \pi_1^M > \pi_1^F$
- 6426 • S 社は女性の方が多く選んでいる $\rightarrow \pi_3^M < \pi_3^F$
- 6427 • 男性は D 社が一番好きで女性は A 社が一番好き $\rightarrow \pi_2^M > \pi_1^M$ かつ $\pi_2^M > \pi_3^M$ かつ $\pi_1^F > \pi_2^F$
6428 かつ $\pi_1^F > \pi_3^F$

6429 というように考えられます。とくに 3 番目の命題は、「一番好き」というのを「他のどの群と比べても大きいとい
6430 う条件が同時に成り立っている」と考える必要がありますし、「男性は…で、女性は…」というときの「で」を
6431 論理的には「かつ」、確率的には「同時に」という意味で捉えないといけないところがポイントです。

6432 このように、生成変数を使うときさまざまな仮説を検証することが可能です。集計表に関するデータは、官公
6433 庁が公開しているものなどを含め、身の回りのさまざまなところで目にできます。これらの資料を見た時に、標
6434 本統計量だけでなく母比率にまで思いを馳せて、検証可能な仮説を色々考えてみると良いでしょう*8。

6435 ところでこの生成量を使った連言命題の検証については、今回のデータから考えられる事後分布に依存
6436 していることを忘れないようにしてください。今回のデータから言える「命題が成立する確率」は、事前分布と
6437 データに基づいている結果ですから、「仮説が正しい確率」とまで言い切るのは危険です。あくまでもモデル
6438 に基づく推定であることを忘れないようにしましょう。

24.4 κ 係数の算出

6440 先ほどのクロス集表の例では、男性と女性という独立した群から作られていましたので、2 つの母集団を
6441 別々に考えることができました。しかしたとえば、「法案 A,B それぞれに対して賛成か反対か」といった調査を
6442 した場合には、「A にも B にも賛成」「A には賛成, B には反対」「A には反対, B には賛成」「A にも B にも
6443 反対」という 2×2 のクロス集計表が得られますが、これは一人の人間が 2 つの質問に答えていますので**対
6444 応のある**データになります。対応のあるデータの場合は、先ほどと同じように独立した確率分布を考えるわけ
6445 には行きません。

*8 たとえばテレビのバラエティ番組や情報番組においても、該当アンケート結果についてコメントが色々物申すことが少なくあ
りません。「エビフライのしっぽは食べるのか、残すのか」とか「焼き鮭の皮は食べるのか残すのか」といった些細なことについて、
無作為抽出とは思えないような該当インタビューをし、僅差でも多数派を勝者と見立てて騒ぎ立てる様をみる度に、筆者は「せめて
ちゃんと検定、推定してから結論づけるべきだ」と思い、常々腹立たしく感じていました。しかし妻に「そういう些細なことを、僅
差でも大袈裟にいうことで、お茶の間に盛り上がる話題を提供しているのであって、学術報告ではない」と諭されたことがあって、
なるほどそういうものかと納得した次第です。

6446 また「ある病院にいったら健康だと診断されたが、別の病院に行ったら異常が発見された」とか、「ある専門
6447 家の見立てでは問題のある行動といわれたが、別の専門家の見立てでは問題ないといわれた」、「ある映像
6448 を見て判断する課題で、審査者 A と B の判断が一致したりしなかったりする」といったシーンを考えてみてく
6449 ださい。これらのシーンでは、 2×2 のクロス集計表が表 24.3 のようになっていると考えられます。ここでは
6450 Yes/No としてありますが、Hit/Miss でも健康/異常でもいいのですが、とにかく 2 つの判断が合致したか
6451 どうか、ということが問題になるシーンです。

表 24.3 対応のあるカテゴリ判断

		A		
		Yes	No	
B	Yes	a	c	
	No	b	d	
				n

6452 このような時も、行と列の判断が独立しているとは言えません。むしろどの程度重複しているかのほうが問
6453 題になるわけです。こういったときは、判断が一致した程度を**カツバ係数 (kappa coefficient)** で表現で
6454 きます*⁹。この係数は、観測された一致度、すなわち

$$p_o = \frac{(a + d)}{n}$$

6455 と、偶然の一致度

$$p_e = \frac{(a + b)}{n} \frac{(a + c)}{n} + \frac{(b + d)}{n} \frac{(c + d)}{n} = \frac{(a + b)(a + c) + (b + d)(c + d)}{n^2}$$

6456 から、次のように計算されます*¹⁰。

$$\kappa = \frac{p_o - p_e}{1 - p_e}$$

6457 この式の分子は「合致した割合から偶然の一致をひいたもの」であり、分母は「偶然ではない割合」になっ
6458 ていることから、一致の程度を表す指標として考えられるのです。この κ 係数は -1 から $+1$ の範囲の値を取
6459 ります。

6460 このように、2 つの群わけが独立ではない場合の確率はどう考えればいいでしょうか。 $a : b : c : d$ を
6461 $\pi_1 : \pi_2 : \pi_3 : \pi_4$ と考えるのでは 4 つのセルが独立だと考えていることになってしまいます。そうではなく、順
6462 番にまず A が確率 α で Yes と判断する、ということを考えましょう。No と判断する確率は $1 - \alpha$ です。次に
6463 B の判断ですが、A が Yes と判断した時に B が Yes と判断する確率を β としましょう。また A が No と判
6464 断した時に B も No と判断する確率を γ と考えます。

6465 そうすると、両者が Yes と答える確率 $\pi_a = \alpha\beta$ 、A が No で B が Yes と答える確率 $\pi_b = (1 - \alpha)\beta$ 、A
6466 が Yes で B が No と答える確率 $\pi_c = \alpha(1 - \gamma)$ 、両者が No と答える確率 $\pi_d = (1 - \alpha)\gamma$ のようにあらわ
6467 すことがで切ようになります (図 24.2)。 κ 係数はこれら $\pi_a \sim \pi_d$ の組み合わせで計算できますから、生成
6468 量を使えば算出できますね。

6469 それではこれをコードにしていきましょう。確率ですから、変数の範囲は 0 から 1 に制限し、事前分布もこ
6470 の範囲の一樣分布にしています。

*⁹ カツバ係数は河童ではなく、ギリシア文字の κ にあたる κ です。

*¹⁰ A が Yes と判断する割合 \times B が Yes と判断する割合を、A が No と判断する割合 \times B が No と判断する割合に足している数字が p_e です。割合 (= 確率) の積なので二つの条件が偶然に一致する割合を、Yes-Yes のケースと No-No のケースで算出し、合わせたものと理解すればいいでしょう。

		Aの判断	
		Yes	No
Bの判断	Yes	$\pi_a = \alpha\beta$	$\pi_c = (1 - \alpha)(1 - \gamma)$
	No	$\pi_b = \alpha(1 - \beta)$	$\pi_d = (1 - \alpha)\gamma$

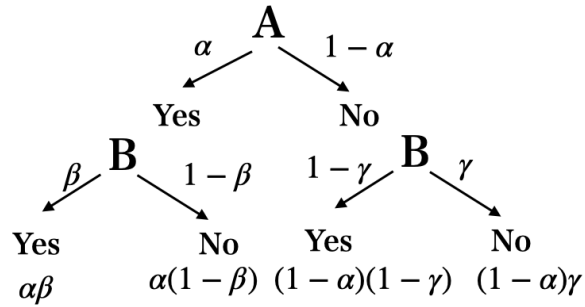


図 24.2 対応のあるカテゴリ判断の確率

code : 24.4 κ 係数を計算するモデル

```

6471 1 data{
6472 2     array[4] int Y;
6473 3 }
6474 4
6475 5 parameters{
6476 6     real<lower=0,upper=1> alpha;
6477 7     real<lower=0,upper=1> beta;
6478 8     real<lower=0,upper=1> gamma;
6479 9 }
6480 10
6481 11 transformed parameters{
6482 12     simplex[4] Pi;
6483 13     Pi[1] = alpha * beta;
6484 14     Pi[2] = (1-alpha) * (1-gamma);
6485 15     Pi[3] = alpha * (1-beta);
6486 16     Pi[4] = (1-alpha) * gamma;
6487 17 }
6488 18
6489 19 model{
6490 20     Y ~ multinomial(Pi);
6491 21     alpha ~ uniform(0,1);
6492 22     beta ~ uniform(0,1);
6493 23     gamma ~ uniform(0,1);
6494 24 }
6495 25
6496 26 generated quantities{
6497 27     real po;

```

```

6499 28   real pe;
6500 29   real kappa;
6501 30   po = Pi[1] + Pi[4];
6502 31   pe = (Pi[1]+Pi[2])*(Pi[1]+Pi[3]) + (Pi[2]+Pi[4])*(Pi[3]+Pi[4]);
6503 32   kappa = (po-pe)/(1-pe);
6504 33 }
6505

```

6506 カテゴリカルな分布はこのように度数に関するモデルを考えることができますし、この分布を応用して潜在的変数がカテゴリカル分布に従うと考えると、分類わけを表現することもできます。さまざまな応用例が思いつくのではないのでしょうか。

6509 24.5 課題

6510 次の2つの課題について、考察を導くための計算をする R/Stan コードとともに、回答を提出してください。Rmd ファイルでの提出が望ましいですが、メモやコメントアウト、Word ファイル、Google ドキュメントなどでの提出も可とします。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

6515 ■多項分布と連言命題 とある棒状のお菓子にはさまざまな味のバリエーションがあるが、コーンポタージュ味、チーズ味、めんたい味、野菜サラダ味、たこ焼き味あたが人気上位に入らしい。そこで実験を行なって、最も美味しい味に投票してもらったところ、コーンポタージュ味が105票、チーズ味が80票、めんたい味が75票、野菜サラダ味が60票、たこ焼き味が45票であった。この時、次の仮説を検証するコードを書き、考察してください。

- 6520 ● コーンポタージュ味がチーズ味よりも好まれているといえるかどうか
- 6521 ● コーンポタージュ味がめんたい味よりも好まれているといえるかどうか
- 6522 ● コーンポタージュ味が他のどの味よりもこのまれているといえるかどうか

6523 ■一致率の判断 新型コロナウイルスに対するスピード臨床検査が開発された。新型コロナに罹っていることがわかっている40名のうち、新しいスピード検査法では34名を陽性と判断できたが、6名はコロナであると判断できなかった。また新型コロナに罹っていないことがわかっている200名の、新しいスピード検査法では190名はコロナでないと正しく判断できたが、10名はコロナであると間違えて判断してしまった。この検査法のデータから一致率を計算し、新しいスピード臨床検査法が有用であるといえるかどうか、事後分布に基づいて考察してください^{*11}。

^{*11} この課題は Lee and Wagenmakers (2013) の Pp.59, 練習問題 5.3.1 を参考にしています

第 25 章

一般化線形モデル

25.1 一般線形モデル

ここまで検定と線形モデルが同じ形をしていることについては、何度か指摘してきました。検定モデルは、

$$Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$$

という形をしており、ここでとくに X_i が統制群に属するか、実験群に属するかを 0/1 で表現したものであったことを再確認しておきましょう。たとえば二水準モデルの場合、統制群の場合 $X_i = 0$ ですから、データは $Y_i = \beta_0 + 0 \times \beta_1 + \varepsilon_i = \beta_0 + \varepsilon_i$ であり全体平均に誤差がついただけの値になります。実験群の場合は $Y_i = \beta_0 + 1 \times \beta_1 + \varepsilon_i = \beta_0 + \beta_1 + \varepsilon_i$ になりますから、この β_1 の部分が効果の大きさを表しているのです。

水準数が増えても同じことです。3 水準モデルの場合はベクトルを使って表記した方がわかりやすいので、データセットを \mathbf{Y} と表すと、

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (25.1)$$

であり、詳しく書くと

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 1 & -1 & -1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

のようになります*1。ここで行列 \mathbf{X} はその群に属しているかどうか（その群の効果が発動するかどうか）を 0/1 で表している**デザイン行列 (design matrix)** であり、有無という状態に一一対応した**名義尺度水準**の数字です。説明変数が名義尺度水準になっているだけで、これが間隔尺度水準以上の**量的なデータ**であれば、**回帰分析 (regression analysis)** になるわけです。回帰分析は一般に、

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \varepsilon$$

あるいはベクトルで書くと

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (25.2)$$

となります。この通り式で書くと 25.1 と 25.2 にはなんら違いがないことがわかりますね。

このように**要因計画と線形モデルは同じである**、ということを含意して、これをまとめて**一般線形モデル (General Linear Model)** と言います*2。

*1 $\beta_3 = 0 - (\beta_1 + \beta_2)$ となっていることに注意。

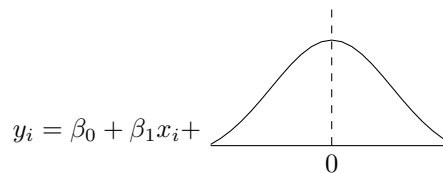
*2 一般、というのは同じだという意味があって、辞書には「私は彼女と同じい罪を犯したも一般だ」といった用法が載っています。

6549 25.1.1 回帰分析の確率モデル

6550 ところで、この一般線形モデルでは誤差が正規分布に従うという仮定を置いていました。すなわち、
6551 $\varepsilon \sim N(0, \sigma)$ というわけです*3。線形モデルの部分、すなわち $\hat{Y}_i = \beta_0 + \beta_1 X_i$ のところは確率的ではありません
6552 せんから、データ全体では次のようになっていることになります。

$$y_i = \underbrace{\beta_0 + \beta_1 x_i}_{\text{確率的でないところ}} + \underbrace{e_i}_{\text{確率的}}$$

6553 誤差 e_i は平均 0 を中心に確率的に散らばりますが、その前のところは確率的ではありませんので、これを
6554 組み合わせた式全体としては、最後に確率的な散らばりがひっついているような形になります。数式ではあり
6555 ませんが、イメージでいうと次のような形でしょうか。



6556

6557 このことから、モデル全体の確率的挙動は、

$$Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma) \quad (25.3)$$

6558 ということになります。これが正規分布の確率モデル的表現です。一般線形モデルですから、説明変数 X が
6559 連続的でも離散的でも同じモデルです。また、この正規分布は平均と SD でその形状が定められるのでし
6560 た。平均は分布の位置を表すので、別名**位置母数 (location parameter)** といいます。ちなみに SD のこ
6561 とは別名**尺度母数 (scale parameter)** と呼ぶことがあります。一般線形モデルは平均的な位置を線形で
6562 予測して、データはそれに付随した誤差が幅 σ でぶれる、と考えます。そういう意味では、図 25.1 のように、
6563 回帰線にそって誤差分布があり、結果として実現値が得られているイメージ図のほうがわかりやすいかもし
6564 れません。

6565 これを確認した上で、では回帰分析をベイズ的に推定するコードを書いてみましょう。設計図も一応準備し
6566 ておきます。図 25.2 では、 μ が $\beta_0 + \beta_1 X_i$ で作られているので、矢印のところは \sim ではなく $=$ になってい
6567 ます。また、未知数 β_0, β_1 については平均 0、SD100 のとても幅の広い正規分布を置いています。一様分
6568 布でもいいのですが、原理的に $\pm\infty$ までの幅があり、どこか 1 つの値に収束するだろう = 単峰性があるだろ
6569 う、ということで正規分布を選んでいきます。誤差 SD σ については、半コーシー分布*4 設計図をそのままコー
6570 ドに起こせば一丁上がりです。事後予測分布も生成するコードを書きました。

code : 25.1 線形モデルの Stan コード

```
6571 1 data{
6572 2   int N;
6573 3   array[N] real X;
6574 4   array[N] int Y;
6575 5 }
```

*3 誤差の平均はゼロであるのは、偶然誤差の仮定から必然的です。

*4 正の値しか取らないように、コーシー分布を 0 のところで半分に折り畳んだ形にした分布です。Stan 上では `lower=0` の制限をかけることで正の範囲からしかサンプリングしない、というように表現します。

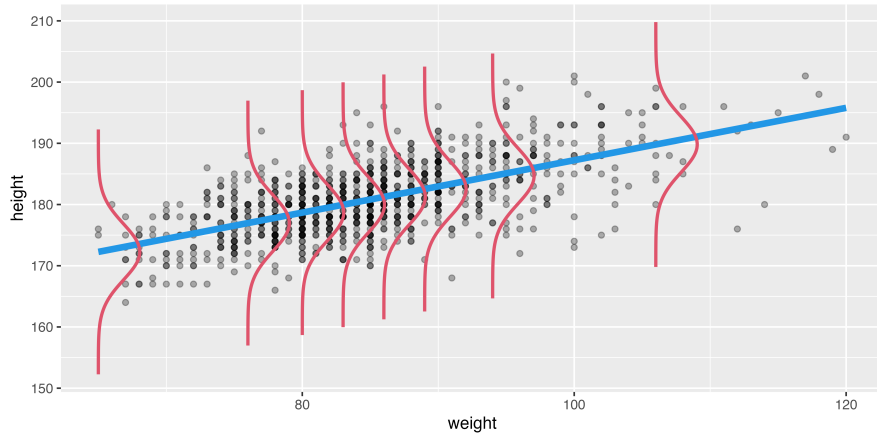


図 25.1 回帰線に伴う誤差分布

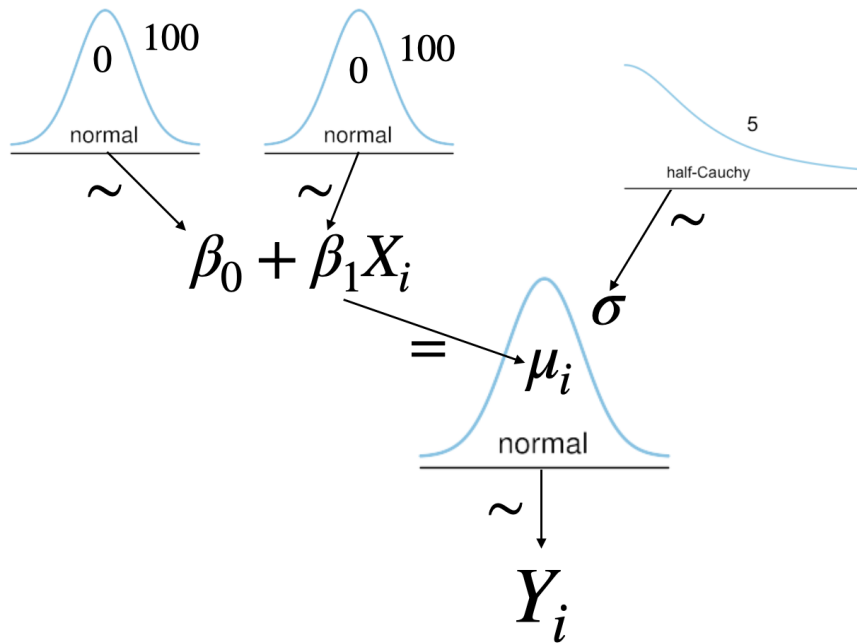


図 25.2 線形モデルの設計図

```

6577 6
6578 7 parameters{
6579 8   real beta0;
6580 9   real beta1;
6581 10  real<lower=0> sig;
6582 11 }
6583 12
6584 13 transformed parameters{
6585 14   array[N] real mu;
6586 15   for(i in 1:N){
6587 16     mu[i] = beta0 + beta1 * X[i];

```

```

6588 17   }
6589 18 }
6590 19
6591 20 model{
6592 21   // model
6593 22   for(i in 1:N){
6594 23     Y[i] ~ normal(mu[i], sig);
6595 24   }
6596 25   // prior
6597 26   beta0 ~ normal(0,100);
6598 27   beta1 ~ normal(0,100);
6599 28   sig ~ cauchy(0,5);
6600 29 }
6601 30
6602 31 generated quantities{
6603 32   real predX[N];
6604 33   for(i in 1:N){
6605 34     predX[i] = normal_rng(mu[i], sig);
6606 35   }
6607 36 }
6608

```

6609 これがどのような挙動をするのか、実際のデータセットを使って検証してみましょう。サンプルデータとして、baseball2020.csv を使います。このデータセットはプロ野球データ Freak^{*5}さんから取ってきたものです^{*6}。

6612 この野球データセットにはさまざまな情報が含まれていますが、身長と体重のデータだけ取り出して回帰分析をしてみましょう。身長を体重から予測することにします。ここで身長は正規分布に従っていると考えられますから、一般線形モデルが適切ですね。データの散布図を書いて、線形関係にありそうかどうか確認しておきます (図 25.3)。

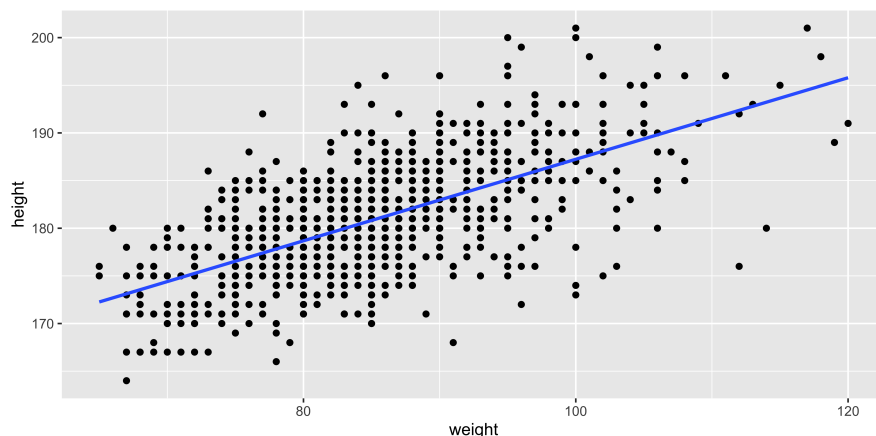


図 25.3 身長を体重で回帰するモデル

6616 そこそこ相関がありそうですね。では回帰モデルを当てはめてみましょう。

*5 <https://baseball-data.com/>

*6 データセットはシラバスのサイトでも提供していますし、著者の個人サイト (<https://kosugitti.github.io/kosugitti10/index.html>) にはスクレイピングのコードを公開しています。

code : 25.2 一般線形モデルを実行するコード

```

6617 1 dat <- read_csv("baseballDecade.csv") %>% filter(Year=="2020年度")
6618 2 dataSet <- list(N = NROW(dat), Y = dat$height, X = dat$weight)
6619 3 ## cmdstanrで実行する場合
6620 4 modelC <- cmdstan_model("LM.stan")
6621 5 fit <- modelC$sample(
6622 6   data = dataSet,
6623 7   chains = 4,
6624 8   parallel_chains = 4,
6625 9   iter_warmup = 1000,
6626 10  iter_sampling = 2000
6627 11 )
6628 12 ## 実行後のファイルをstanfit形式に置き換えておく
6629 13 fit.stanfit <- fit$output_files() %>% rstan::read_stan_csv()
6630 14 ## 事後予測を取り出す
6631 15 predY <- rstan::extract(fit.stanfit)$predY
6632 16 # 事後予測分布の描画
6633 17 bayesplot::ppc_dens_overlay(y = dataSet$Y, yrep = predY[1:10, ])
6634 18 bayesplot::ppc_intervals(
6635 19   y = dataSet$Y,
6636 20   yrep = predY,
6637 21   x = dataSet$X,
6638 22   prob = 0.5,
6639 23   prob_outer = 0.95
6640 24 )
6641
6642

```

6643 ■コード解説

6644 1行目 データを読みこみます。

6645 2行目 stan に与えるデータセットをリストで作ります。

6646 3-11行目 cmdstanr でコンパイルする例です。rstan パッケージを使っても構いません

6647 13行目 実行後のファイルを stanfit 形式に置き換えています。rstan パッケージを使っている人はこの行
6648 を実行する必要がありません。

6649 15行目 事後予測分布の一部です。予測値は MCMC サンプル数 (ここでは 8000 点) × サンプルサイズ
6650 (野球選手 959 人) ありますが、このうち最初の 10 人についての実際のデータ y と予測値 yrep を
6651 重ねて表示しています。全員分表示させると非常にサイズが大きくて重たいからです。

6652 16-22行目 データの散布図に予測値とその 50% および 95% 確信区間をプロットしました。

6653 今回の分析結果を単純に表示するなら、次のようになります (出力 3)。

cmdstanr の出力 3: ベイズ推定の結果

variable	mean	median	sd	mad	q5	q95	rhat	ess_bulk	ess_tail
beta0	144.43	144.42	1.46	1.46	141.97	146.83	1.00	2344	2554
beta1	0.43	0.43	0.02	0.02	0.40	0.46	1.00	2339	2507
sig	4.57	4.57	0.11	0.10	4.40	4.75	1.00	3171	2952

6654
6655 これまで習ってきた**最尤法 (Maximum Likelihood method)** の結果と出力の違いを比較しましょう
6656 (出力 25.1)。

R の出力 25.1: ML 推定の結果

```

Call:
lm(formula = height ~ weight, data = dat)

Residuals:
    Min       1Q   Median       3Q      Max
-16.3680  -2.8245  -0.0963   3.0368  14.9037

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept) 144.45986    1.44682   99.85  <2e-16 ***
weight       0.42775    0.01693   25.27  <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.573 on 957 degrees of freedom
Multiple R-squared:  0.4002, Adjusted R-squared:  0.3996
F-statistic: 638.6 on 1 and 957 DF,  p-value: < 2.2e-16

```

6657

6658 最尤法では、切片が 144.45986、傾きが 0.42775 という点推定だけが表示される形になっていますが、
 6659 MCMC を使ったベイズ推定では切片の平均 (EAP 推定値) が 144.43, 中央値 (MED 推定値) が
 6660 144.42, 95% 確信区間が [141.97,146.83] となっています。同じく傾きの平均が 0.43, 中央値も 0.43, 95%
 6661 区間が [0.40,0.46] となっています。結果が分布として得られますから、初めから幅を持った推定になってい
 6662 ることが確認できます。

6663 事後予測分布 (図 25.4) をみると、いくつか予測幅から外れているデータ点はありますが、概ね回帰線の
 6664 範囲に捉えられており、線形モデルの当てはまりがそこそこ良いことが確認できます。事後予測分布が元デー
 6665 タとよく当てはまるということは、データ生成メカニズムとして正規分布や線形モデルを想定したわれわれの
 モデルが悪くないんだ、と解釈していいでしょう。

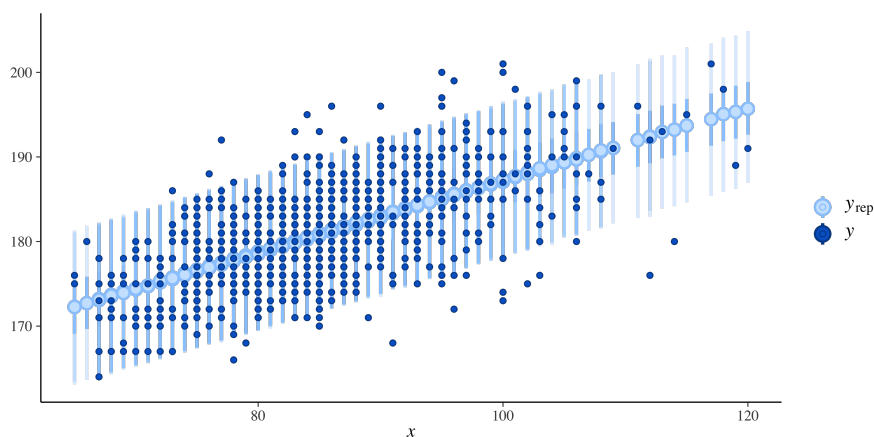


図 25.4 事後予測分布

6666

25.2 データに合わせた確率分布

今回は身長が従属変数、体重が独立変数でした。ここで独立変数を名義尺度水準のもの、すなわち離散変数だと考えると要因計画のようになります。Stan のコードは書き換える必要がありますが、線形モデルとして表現できる場所は同じです*7。

では従属変数を離散変数に置き換えたらどうなるでしょうか。野球のデータの例で考えます。野球は野手と投手に大きく分けることができます。投手は野手と違って毎日登板するということはありませんから、当然のことながら出場試合数は野手と比べ物になりません*8。図 25.5 には、年収 5000 万円を超える一流の選手たちに限定し、試合数で野手か投手かを予測するような分析ができないかと、回帰直線を引いたものを示しています。この直線はどうみても、データにうまく当てはまっているとは言えませんね。

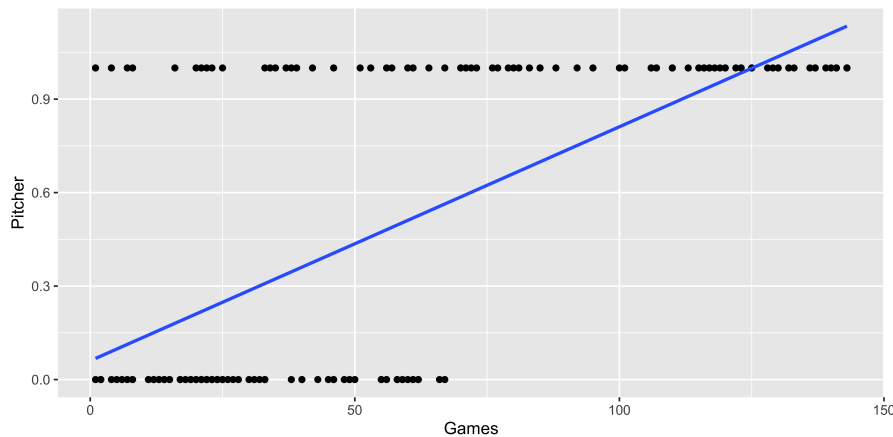


図 25.5 離散変数が目的変数だとおかしなことになる

従属変数が離散的、とくに今回は野手か投手か、という二択の問題になっていますから、データとしては 0/1 しかないわけです。この二種類の数字が正規分布から出てくると考えるのは、流石におかしなことになってしまうでしょう。

そこで前回学んだ離散確率分布の登場です。データが二値しか取らないというのは、コインの表と裏のようなものですからベルヌーイ分布 (Bernoulli Distribution) に従うと考えるのが自然です。ベルヌーイ分布はパラメータ θ で表 (1) が、 $1 - \theta$ で裏 (0) がでるというように、パラメータが 1 つで、このパラメータがそもそも「表が出る確率」を表しているわけです。今回は野手であれば 1、投手であれば 0 というように考えることにしましょう。野手か投手かを表す変数で、野手になる確率 θ が試合数によって変わる (予測できる) ということを考えるのです。この予測の仕方については、「試合数が多いのは野手である確率が高い」という線形的な関係が考えられますから、 $\theta_i = \beta_0 + \beta_1 X_i$ という式を考えることにします。ここで X_i は試合数を表す変数ですね。

ところがここで、注意すべきことが 1 つあります。 θ は野手になる (表が出る) 確率ですから、0 から 1 の範囲の数字しか取れません。プロ野球の試合数は 1 試合から 140 試合ぐらいありますので、この数字をその

*7 実際、R の最尤推定関数の `lm` は、説明変数が `factor` 型であれば関数を書き換えることなく、分散分析のような結果出力をしてくれるのでした。

*8 メジャーリーガーの大谷選手は特殊すぎるケースです。投手は一試合で 100 球程度投げますが、当然徐々に握力が弱くなり、身体的な疲労も翌日、翌々日まで残り続けるものなんでしょう。彼は投げた次の日に DH で打席に立ちますから、とんでもないことです。

6689 まま伸ばしたり (β_1 倍) ずらしたり (β_0 を加える) しても, 0 から 1 の範囲に収めるのは難しそうです。しかも
 6690 上手くその範囲に収めたとしても, そもそも直線なのでどうあってもデータと合致しないことになるでしょう。
 6691 そこで, この数字をロジスティック関数で変換することを考えます。

6692 ロジスティック関数とは次のような式で表されるものです。

$$y = \frac{1}{1 + \exp(-x)}$$

6693 ここで $\exp(x) = e^x$ であり, 正規分布の式の中にも出てくる関数です*9。関数の形は図にするとよくわかりま
 6694 す (図 25.6)*10。

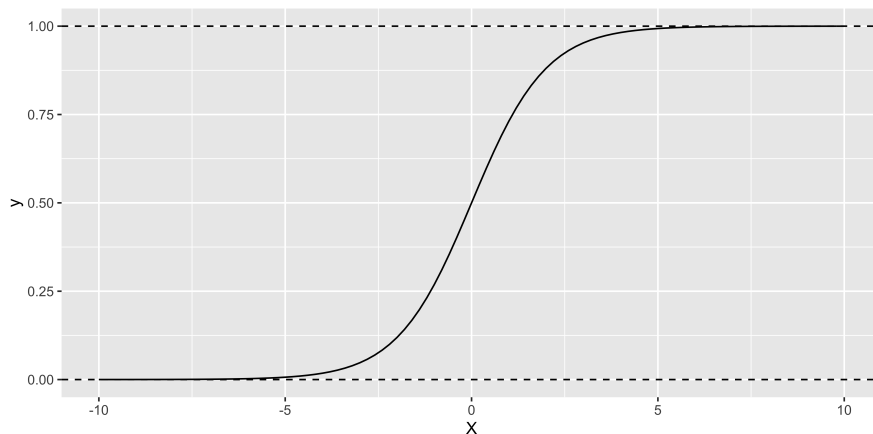


図 25.6 ロジスティック関数

6695 この S 字カーブは, x がどのように変わっても y は必ず 0 から 1 の範囲に入りますので, 今回の目的に
 6696 もってこいです。つまり, 今回「野手か投手か」を判断する結果変数 Y_i は, 次のような分布に従っていると考
 6697 えます。

$$Y_i \sim \text{Bernoulli}(\theta) = \text{Bernoulli}(\text{logistic}(\beta_0 + \beta_1 X_i))$$

6698 これを使ってモデルを書いてみましょう。まずは設計図です (図 25.7)

6699 ここまで書ければ, Stan コードに起こすことも簡単ですね。Stan には `inv_logit` という関数があります
 6700 のでそれを使うと良いでしょう。

code : 25.3 ロジスティック回帰分析のコード

```

6701 1 data{
6702 2   int N;
6703 3   array[N] int<lower=0, upper=1> Y;
6704 4   array[N] real X;
6705 5 }
6706 6
6707 7 parameters{

```

*9 e はネイピア数とも呼ばれる自然対数の底, 約 2.7182... の実数です。

*10 すでに第 4 講で出てきた関数ですが, あの時はいくつかの関数として紹介していました。今回は $\pm\infty$ の範囲の数字を 0 から 1 に変換するものとして改めて紹介しています。なお, 今回は累積正規分布への近似を目的としているわけではないので, 係数 1.7 は外してあります。

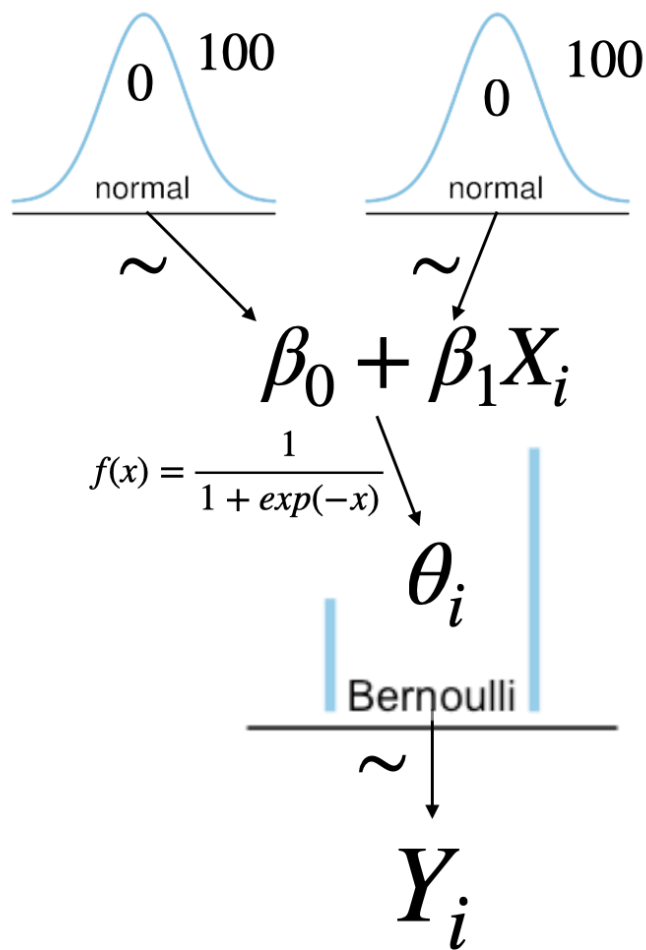


図 25.7 ロジスティック回帰分析の設計図

```

6709 8   real beta0;
6710 9   real beta1;
6711 10  }
6712 11
6713 12  transformed parameters{
6714 13   real theta[N];
6715 14   for(i in 1:N){
6716 15     theta[i] = inv_logit(beta0 + beta1 * X[i]);
6717 16   }
6718 17 }
6719 18
6720 19  model{
6721 20   // model
6722 21   for(i in 1:N){
6723 22     Y[i] ~ bernoulli(theta[i]);
6724 23   }
6725 24   // prior

```

```

6726 25   beta0 ~ normal(0,100);
6727 26   beta1 ~ normal(0,100);
6728 27   }
6729

```

Code25.3 がこのモデルになります。ベルヌーイ分布を使っているのでベルヌーイ回帰と読んでもいいのですが、一般には**ロジスティック回帰 (logistic regression)**と呼ばれています。これに適切なデータセットを与えて (R コード 25.4), 推定した結果が出力 3 になります。

code : 25.4 ロジスティック回帰分析を実行する R コード

```

6733 1  dat2 <- dat %>%
6734 2    dplyr::mutate(Pitcher = if_else(position == "投手", 0, 1)) %>%
6735 3    dplyr::filter(salary > 5000) %>%
6736 4    dplyr::select(Games, Pitcher) %>%
6737 5    na.omit()
6738 6
6739 7  model <- cmdstanr::cmdstan_model("cmdstan/logistic.stan")
6740 8  dataSet <- list(N = NROW(dat2), Y = dat2$Pitcher, X = dat2$Games)
6741 9
6742 10 fit.logistic <- model$sample(
6743 11   data = dataSet,
6744 12   chains = 4,
6745 13   parallel_chains = 4,
6746 14   iter_warmup = 1000,
6747 15   iter_sampling = 5000
6748 16 )
6749 17
6750 18 ## 簡易表示
6751 19 fit.logistic$print(c("beta0", "beta1"))
6752
6753

```

6754 ■コード解説

6755 1-5 行目 先ほど読み込んだデータに Pitcher 変数を作ります。これは投手であれば 0, 野手であれば 1
6756 になる変数です。その後、サラリーが 5000 万円以上のデータにし、変数 Game, Pitcher と必要なも
6757 のだけ抜き出します。最後に na.omit で欠損値のあるデータは除外しました

6758 7 行目 cmdstanr でコンパイルする例です。rstan パッケージを使う場合は適宜変更してください。

6759 8 行目 stan に与えるデータセットをリストで作ります。

6760 10-13 行目 サンプリングするコードです。

6761 19 行目 結果の出力です。出力 3 のように得られます。

rstan の出力 3: ロジスティック回帰分析の結果

variable	mean	median	sd	mad	q5	q95	rhat	ess_bulk	ess_tail
beta0	-2.77	-2.75	0.40	0.41	-3.47	-2.13	1.00	4978	5537
beta1	0.05	0.05	0.01	0.01	0.04	0.07	1.00	4832	5417

6762
6763 出力結果はグラフにした方がわかりやすいでしょう。図 25.8 にあるように、上下の黒い点がデータ点、真ん
6764 中を通る S 字カーブが回帰線になります。回帰線は MAP 推定値で描きましたが、50%, 95% の幅であり
6765 得る可能性の範囲についても薄く書き足しています。これを見ると、ちょうど 50 試合目ぐらいの時にカーブが

6766 Y 軸 0.5, すなわち表裏半々の点を通りますから, 50 試合以上出ている人はより野手と判断されがち, と考
えれば良いでしょう。

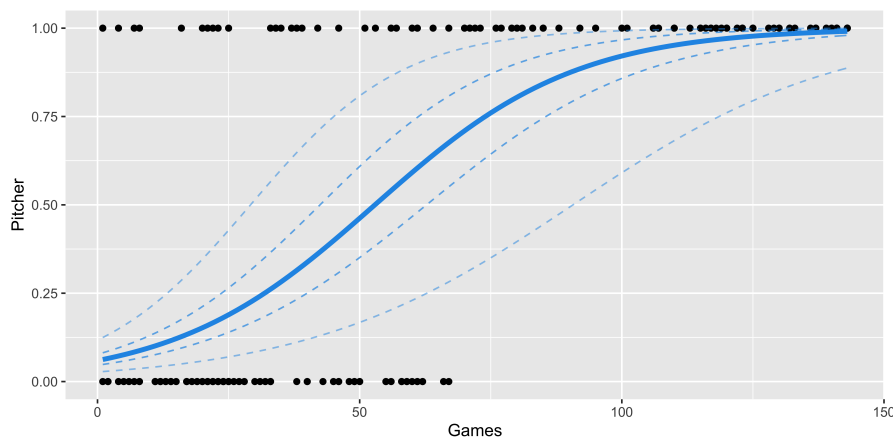


図 25.8 ロジスティック回帰分析の結果

6767

6768 さて, このように従属変数が離散的である場合に, 離散確率分布をつかって線形モデルを当てはめること
6769 ができるのをみてきました。図 25.8 に示されているように, 出力される関数は S 字のカーブ型をしています
6770 が, 予測モデルとしては線形です。今回 $\theta_i = \beta_0 + \beta_1 X_i$ としたように, θ は線形モデルであり, 予測変数 X_i
6771 が線型結合されているので, これも広い意味で線形モデルになるわけです。線型結合したものを, そのまま確
6772 率分布に入れることができなかったので, 変数を変換してモデルの中に組み込むという工夫が必要でした。こ
6773 の変換関数のことをとくにリンク関数 (Link function) と言います。線型結合とパラメータの形をつなげる
6774 (link する) 関数だからです。

6775 このように, リンク関数をかませてさまざまな確率分布に対応させる線形モデルを総称して, 一般化線形モ
6776 デル (Generalized Linear Model; GLM) と言います。この章の初めに出てきた一般線形モデルとは
6777 違いますので注意してくださいね。日本語では「化」があるかないか, 英語では “-ed がつくかどうかだけの
6778 違いですが, 一般線形モデルが正規分布一般だったものに対し, それ以外の確率分布にまで一般化したのが
6779 こちらです。

6780 ちなみに今回はベイズ推定で一般化線形モデルを実践してみましたが, 最尤法による推定も可能です。可
6781 能ではあるのですが, ベイズ推定の方が拡張性が高いこと, 確率モデルとしてまとめて理解しやすいと筆者は
6782 考えています。久保 (2012) は線形モデルの発展段階を推定法とともに, 図 25.9 のように表しています。一
6783 般線形モデルは最小二乗法による推定でも十分だったのですが, 確率モデルとして考えたときに最尤法によ
6784 る推定が必要になってきました。さらに複数の分布を混ぜるような一般化線形混合モデル (Generalized
6785 Mixed Linear Model) にまで発展すると, 最尤法でカバーしきれないところが出てきました (網掛けが
6786 部分的になっていることでそれを表しています)。分布の上に分布がくるような, 幾重にも層が厚くなるモデ
6787 ルのことを階層モデルと言いますが, そうなってくるとはや最尤法では推定できません。ここで「形はわか
6788 らなくても乱数は発生させられる」という MCMC の方法が出てきます。MCMC はすでに説明した通り
6789 (→Pp.178), 乱数による事後分布の近似法です。この方法はベイズ統計学に基づいた方法ですから, ベイズ
6790 推定ともいわれるわけです。ベイズ推定は最尤推定の結果と解釈の仕方が少し異なりますから, そこに注意
6791 する必要がありますが, ほぼ無敵の推定法と言ってもいいでしょう。

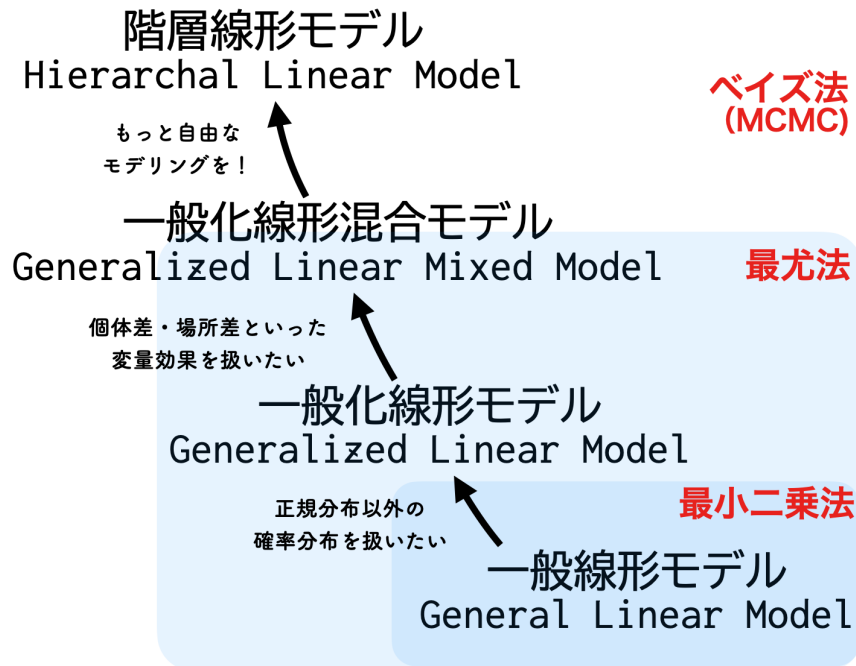


図 25.9 久保 (2012) による線形モデルの発展の図

25.3 リンク関数とパラメータの解釈

さてここでもう一度、先ほどのロジスティック回帰分析の結果をみてみましょう。EAP 推定値でいくと、 $\beta_0 = -2.77, \beta_1 = 0.05$ でしたから、 $\hat{Y}_i = -2.77 + 0.05X_i$ という関係だったことがわかります。

ところで一般線形モデルのうち、この傾きは説明変数が一単位増えた時に、従属変数がどれくらい増えるのかを表す数字でした。これが正規分布に従うモデルであれば、出場試合数が一試合増えれば野手に 0.05 ポイント近づく、という考え方をするのでした。しかし今回は事情が違います。この \hat{Y}_i はこのあと、 $\frac{1}{1 + \exp(-\theta_i)}$ と変換されてしまいますので、 X_i が $X_i + 1$ になった値も変換されて確率の数字に変わっていくはずなんです。ここで少し、数式的な展開を追って何が起きているか考えてみましょう。

まずはロジスティック関数の表記を、少し改めます。同じことなのですが、 $\exp(x) = e^x$ に書き改めて

$$\frac{1}{1 + \exp(-x)} = \frac{1}{1 + e^{-x}} = \frac{1}{1 + \frac{1}{e^x}}$$

とします*11。ここで最後の形に $\frac{e^x}{e^x}$ をかけて、

$$\frac{1}{1 + \frac{1}{e^x}} = \frac{e^x}{e^x + 1}$$

としましょう。さらに、 $\hat{Y} = \beta_0 + \beta_1 X$ で、ベルヌーイ分布のパラメータ θ はこれによって変換されたものを代

*11 指数計算のルールで、マイナス乗は逆数になります。 $a^{-n} = \frac{1}{a^n}$ です。

6803 入ることになりますから、

$$\begin{aligned}\theta &= \frac{e^{\hat{Y}}}{e^{\hat{Y}} + 1} \\ (e^{\hat{Y}} + 1)\theta &= e^{\hat{Y}} \\ \theta e^{\hat{Y}} + \theta &= e^{\hat{Y}} \\ \theta &= e^{\hat{Y}} - \theta e^{\hat{Y}} \\ \theta &= e^{\hat{Y}}(1 - \theta) \\ \frac{\theta}{1 - \theta} &= e^{\hat{Y}} \\ \log\left(\frac{\theta}{1 - \theta}\right) &= \hat{Y}\end{aligned}$$

6804 という関係が成り立ちます*12。つまり、

$$\log\left(\frac{\theta}{1 - \theta}\right) = \beta_0 + \beta_1 X$$

6805 ということですね。

6806 この式の右辺は線形モデルです。左辺は**ロジット関数 (logit function)**と呼ばれるものです。線形モデルはロジット関数を意味し、ベルヌーイ分布のパラメータ θ は確率を表す数字でした。ですから係数の解釈は、線形モデルの方で一単位増えたことは、ロジット関数が一単位増えたことを意味するわけです。ちょっとややこしいですが、線形モデルを確率のパラメータにするには、ロジット関数の逆をすることですから、**ロジスティック関数は逆ロジット関数**とも呼ばれています。そう言えば、Stan の関数も `inv_logit` でしたね*13。

- 6811 • 確率の数字 → 線形モデルにつなげる;ロジット関数; $\log\left(\frac{\theta}{1 - \theta}\right)$
- 6812 • 線形モデル → 確率の範囲に収める;ロジスティック関数あるいは逆ロジット関数; $\frac{1}{1 + \exp(-x)}$

6813 また、線形モデルが表していたのは確率の比ですが、その対数 (log) をとったものでした。対数の中身は確率 θ と確率 $1 - \theta$ の比率ですから、裏が出る確率に比べて何倍表が出やすいか、といったことを表しているのです。これをとくに**オッズ (Odds)**と言います。賭け事の勝率などを表す言葉ですね。オッズの対数がロジットです。ロジットが正の数であれば分子の方が分母より大きい、つまりより表が出やすいことになり、負の数であれば分母の方がより大きい、つまり裏の方が出やすいということになります。

6818 さて話は戻って、線形モデルで説明変数が一単位増えた時のことを考えてみましょう。わかりやすくするために、増える前を A 、増えた後を B とします。

$$\begin{aligned}A &= \frac{1}{1 + \exp(-(\beta_0 + \beta_1 X))} \\ B &= \frac{1}{1 + \exp(-(\beta_0 + \beta_1(X + 1)))}\end{aligned}$$

6821 これはロジット関数を使って元に戻せますから、

$$\log\left(\frac{A}{1 - A}\right) = \beta_0 + \beta_1 X$$

*12 最後の一行は、指数の反対が対数で、 $\log_e^x = a$ は $e^a = x$ をあらわしていますから、指数を取るために両辺に \log をかけています。

*13 `inv` は inverse, 「逆」の意味です。

6822

$$\log\left(\frac{B}{1-B}\right) = \beta_0 + \beta_1(X+1) = \beta_0 + \beta_1 X + \beta_1$$

6823 ここで傾き β_1 のことを考えたいのですが、 β_1 はこの式の関係から、

$$\beta_1 = \log\left(\frac{B}{1-B}\right) - \log\left(\frac{A}{1-A}\right)$$

6824 ということになりますね。ここで対数の差分は割り算ですから*14、

$$\beta_1 = \log\left(\frac{\frac{B}{1-B}}{\frac{A}{1-A}}\right)$$

6825 となります。右辺が非常にややこしい形をしていますので、 \log を外してあげましょう。

$$\exp(\beta_1) = \frac{B}{1-B} / \frac{A}{1-A}$$

6826 右辺は A のオッズに対する B のオッズになりました。オッズの比なのでこれを**オッズ比 (Odds ratio)**と言
6827 います。左辺はロジスティック回帰分析の係数を \exp 関数に入れたものですが、これが説明変数の増加前の
6828 オッズと増加後のオッズの比率です。つまり一単位増加した時にどれくらい生じやすくなるかの比の上昇率
6829 が、 $\exp(\beta_1)$ だというわけです。

6830 ややこしいですね！具体的な数字で考えてみますと、今回 $\beta_1 = 0.05$ だったわけですから、 $\exp(0.05) =$
6831 1.051271 です。つまり、一試合多く出ている選手はそうでない選手に比べて、1.05 倍野手と判断される確率
6832 が高い、ということになります。

6833 25.4 まとめ

6834 今日の話を中心にまとめておきます。

- 6835 • 一般線形モデルは、正規分布を使った線形モデルで、説明変数が連続変数でも離散変数でも同じ形
6836 で表すことができることを示している。
- 6837 • 一般線形モデルの回帰係数をベイズ推定するモデルは、正規分布の位置母数に線形モデルを入れる
6838 だけ。
- 6839 • 従属変数が離散変数になると、確率モデルは離散確率分布を使う必要がある。そのためには、説明変
6840 数の線型結合をリンク関数を使って変形し、適切な形にして推定する必要がある。これを一般化線形
6841 モデルといいます。
- 6842 • 従属変数がベルヌーイ分布に従う例として、ロジット関数をリンク関数としたロジスティック回帰分析の
6843 実例をやってみました。
- 6844 • 一般化線形モデルの場合は説明変数がリンク関数で変換されているから、そのまま解釈するのではな
6845 く、変換の意味を理解しながら注意深く解釈しよう。

6846 一般化線形モデルを導入したことで、要するにリンク関数がわかっていれば色々な確率分布が使えるんだな、
6847 ということが見えてきたのではないのでしょうか。

*14 $\log_a \frac{M}{N} = \log_a M - \log_a N$ です。たとえば $8 = 2^3$, $16 = 2^4$ ですから、 $\log_2 168 = \log_2 2 = 1$ は $\log_2 16 - \log_2 8 = 4 - 3 = 1$ と同じことです。

6848 これは従属変数の形に関わることで、心理学的実践の上でも注意が必要です。たとえば「記憶のト
6849 レーニングをすると思い出せる単語の数が増える」といった研究をする時、従属変数は「思い出した単語の
6850 数」になると思いますが、実験群と統制群とでこれを t 検定する、というのが間違っていることはもうお分かり
6851 ですね。「思い出した数」のようにカウントできる数字はマイナスになることはなく、0 以上の正の整数しか取り
6852 ません。これは正規分布ではなく、**ポアソン分布 (poisson distribution)** に従うことになります。またたと
6853 えば「記憶のトレーニングをすると、50 の単語リストのうち思い出せる単語の割合が増える」と言った研究を
6854 する時、従属変数は今度は「思い出せた率」ということになると思います。これを実験群と統制群で t 検定す
6855 る、というのが間違っているのも明らかですね。「思い出せた率」という比率のデータは 0 から 1 の範囲の実
6856 数で、正規分布ではないからです。これは前回やりました**二項分布 (binomial distribution)** に従うこと
6857 になります。

6858 ポアソン分布に従うような変数の場合は、リンク関数として対数関数 \log が用いられます (逆リンク関数
6859 は指数関数 \exp)。また二項分布に従う変数の場合はリンク関数としてロジット関数 (逆リンク関数はロジス
6860 テック関数) を用います。何気ない「度数」「比率」のようなデータを分析しようという時に、とりあえず数字が
6861 得られているから正規分布でいいや、と考えるのは間違いです。データの性質に応じたモデルを描くように心
6862 がけましょう。

6863 25.5 課題

6864 野球選手のデータの中で、とくに野手に限って考えます。

6865 年俸の高い選手は、成績が良いから高い契約金が得られるのでしょうか。言い換えると、年俸で成績がある
6866 程度予測できるかもしれません。そこで打率を従属変数に、年俸を独立変数にした一般化線形モデルを考え
6867 たいと思います。選手によっては試合への出場回数が違いますので、打席数 N で安打数が K 本である、と
6868 いう情報を用いましょう。これは二項分布に従うモデルということになります。

6869 また、独立変数の年俸 (変数 `saraly`) は個人差が非常に大きいので、標準化してから使うようにしましょ
6870 う。元のデータを加工するコードは提供しますので、二項分布に従う一般線形モデルの Stan コードを書いて
6871 ください。また結果から、年俸が 1 単位上昇するとどういことが言えるでしょうか。併せて報告してください。

6872 考察を導くための計算をする R/Stan コードとともに、回答を提出してください。Rmd ファイルでの提出
6873 が望ましいですが、メモやコメントアウト、Word ファイル、Google ドキュメントなどでの提出も可とします。な
6874 お提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの書き
6875 方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

6876 ■ヒント 打席数と安打数とのデータがありますから、打率は計算したら出るように思えます。しかしここで
6877 は、打率を直接データとして与えるわけではありません。安打数が打席数に伴う二項分布から出現したもの
6878 とし、打率は推定するべきパラメータとして考えます。その打率に線形モデルの構造を入れるようにしてくだ
6879 さい。

6880 打率を安打数と打席数から計算するのは、確率変数の実現値から計算する標本平均を使うことと同じで
6881 す。つまり、標本平均が母平均の推定値になるという**モーメント法**によるアプローチなのです。この方法では、
6882 シーズン初日、一打席だけ出場して 1 本ヒットを打つと、打率 100% だということと同じです。あるいはその初
6883 打席でヒットが出なければ、その人は打率 0% なののでしょうか？本当にそんなことを信じている人はいません
6884 よね。つまり 1 回の実現値だけで全体の推定値とするには偏りがある、ということはわかっているわけです。
6885 モーメント法による推定、あるいは**頻度主義統計学**による考え方では、無限回の試行の平均値、無限に打席
6886 に立ち続けた時に収束していく先の値として確率が定義されています。これに対して**ベイズ統計学**では、**確実**

- 6887 さの指標として確率を用いますから, 初打席の結果がどうであれ「本当はこんなもんじゃないよな」という事前
6888 分布との組み合わせで, 真の打率を推定することになります。

第 26 章

階層線形モデル

前回は一般線形モデルから、一般化線形モデルへ、すなわち正規分布以外の確率分布を扱えるモデルへと展開しました。今回はさらにモデルを展開し、一般化線形混合モデル (Generalized Linear Mixed Model; GLMM) へと展開します。混合 mixed という言葉があるように、複数の確率分布の混ぜ合わせが含まれることが、GLMM の特徴です。

とはいえ、皆さんはすでに第 23 講でこの問題を扱っています。Within 計画・多水準のモデルがそれに当たります。Within 計画のモデルでは、個人差 μ_i に効果 δ_j が加わるという形でモデルを形成しました。ここで個人差 μ_i も正規分布に従うものと考えていました。個人差が従う正規分布、誤差が従う正規分布が混じりあってデータになりますから、これも Mix されたモデルと考えることができるのです。

ただしこの時は、いずれも正規分布のモデルでした。このモデルを正規分布以外のデータに対応させたい、それが今回のモデルになります。

26.1 一般化線形混合モデル

線形モデルを正規分布以外のモデルに対応させるためには、確率分布のパラメータの形にあうようにリンク関数 (link function) で接続してやる必要があるのです。たとえばベルヌーイ分布の場合、線形モデルは次のような確率の比 (オッズ) を表すものと考えます。

$$\log\left(\frac{\theta}{1-\theta}\right) = \beta_0 + \beta_1 X$$

このリンク関数の逆関数を考えてやると、

$$\theta = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 X))}$$

となり、ロジスティック関数と呼ばれるこの変換を通じて $0 \leq \theta \leq 1$ の範囲に入るようにしてから、ベルヌーイ分布に入れてやるのでしたね。

リンク関数と逆リンク関数、この違いがちょっとわかりにくいかもしれません。基本となる線形モデルは $\beta_0 + \beta_1 X$ ですので*1、これが確率分布とセットになる方がリンク関数です。あるいは、線形モデルを変えてしまう方が逆リンク関数だ、と考えてください。

*1 より正確に、より一般的にいうと、説明変数は複数あっても構いませんから、 $\beta_0 + \beta_1 X_1 + \beta_2 X_2 \cdots \beta_n X_n$ と書く方がいいのかもしれませんが、もちろんベクトルを使って $\mathbf{X}\boldsymbol{\beta}$ とすればいいでしょう。この時、係数ベクトル $\boldsymbol{\beta}$ には切片項 β_0 も含まれていることを思い出してください。また、説明変数が複数になった場合は、交互作用 (interaction) 項を考えることもありますが、本書では簡略化のために含めずに考えています。

6911 リンク関数の例をもう 1 つ。**ポアソン分布 (Poisson distribution)** という離散分布があります*2。これ
 6912 は 0 以上の整数を取る**カウント変数**に使われる関数です。たとえば野球選手のホームランの数, たとえば青
 6913 年期における親友の数や交際経験の人数, 記憶実験で思い出せた単語の数などは, 負の数をとることも小数
 6914 点を持った実数になることもありません。こうした数を数えるような変数が従うのがポアソン分布ですが, これ
 6915 のパラメータ λ は正の実数です。ということは, 線形モデルで算出される値を必ず正の数にしなければなりま
 6916 せんから, 逆リンク関数は \exp を使うこととなります*3。

$$\lambda = \exp(\beta_0 + \beta_1 X)$$

6917

$$\log(\lambda) = \beta_0 + \beta_1 X$$

6918 ですから, リンク関数は \log になります。

6919 表 26.1 に確率分布, 逆リンク関数, リンク関数の関係をまとめました。さらに Stan ではそれぞれの確率分
 6920 布に加え, リンク関数とセットになった関数がすでに用意されています。transformed parameters ブロッ
 6921 クで逐一変換しなくても, この関数を使うと model ブロックで直接線形モデルを書き込むことができますの
 6922 で, 便利です。ここで線形モデルで表されるものを μ, θ, λ などさまざまなギリシア文字を使っていますが, こ
 6923 のギリシア文字の違いにとくに意味はありません。一般的に正規分布の**位置母数**には μ が, ベルヌーイ分布
 6924 のそれには θ が, ポアソン分布のそれには λ が使われることが多く, その慣例に習っただけです。いずれも確
 6925 率分布の中心的な位置を表す母数であるということだけが重要です。また正規分布のところを見ると, リンク
 6926 関数も逆リンク関数も同じ形をしていることがわかります。こうした**恒等式**で済んだので複雑なことを考える
 6927 必要がなかったんですね。

表 26.1 確率分布とリンク関数の関係

確率分布	逆リンク関数	リンク関数	Stan の専用関数
正規分布	$\mu = \beta_0 + \beta_1 X$	$\beta_0 + \beta_1 X = \mu$	normal
ベルヌーイ分布	$\theta = \frac{1}{1 + \exp(-(\beta_0 + \beta_1 X))}$	$\beta_0 + \beta_1 X = \log\left(\frac{\theta}{1 - \theta}\right)$	bernoulli_logit
ポアソン分布	$\lambda = \exp(\beta_0 + \beta_1 X)$	$\beta_0 + \beta_1 X = \log(\lambda)$	poisson_log

6928 では一般化線形モデルに「個人差」を混合するモデルを考えてみましょう。ここでは前の章でも扱った野球
 6929 のデータを例に説明します。図 26.1 には Swallows で 2011 年から 2020 年の間, 8 年以上在籍した投手の
 6930 年俸 (単位は千万円) と勝利数の関係を表しています。年俸が高いと勝ち星が増える, という関係にあると考
 6931 えてモデルを組みます。

6932 線形モデルですから, 勝ち星を Y_i , 年俸を X_i として $\hat{Y}_i = \beta_0 + \beta_1 X_i$ という関係を考えます。とはいえ,
 6933 選手の年俸というのはそれまでの業績によるどころも大きく, そもそも個人差があり得るところです。そこで個
 6934 人差 μ_i を考えて, $\hat{Y}_i = \beta_0 + \beta_1 X_i + \mu_i$ とした線形モデルを考えましょう。また, 勝利数は 0 以上の整数し
 6935 か取りませんから, これはポアソン分布に従うと考えられます。そこでポアソン分布に合わせるために, 逆リン
 6936 ク関数を使って $\lambda = \exp(Y_i) = \exp(\beta_0 + \beta_1 X_i + \mu_i)$ とした上で, $Y_i \sim \text{poisson}(\lambda_i)$ という確率モデルに
 6937 入れることにしましょう。このモデルの設計図は図 26.2 のようになります。

*2 後学のためにきちんと定義式を書いておくと, 変数 X が k の値を取る確率を $P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}$ とした分布です。ここで e はいつものネイピア数であり, $k!$ の ! は階乗を表す記号です。

*3 この関数はよく出てきますから周知のことと思いますが, $\exp(x) = e^x$ でこの e もいつものネイピア数です。この関数で x が負の数になっても, $x^{-1} = \frac{1}{x}$ という指数計算の決まりがありますから, 結果が負になることはありません。

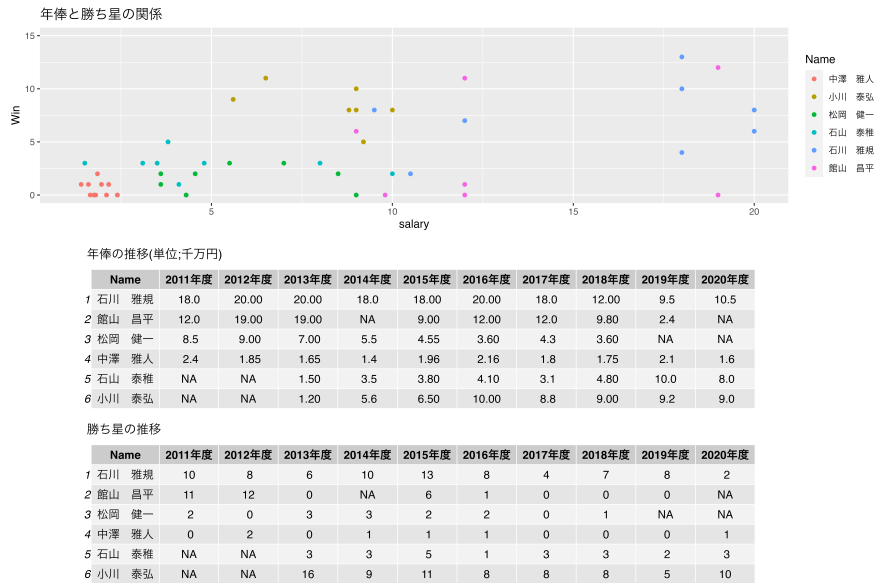


図 26.1 一流選手の年俸と勝利数

6938 それではこれをもとに、GLMM をやってみましょう。実行するための Stan のコードはコード 26.1 のよう
6939 になります。

code : 26.1 個人差を含んだポアソン回帰モデル

```

6940
6941 1 data{
6942 2   int L;
6943 3   int N;
6944 4   array[L] real X;
6945 5   array[L] int Y;
6946 6   array[L] int index;
6947 7 }
6948 8
6949 9 parameters{
6950 10  real beta0;
6951 11  real beta1;
6952 12  array[N] real mu;
6953 13 }
6954 14
6955 15 model{
6956 16  // model
6957 17  for(l in 1:L){
6958 18    Y[l] ~ poisson_log(beta0 + (beta1 * X[l]) + mu[index[l]]);
6959 19  }
6960 20  // prior
6961 21  beta0 ~ normal(0,10);
6962 22  beta1 ~ normal(0,10);
6963 23  mu ~ normal(0,10);
6964 24 }
6965

```

6966 ■コード解説

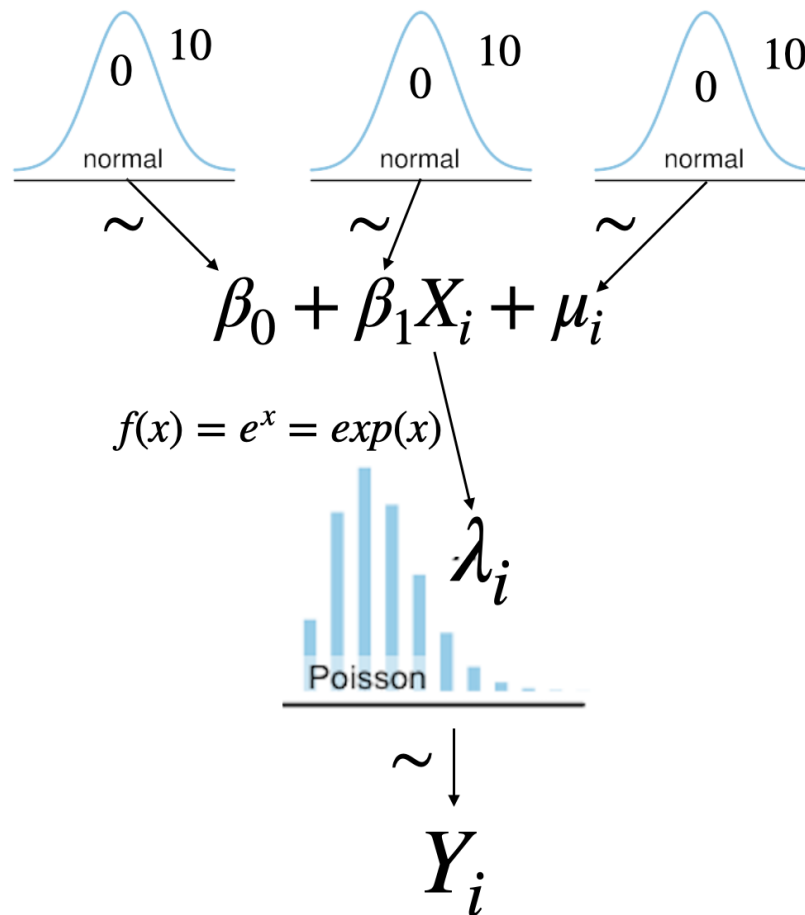


図 26.2 ポアソン分布を使った GLMM モデルの設計図

6967 data ブロック データ長 L と扱う選手数 N , 独立変数 X , 従属変数 Y , 個人差インデックスを宣言して
6968 います。ポアソン分布に従う確率変数ですので従属変数は `int` 型にしています。

6969 parameters ブロック 線形モデルの係数と個人差 μ_i を人数分用意しました。

6970 model ブロック モデル尤度のところは `poisson_log` を使っているの、線形のコードをそのまま書き込
6971 むことができます。事前分布は正規分布にしました。

6972 これを R で呼び出すコードは Code::26.2 の通りです。少しテクニカルなコードになっていますので自分で
6973 書けなくても構いませんが^{*4}, 各ステップで何をしているか読めるようになっておくといいでしょう^{*5}。

code : 26.2 ポアソン回帰のコード

```
6974
6975 1 baseball <- read_csv("baseballDecade.csv")
6976 2
6977 3 dat <- baseball %>%
6978 4   filter(position == "投手") %>%
6979 5   filter(team == "Swallows") %>%
```

*4 コードの書き方は単一の正解というものはありません。ここに書いてあるのは著者流のやり方だというだけであって、同じ結果が出るのであれば別のコードの書き方であっても構いません。今回のこのコードはシラバスのサイトで配布しています。

*5 パイプ演算子 `%>%` を使って変換プロセスをつなげて言っていますが、その途中でデータがどのように加工されているか知りたい場合、`%>% print %>%` と `print` 関数を挟むといいでしょう。ステップごとのデータの変形プロセスがよくわかると思います。

```

6980 6   group_by(Name) %>%
6981 7   nest() %>%
6982 8   mutate(
6983 9     n = purrr::map_dbl(data, ~ NROW(.)),
6984 10    FLG = purrr::map_lgl(data, ~ anyNA(.$Win))
6985 11   ) %>%
6986 12   filter(n > 7) %>%
6987 13   filter(!FLG) %>%
6988 14   unnest(data) %>%
6989 15   select(Year, Name, salary, Win)
6990 16
6991 17   dat.tmp <- dat %>%
6992 18   mutate(salary = salary / 1000) %>%
6993 19   mutate(ID = as.factor(Name)) %>%
6994 20   mutate(ID = as.numeric(ID))
6995 21
6996 22   dataSet <- list(
6997 23     L = NROW(dat.tmp),
6998 24     X = dat.tmp$salary,
6999 25     Y = dat.tmp$Win,
7000 26     index = dat.tmp$ID
7001 27   )
7002 28
7003 29   model_pois <- cmdstanr::cmdstan_model("glm_poisson.stan")
7004 30   fit <- model_pois$sample(
7005 31     data = dataSet,
7006 32     chains = 4,
7007 33     parallel_chains = 4,
7008 34     iter_warmup = 1000,
7009 35     iter_sampling = 3000
7010 36   )
7011

```

7012 ■コード解説

7013 1 行目 ファイルからデータを読み込みます。このデータファイル `baseballDecade.csv` には野手、投手の
7014 10 年分のデータが入っています。

7015 3-15 行目 最初のブロックとして、データの加工・変改を行い、それを `dat` オブジェクトに格納します。

7016 4 行目 フィルターをかけて、投手だけのデータに絞り込みます。

7017 5 行目 フィルターをかけて、`Swallows` の選手に限りませ

7018 6 行目 選手名でグルーピングします。

7019 7 行目 `nest` 関数でグループ化変数に沿ってデータを畳み込みます。引数をとくに指定しなければ、
7020 変数はすべて `data` という変数名に含まれます。この段階でデータフレームは出力 26.1 のように
7021 なっています。`nest` 関数のイメージを図 26.3 に示しました。

7022 8 行目 畳み込まれたデータセットを使って、そのデータの行数 (何年分のデータがあるか^{*6})、勝利

^{*6} ここで `data` とあるのはまとめられたミニデータフレーム全体を引数にしています。それに対して `NROW` という関数をあてがうことで行数を数えています。`NROW` 関数の引数は、ですが、これは `map` 関数の第一引数全体を再度参照するときの略記号です。`purrr::map` というのは `purrr` パッケージの `map` 関数で、同じ関数操作を繰り返し行うときに使います。

7023 数情報の欠損値の有無を作る変数^{*7}を作っています。

7024 9 行目 フィルターをかけて、データの行数 n が 7 より大きいもの、つまり 8 年以上登板している一

7025 流の選手だけに絞ったのです。

7026 10 行目 フィルターをかけて、欠損値がないデータだけに絞っています。

7027 11 行目 畳み込みを解除します。

7028 12 行目 年度、選手名、年俸、勝利数だけ変数をセレクトして取り出します。

7029 17-20 行目 このブロックでは Stan で使うためにさらに一時的なデータ加工をし、`dat.tmp` オブジェクトに

7030 格納しています。

7031 18 行目 年俸のデータは数字が大きいので、1000 万円単位にします。分析の際、数字が大きすぎる

7032 と計算がオーバーフローしてしまうことがありますから、適当な単位にしておくことは実践上の有

7033 効なテクニックの 1 つです。

7034 19 行目 名前の変数は文字型ですが、そのまま Stan に渡すことはできないので、Factor 型にした

7035 ID という変数に作り替えます。

7036 20 行目 先ほど作った ID は Factor 型ですが、Stan では質的な違いを表す数字だけあれば良い

7037 ので、数字型に変形します。

7038 22-27 行目 Stan に与えるためのデータセットにするブロックです。

7039 29-36 行目 `cmdstanr` でサンプリングするためのコードです。`rstan` の場合は `sampling` 関数を使うこと

7040 になります。

R の出力 26.1: データをグループ化変数でネストする

```
# A tibble: 295 × 2
# Groups:   Name [295]
  Name      data
  <chr>    <list>
1 永川 勝浩 <tibble [3 × 16]>
2 前田 健太 <tibble [5 × 16]>
3 シュルツ <tibble [1 × 16]>
4 大竹 寛   <tibble [8 × 16]>
5 横山 竜士 <tibble [2 × 16]>
6 ジオ     <tibble [2 × 16]>
7 バリントン <tibble [5 × 16]>
8 藤川 球児 <tibble [7 × 16]>
9 久保 康友 <tibble [7 × 16]>
10 小林宏   <tibble [1 × 16]>
# ... with 285 more rows
```

7041

7042 さて、このコードを実行し、結果からモデルを図にしたのが図 26.4 です。図中の赤い転線は、 $+\mu_i$ をお

7043 ずに全体を混ぜ合わせたときのポアソン回帰です。個人差を表す項目を追加することで、より個々のデータに

7044 沿った回帰プロットが実践できていることがわかると思います。

7045 今回は個人差 μ_i だけを考えましたが、たとえば野球の弾の規定が変わってピッチャーが有利になる年とそ

7046 うでない年があったとか、登板する球場のサイズによって有利不利がある、といったさまざまな個別事情に対

7047 し、それを表す変数の項を付加することで、よりデータの詳細を細かくモデルに組み込んでいくことができま

^{*7} `purrr::map` の使い方は先ほどと同じで、`.$Win` は `data$Win` と同じ意味です。`anyNA` は NA があるかどうかを TRUE/FALSE で返す関数です。

Year	Name	V1	V2
2011年度	A	19	1
2011年度	B	31	10
2011年度	C	19	0
2011年度	D	6	1
2012年度	A	2	0
2012年度	B	18	3
2012年度	C	30	13
2012年度	D	56	3
2013年度	A	20	8
2013年度	B	42	1
2013年度	C	32	12
2013年度	D	13	5
2014年度	A	15	3
2014年度	B	18	1
2014年度	C	19	2
2014年度	D	23	4

Name	Data															
A	<table border="1"> <thead> <tr><th>Year</th><th>V1</th><th>V2</th></tr> </thead> <tbody> <tr><td>2011年度</td><td>19</td><td>1</td></tr> <tr><td>2012年度</td><td>2</td><td>0</td></tr> <tr><td>2013年度</td><td>20</td><td>8</td></tr> <tr><td>2014年度</td><td>15</td><td>3</td></tr> </tbody> </table>	Year	V1	V2	2011年度	19	1	2012年度	2	0	2013年度	20	8	2014年度	15	3
	Year	V1	V2													
	2011年度	19	1													
	2012年度	2	0													
2013年度	20	8														
2014年度	15	3														
B	<table border="1"> <thead> <tr><th>Year</th><th>V1</th><th>V2</th></tr> </thead> <tbody> <tr><td>2011年度</td><td>31</td><td>10</td></tr> <tr><td>2012年度</td><td>18</td><td>3</td></tr> <tr><td>2013年度</td><td>42</td><td>1</td></tr> <tr><td>2014年度</td><td>18</td><td>1</td></tr> </tbody> </table>	Year	V1	V2	2011年度	31	10	2012年度	18	3	2013年度	42	1	2014年度	18	1
	Year	V1	V2													
	2011年度	31	10													
	2012年度	18	3													
2013年度	42	1														
2014年度	18	1														
C	<table border="1"> <thead> <tr><th>Year</th><th>V1</th><th>V2</th></tr> </thead> <tbody> <tr><td>2011年度</td><td>19</td><td>0</td></tr> <tr><td>2012年度</td><td>30</td><td>13</td></tr> <tr><td>2013年度</td><td>32</td><td>12</td></tr> <tr><td>2014年度</td><td>19</td><td>2</td></tr> </tbody> </table>	Year	V1	V2	2011年度	19	0	2012年度	30	13	2013年度	32	12	2014年度	19	2
	Year	V1	V2													
	2011年度	19	0													
	2012年度	30	13													
2013年度	32	12														
2014年度	19	2														
D	<table border="1"> <thead> <tr><th>Year</th><th>V1</th><th>V2</th></tr> </thead> <tbody> <tr><td>2011年度</td><td>6</td><td>1</td></tr> <tr><td>2012年度</td><td>56</td><td>3</td></tr> <tr><td>2013年度</td><td>13</td><td>5</td></tr> <tr><td>2014年度</td><td>23</td><td>4</td></tr> </tbody> </table>	Year	V1	V2	2011年度	6	1	2012年度	56	3	2013年度	13	5	2014年度	23	4
	Year	V1	V2													
	2011年度	6	1													
	2012年度	56	3													
2013年度	13	5														
2014年度	23	4														

図 26.3 データの畳み込み関数の挙動。左側がもとのデータで、グループ化変数を使って畳み込むと、データフレームの中にグループごとのミニ・データフレームが畳み込まれる。

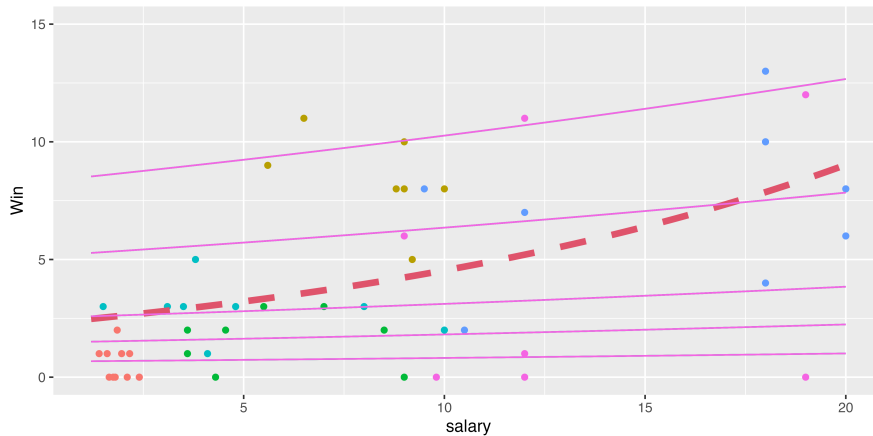


図 26.4 ポアソン分布を使った GLMM の結果

7048 す。線形モデルは線形であるという制限から逃げることはできませんが、確率分布や個別の事情を組み込む
 7049 ことで、より上手く個々のデータの特徴を表現できるようになることが、イメージできたのではないのでしょうか。

7050 26.2 ネストされたデータ

7051 さて先ほど、nest 関数を使ってデータをネストするという例を見ました。
 7052 ネストとは「入れ子にする」という意味で、マトリョーシカのようにデータフレームの中のデータフレーム、デー
 7053 タフレームの中のデータフレームの中のデータフレーム、というように、同じ形のものをに入れていくイメージで
 7054 す。このネストされたデータは、別名階層性を持ったデータだ、ということが出来ます。
 7055 たとえば今回の野球データの場合、選手はどこかのチームの一員ですから、チームという上位階層の中に
 7056 選手がいることになります。野球チームは 12 球団ありますが、それぞれセ・リーグ、パ・リーグに分かれて戦い

7057 ます*8。つまりリーグという上位階層があるわけです。

7058 逆に、チーム傘下の選手も、年間 140 試合ぐらいありますから、試合ごとに調子が良かったり悪かったりす
7059 るでしょう。試合の中でも、野球は 9 回のターン制バトルですから、毎回の活躍というのがあるかもしれませ
7060 ン。すなわち、ある選手の中に試合や回がネストされていると考えることもできます。私たちはデータのある側
7061 面で区切ってみていることになります。

7062 こうしたデータの階層性は、何も野球選手のデータだけではありません。たとえば調査研究を行う場合、あ
7063 る小学校でデータをとったというときに、児童は学年でネストされ、学校でネストされ、地区でネストされ、市
7064 区町村でネストされ、都道府県でネストされている、ということになります。もちろん 1 つの小学校だけでデー
7065 タを取るのであれば比較検証はできませんが、ある程度規模が大きくなってくると何らかの階層的な情報を
7066 持っていると考えerべきでしょう。

7067 またあるいは、心理学の実験のような少数数からしかデータを得ないという場合であっても、条件 A で n
7068 回試行、条件 B で m 回試行する、といった Within デザインの場合は**反復測定 (repeated measure)**
7069 であり、これも被験者ごとに各データがネストされている、と考えることができるわけです。Within データの
7070 場合は個人差を考え、GLMM で個人差以外の要因もモデルに組み込むことができるようになったわけです
7071 から、データの階層性について考えないわけにはいきません。

7072 何がそんなに問題なのでしょう。図 26.5 をみてください。左右にあるのは同じデータの分布ですが、左側
7073 は階層の区別をしない場合で、右側はその区別をした場合です。線形回帰の直線を引いておきましたが、こ
7074 れを見ると明らかなように階層性を考えなければ、私たちは正の相関があるデータだと判断してしまうでしょ
7075 う。しかし右側、群ごとに区別してモデルを書いてみると、回帰直線はすべて右下がり、つまり負の相関がある
7076 データだったことがわかります。階層性や群ごとの特徴を無視して分析してしまうと、結果が真逆になってしま
う可能性を孕んでいるわけです。

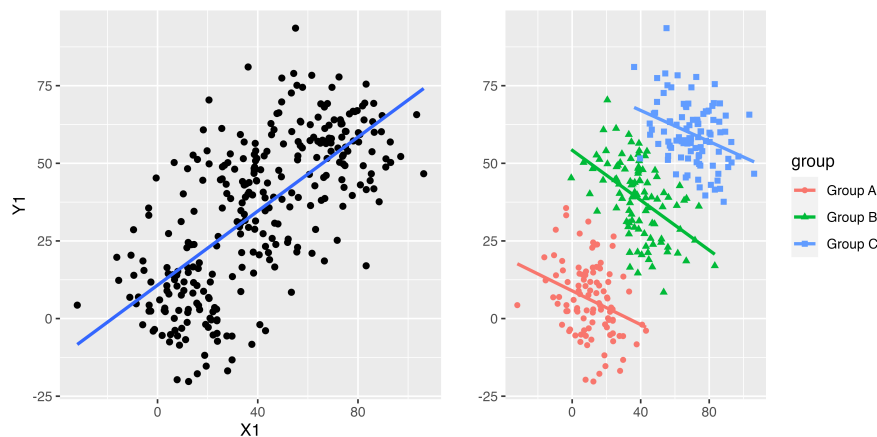


図 26.5 階層データの分析が必要な理由

7077

7078 **データは図にすることの重要性**は何度指摘してもしすぎることはないと思いますが、ことほど左様にデータ
7079 の持っている性質はなるべく丁寧に見定めなければなりません。そして階層に区分して分析する必要がある

*8 セントラル・リーグ Central League とパシフィック・リーグ Pacific League の略であり、セ・リーグには巨人 (Giants)、阪神 (Tigers)、広島東洋カープ (Carp)、中日 (Dragons)、ヤクルト (Swallows)、DeNA (Baystars) が含まれます。パ・リーグにはソフトバンク (Hawks)、西武 (Lions)、日本ハム (Fighters)、ロッテ (Marines)、オリックス (Buffaloes)、楽天 (Eagles) からなります。2021 年現在のチーム名であり、かつては大洋ホエールズとか、阪急ブレーブスとか、南海ホークスなどのチームがありました。

7080 のであれば、モデルもそのようにレベルアップさせなければなりません。そこで出てくるのが**階層線形モデル**
7081 (**Hierarichal Linear Model**) です。

7082 26.3 階層線形モデル

7083 **階層線形モデル**は、ネストされたデータに対して群ごとの特徴を記述できるよう、線形モデルを拡張したも
7084 のです。

7085 野球の例で考えてみましょう。階層化せずにデータを眺めて、「グラウンドにはゼニが落ちてるんだ」とばかり
7086 りに、たくさん年俸をもらっているバッターはホームランを打つし、ピッチャーは勝利数を稼ぐ、という傾向を分
7087 析したいとします。年俸が説明変数で、成績が従属変数になるモデルですね^{*9}。しかし、ソフトバンクや読売巨
7088 人のように、資本力のある親会社のもとでの球団であればそういう傾向もあるかもしれませんが、広島カーブ
7089 のような市民球団^{*10}ではそこまで大盤振る舞いできるほどではないよ、といった球団ごとの差異があるかも
7090 しません。

7091 年俸が成績に影響する、というのが全体的な傾向としてあり、影響の程度は群ごとの特徴が多少入り込
7092 む、というこの状況をモデルで表現することを考えましょう。変数の影響力の大きさは、 β_0, β_1 などの回帰係
7093 数で表現されます。全体的な傾向はこれで表現できるとして、個別の効果はこの係数が群 j ごとに違ってい
7094 る、ということになります。この群ごとの違いを、群平均からの差異で表現しましょう。

7095 球団 j に属するプレイヤー i の成績を Y_{ij} と表現したとします。データ全体としては、 $\hat{Y}_{ij} = \beta_{0j} + \beta_{1j}X_{ij}$
7096 なのですが、ここでたとえば切片 β_{0j} に群 j ごとの効果があるとすれば、 $\beta_{0j} = \gamma_{00} + u_{0j}$ のように、切片
7097 の平均 γ_{00} と、切片に加わる群の効果 u_{0j} という、**係数を説明する線形モデル**を想定することと同じです。
7098 切片だけでなく傾き β_{1j} の方も、 $\beta_{1j} = \gamma_{10} + u_{1j}$ と線形モデルで考えましょう。個人差 u_{1j} は平均 0 の正
7099 規分布に従うと考えます。ここで全体的な傾向、すなわち個人差の平均である γ_{00}, γ_{10} のことを**固定効果**
7100 (**Fixed Effect**) といい、個別の効果 u_{0j}, u_{1j} のことを**変量効果 (Random Effect)** と呼ぶことがあり
7101 ます。

7102 式を整理するとこうなります。

$$\begin{aligned} \hat{Y}_{ij} &= \beta_{0j} + \beta_{1j}X_{ij} \\ \beta_{0j} &= \gamma_{00} + u_{0j} \\ \beta_{1j} &= \gamma_{10} + u_{1j} \end{aligned} \tag{26.1}$$

$$\begin{pmatrix} u_{0j} \\ u_{1j} \end{pmatrix} \sim MN \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \tau_{00} & \tau_{01} \\ \tau_{10} & \tau_{11} \end{pmatrix} \right)$$

7103 あらまあ。何だか訳がわからなくなりそうですね。いやもちろん、じっくり考えればお分かりいただけると信
7104 じていますが、直感的にわかりにくいかもしれません。そこで同じことを違う角度から説明してみたいと思
7105 います。

7106 階層性のあるデータとはいえ、たかだか 12 球団ですから、12 回ぐらいの繰り返し分析は頑張ればできま
7107 すよね。実際それをやったのが図 26.6 にあります。えっへん。この結果を解釈するとき、球団ごとにバラバ
7108 ラの分析をしていることと同じですから、12 球団ごとの切片 β_{0j} と傾き β_{1j} が出てきます。ちょっと面倒です
7109 が、データがそういうものなだから仕方ありません。

^{*9} ちょっと待って、逆じゃないか、と思った人もいるかもしれません。つまり成績がいいから年俸が高くなるのでは、ということですよ
ね。その場合は、ある年 X の年俸が翌年 $X + 1$ 年の成績を予測する、という形にデータを変えなければなりません。今回は年
度ごとの年俸と成績からなるデータセットになっていますから、年俸が成績に影響する、という仮説になっています。

^{*10} 広島カーブは戦後復興の名目で公的資金が投入されて作られた球団であり、その意味で市民の資本から作られた球団と言
えるでしょう。現在も、特定の親会社はなく複数の企業が連携して運営するスタイルをとっています。

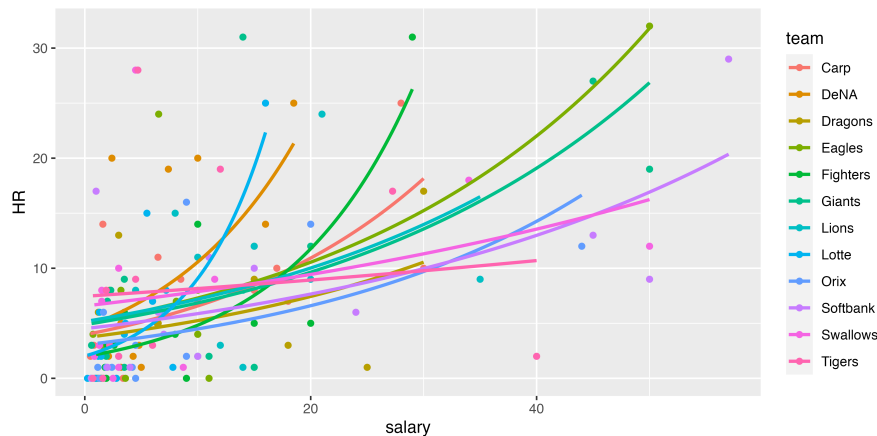


図 26.6 チームごとにバラバラの回帰分析を実施する。できなくはないけど、要するにどういふこと、と言いたくもなる。

7110 でもこれが、12 球団ではなく、100 球団、1000 球団あったらどうしますか。1000 個の切片と傾きをず
 7111 ら一つと数表に並べられても、目が痛くなりますね。辛いだけです。せっかく大量のデータを手に入れている
 7112 のに、分析結果を細かくみてくれなかったら悲しい。かといって「で、結局どういふことなの」と言われた時に、
 7113 「えー、第 1 球団の切片と傾きは…で、第 2 球団の切片と傾きは…で、第 998 球団の切片と傾きは…」と
 7114 読み上げることもないでしょう。ではどうするか。傾きと切片の平均と SD で報告するのではないのでしょうか。
 7115 傾きの平均はこれぐらいです、球団による違いは ± これぐらいです、たとえば、1000 球団の情報を要約し
 7116 て言えることになります。それが $u_{0j} \sim N(0, \tau_0)$ や $u_{1j} \sim N(0, \tau_1)$ が言っていることなんですね。散らばり
 7117 の平均は γ_{00}, γ_{10} であり、それが τ_0, τ_1 ぐらい散らばりますよ、ということです。

7118 これは見方を変えると、階層線形モデルを使う場合は、階層性をまとめる必要があるかどうかを考えてから
 7119 でも構わないということです。12 球団しかなくて、それぐらいならじっくりみましようというのであれば、無理
 7120 に階層モデルにする必要はありません。また 100 球団、1000 球団ある世界であっても、ほとんど球団ごとの
 7121 違いがないようであれば、全部まとめて分析してしまっても構いません。階層線形モデルで表現する場合は、
 7122 ネストされているデータに一定の傾向はあるものの、考慮する必要がある程度に散らばりがある場合、という
 7123 ことになります*11。

7124 またさらに別の角度から考えてみましょう。今度はこの線形モデルの、設計図を書いてみるのです。図 26.7
 7125 にそれを書いてみました。

7126 図 26.7 では、話を簡単にするために係数の共分散を想定しない ($\tau_{10} = \tau_{01} = 0$)、1 次元正規分布から
 7127 係数が来ているモデルにしています。これをみると、

- 7128 1. データはポアソン分布からきている; $Y_{ij} = \text{poisson}(\lambda_{ij})$
- 7129 2. ポアソン分布は線形モデルから逆リンク関数で接続されてきている; $\lambda_{ij} = \exp(\beta_{0j} + \beta_{1j}X_{ij})$
- 7130 3. 線形モデルの切片が別の分布からきている; $\beta_{0j} \sim N(\gamma_{00}, \tau_{00})$
- 7131 4. 線形モデルの傾きも別の分布からきている; $\beta_{1j} \sim N(\gamma_{10}, \tau_{11})$
- 7132 5. γ_{00} がどういふ数字かわからないので、確率分布 (事前分布) でそのわからなさを表現; $\gamma_{00} \sim$
 7133 $N(0, 10)$

*11 これも分散分析的発想で、群内の類似性と群間の類似性の比率から、十分に大きな群間変動があるようなら階層モデルにするべき、と判断する基準があります。この比率のことを級内相関 (Intra-Class Coefficients) と言います。詳しくは清水・荘島 (2017) を参考にしてください。

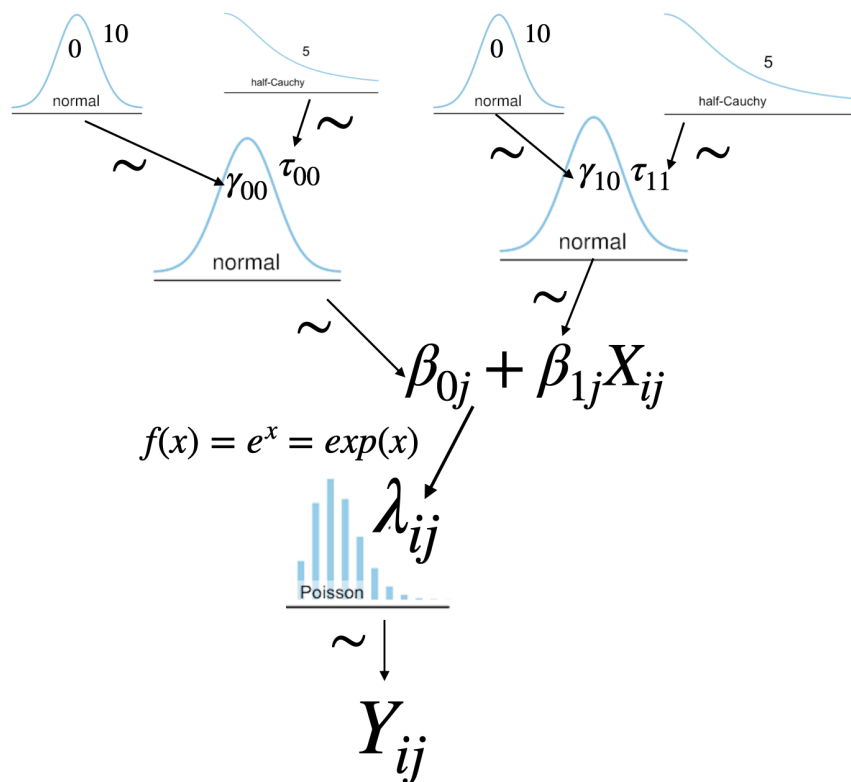


図 26.7 チームごとにバラバラの回帰分析を実施する

- 7134 6. γ_{10} がどういう数字かわからないので、確率分布でそのわからなさを表現; $\gamma_{10} \sim N(0, 10)$
 7135 7. τ_{00} がどういう数字かわからないので、確率分布 (事前分布) でそのわからなさを表現。SD なので
 7136 cauchy 分布にしてみた; $\tau_{00} \sim \text{cauchy}(0, 5)$
 7137 8. τ_{11} がどういう数字かわからないので、確率分布 (事前分布) でそのわからなさを表現。SD なので
 7138 cauchy 分布にしてみた; $\tau_{11} \sim \text{cauchy}(0, 5)$

7139 という順番で読み解くことができます。データに分布を仮定する、これがベイズ推定のスタートでしたが、仮定
 7140 された分布のパラメータに、線形モデルという数学的構造が入っています。このパラメータに数式を入れるこ
 7141 とをモデリング (modeling) というわけです。そしてそのモデルに含まれている未知数 (パラメータ) に、さ
 7142 らに別の分布が入っています。分布に分布が入っているから、混合 mix しているわけです。しかもこの混合
 7143 加減は、上に上にとレベルが上がっていきますので、階層的なモデリングになっています。パラメータのパラ
 7144 メータは、ハイパーパラメータ (Hyper parameter) と呼んだりします。かっこいい。

7145 さて、これをみればわかるように、階層線形モデルは、12 の球団にある違いに正規分布というカバーを
 7146 かけ、平均とそこからの散らばりという形で表現することでもあります。このように、係数たちに正規分布と
 7147 というカバーを上からかけると、12 球団バラバラで推定した時よりも、その係数が少し小さくなることがあり
 7148 ます。上位階層の正規分布の平均値に、少し引っ張られて推定されてしまうわけです。これを係数の縮小
 7149 (shrinkage) と言います。そうした問題はありますが、事前分布や階層性のおかげで、下位レベルの推定値
 7150 が安定して算出できることにもなります。

7151 さあ設計図が書けたら、もうコードに落とすことはできますね。ポアソン分布を使った階層線形モデルのコー
 7152 ド例を Code::26.3 に挙げておきます。

code : 26.3 階層ポアソン回帰モデル

```

7153 1 data{
7154 2   int L;
7155 3   int G;
7156 4   array[L] int Gindex;
7157 5   array[L] real X;
7158 6   array[L] int Y;
7159 7 }
7160 8
7161 9 parameters{
7162 10  array[G] real beta0;
7163 11  array[G] real beta1;
7164 12  real gamma0;
7165 13  real gamma1;
7166 14  real<lower=0> tau0;
7167 15  real<lower=0> tau1;
7168 16 }
7169 17
7170 18 transformed parameters{
7171 19  array[L] real<lower=0> lambda;
7172 20  for(l in 1:L){
7173 21    lambda[l] = exp(beta0[Gindex[l]] + (beta1[Gindex[l]] * X[l]));
7174 22  }
7175 23 }
7176 24
7177 25 model{
7178 26  // model
7179 27  for(l in 1:L){
7180 28    Y[l] ~ poisson(lambda[l]);
7181 29  }
7182 30
7183 31  for(g in 1:G){
7184 32    beta0[g] ~ normal(gamma0,tau0);
7185 33    beta1[g] ~ normal(gamma1,tau1);
7186 34  }
7187 35
7188 36  //prior
7189 37  gamma0 ~ normal(0,10);
7190 38  gamma1 ~ normal(0,10);
7191 39  tau0 ~ cauchy(0,5);
7192 40  tau1 ~ cauchy(0,5);
7193 41 }
7194
7195

```

7196 ■コード解説

7197 **data ブロック** データは整然データの形で渡します。データの長さ L, 説明変数 X, 被説明変数 Y の他に,
7198 群の数 G, どの群に属するのかを表すインデックス変数 G[L] を用意しています。

7199 **parameters ブロック** 群の数だけ係数 β_0, β_1 が必要です。加えて, これら切片と傾きの平均値 γ と散らば
7200 り τ を用意しています。

7201 **transformed parameters ブロック** poisson_log を使ってもいいのですが, 話をわかりやすくするために

7202 いったん λ_{ij} を作っています。係数 $\text{beta0}[g]$ のグループを表す g のところは、 L 行目のデータの所属先を表すインデックス、 $\text{Gindex}[1]$ で代入しています。 beta1 についても同様です。

7203

7204 **model ブロック** データの各行は、行ごとの線形モデルで推定されますのでスッキリしたものです。

7205 $\text{beta0}[g]$ や $\text{beta1}[g]$ はハイパーパラメータを持つ上位の正規分布にカバーされています。

7206 パラメータの事前分布は設計図通りです。

7207 これを使って、年俵でホームランの数を予測する階層線形モデルを推定させましょう。結果を出力 9 に示してあります。

7208

MCMC の結果 9: ハイパーパラメータの推定値

```
# A tibble: 4 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 gamma0  1.446    1.448    1.449    0.120    1.200    1.680
2 gamma1  0.046    0.046    0.044    0.010    0.026    0.068
3 tau0    0.367    0.348    0.324    0.114    0.201    0.642
4 tau1    0.033    0.031    0.028    0.010    0.018    0.057
```

7209

7210 これをみると、12 球団全体で考えると、平均的に年俵が 1000 万上がると $\exp(0.046) = 1.047074$ 、つまり 1 本ホームランが増えるわけですね。ただ年俵の効果は球団ごとの違いがありますから、 $\pm 1\text{SD}$ の範囲だと係数は $0.046 + 0.033 = 0.0793$ から $0.0461 - 0.033 = 0.0127$ までぶれることがあるようです^{*12}。それよりは切片、平均 4000 万ではありますが、2900 万から 6000 万までの差がありますので、やはりどの球団に属するかの方が効果が大きいかなあ、と思ったりもするわけです。

7215 最後にこの結果をプロットしてみました。球団ごとの 50% 確信区間も合わせてプロットしてあります。これを見ると傾きや切片の球団ごとの違いがイメージしやすくなるかと思います。

7217 今日は分布を混ぜるといところから、線形モデルを階層的に組み上げていくことを説明してみました。数式的には難しいように感じるかもしれませんが、設計図を書いてコードにできれば、分析イメージは掴めるのではないのでしょうか。データがどういうところから出てきて、どういう潜在的な影響があるのか、個人差や群の差はどの程度あるのか、といった状況に応じて、モデルを柔軟にカスタマイズしていくことができます。よりデータに適切なモデルへ、より柔軟で複雑な表現を目指すことは、とてもクリエイティブで楽しいことではありませんか！

7222

26.4 課題

7224 以下のモデルを分析する R/Stan コードを提出してください。結果の解釈などを、スクリプトのコメントアウトや別添ファイルなどで提供してもらえるとなお良いです。もちろん Rmd ファイルでの提出であれば完璧です。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

7228 **■個人差を入れたモデルを考えてみよう** 野球のデータで、野手の長打率を年俵で予測するモデルに個人差を入れたモデルを考えてみましょう。データはソフトバンクの、年俵 2500 万円以上の人を対象にします。今回、長打率は安打数 Hit のうち、ホームランの本数 HR から考える割合とします。割合情報の分析ですの

7229

7230

*12 傾き β_1 の平均が γ_1 、標準偏差が τ_1 であり、それぞれの **EAP 推定値** を見えています。 γ_1 の SD、0.0120 は推定値の振幅、すなわち標準誤差ですのでご注意ください。

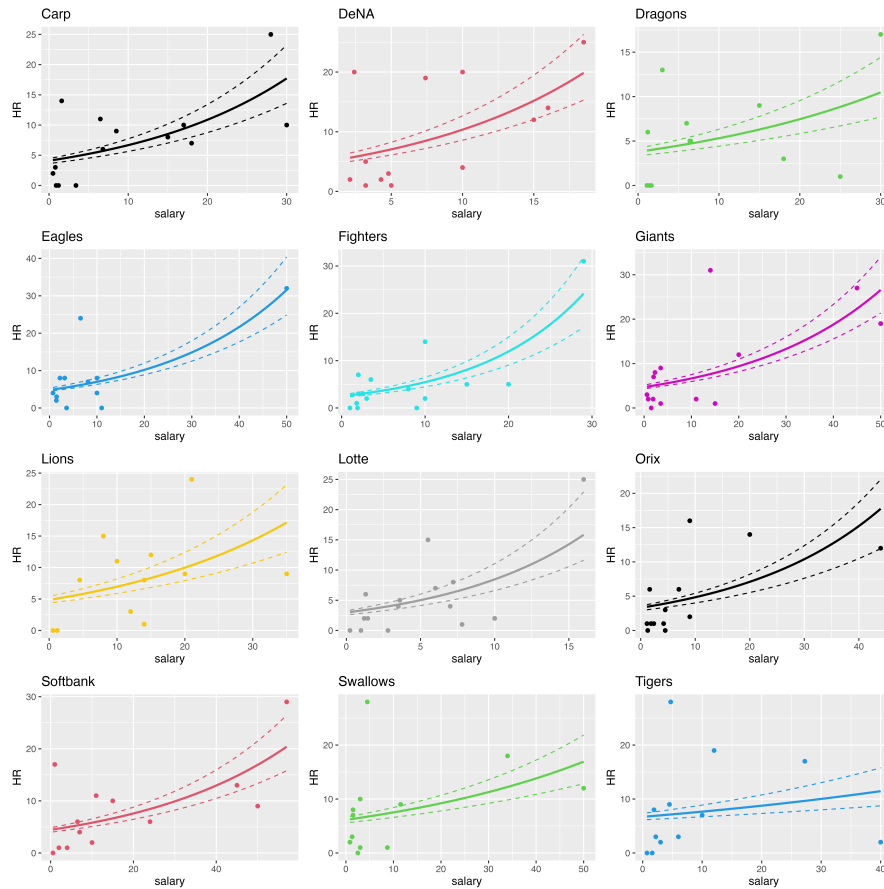


図 26.8 階層線形モデル;チームごとのポアソン回帰分析

7231 で、二項分布で表現できます。二項分布の確率 θ の線形モデルに個人差を入れ、分析してみましょう。コード
7232 [26.1](#) が参考になるはずです。

7233 ■階層線形モデルのコードを書いてみよう 野球のデータで、年俸でピッチャーの勝利数を予測する階層
7234 線形モデルを考えてみましょう。12 球団の平均や、球団ごとの散らばりはどれくらいになるでしょうか。ピッ
7235 チャーの勝利数はポアソン分布に従うと考えられますので、コード [26.3](#) が参考になるでしょう。データ加工の
7236 コードはシラバスのサイトで提供しますが、ご自身でコードを書いてもう一向に構いません。

第 27 章

混合分布モデル

さて LM→GLM→GLMM→HLM と、線形モデルを展開させてきました。後半の GLMM や HLM では分布を混ぜるということをしてきたわけですが、今回も分布を混ぜるモデルについて見ていきたいとおもいます。

今回の分布の混ぜ方は、パラメータの構造が複雑化するのではなく、そもそも違うパラメータの分布が混ざり合っているというパターンです。具体的な例から見ていきましょう。

27.1 混合分布モデル

図 27.1 を見てください。これは野球選手データから Tigers の年俸をヒストグラムにしたものです。年俸のデータはそのままだと左に歪んでいるので、対数をとったものを表示しています。ヒストグラムを見ると、どう

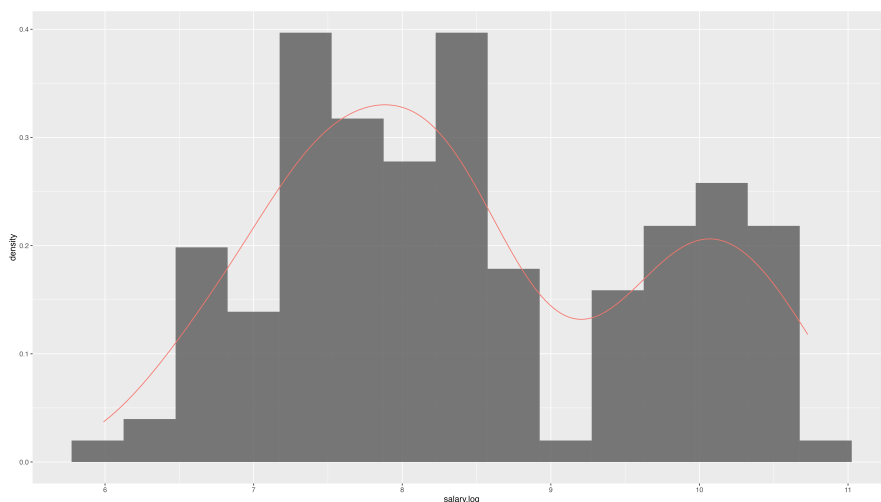


図 27.1 Tigers の選手年俸 (対数) の分布

も 2 つの山があるように見えますね。赤いラインは密度のカーブを書いたものですが、これはどう見ても 2 つのピークを持っています。このようなデータに対して、正規分布のモデルを当てはめるのはおかしいことです。正規分布はご存知の通りベルカーブ、すなわちピークは 1 つ (単峰) で、両裾になだらかに下がっていくものですから、このデータと違う形になっていることは間違いありません。

赤い密度カーブの方が暗示するように、この分布は 2 つの正規分布が混じり合っているのではないかと、思われます。つまり、出どころや種類の違う分布があるわけです。これまで分析してきたデータは、基

7253 本的に同じ分布からデータが出てきたと考えられてきたわけですが、そうでない分布の場合はどうするか、
7254 という問題です。

7255 こうしたモデルを考えるのが**混合分布モデル**と呼ばれるものです。ここでは複数の正規分布を混ぜ合わせ
7256 た**正規混合モデル (normal mixture model)**を扱います。

7257 データが複数の出自を持つという時に、その違いを示す変数が明確なのであれば困ることはありません。
7258 たとえば男性と女性で分布が違うとか、東日本と西日本で分布が違うという既有知識があって、それがデー
7259 タに含まれていれば、その変数でデータセットを分割して階層モデルにすれば良いだけです。しかし今回のよ
7260 うに、阪神の選手という意味で同じ種類のメンバーであるはずなのに、データの表れ方が違う場合、「この 2
7261 つの分布に分かれる潜在的な変数は何か」ということを探す必要があります。言い換えれば、データの表面的
7262 な特徴からグルーピングをしなければなりません。これはデータの類似性だけを用いて、外的基準を持たずし
7263 て分類する、**クラスタリング (Clustering)**という分析方法と関係します。そこで少し寄り道になりますが、
7264 クラスタリングの手法について見ておきましょう。

7265 27.1.1 クラスタリングいろいろ

7266 クラスタリングという分析手法は、先に述べたように別の変数 (外的基準) を使わずに*1、データの分類を
7267 行う方法です。分類は科学の基本で、まず同じものか違うものか、何が同じで何を違うものか考えるか、とい
7268 うことがすべてのスタートになります。クラスタリングではこれをデータから行う必要があります、データの類似性
7269 をもとに分類することが基本です。

7270 データの**類似度**あるいは同じことですが**非類似度**を表す指標の基本は**距離 (distance)**です。距離と聞
7271 くと普通は 2 点間を直線で結んだ**ユークリッド距離 (Euclidean distance)**を想像されると思いますが、
7272 実は距離には他にもいろいろな種類があります (詳しくは第 15 講, セクション 15.2, Pp.153 を参照)。たと
7273 えば**相関係数**も絶対値を取ればデータの類似度を表す距離情報の一種と考えることができます。この距離を
7274 使って、距離が短い=類似している=同じクラスター、という考え方から分類していくことになります。

7275 こうした分類, クラスタリングの手法はいくつかあって、大きく分けると**階層的クラスタリング (hierarchical
7276 clustering)**と**非階層的クラスタリング**の二種類に分かれます。階層的というのは HLM の発想と同じで、
7277 クラスターが積み重なっていくようなイメージです。たとえば要素 $\{a, b, c, d, e\}$ をクラスタリングする時、距離
7278 が近いものを使ってまず $\{a\}, \{b, c\}, \{d, e\}$ と分割したとします。次に $\{b, c\}$ を 1 つの新しいデータ点だと
7279 考えて、 a と $\{b, c\}$, $\{b, c\}$ と $\{d, e\}$ の距離を考え、近い方をさらに次のクラスターにまとめます。すなわち
7280 $\{a, \{b, c\}\}, \{d, e\}$ というようにです。それをさらに上位階層でまとめ $\{a, \{b, c\}, \{d, e\}\}$ とすべてが 1 つのク
7281 ラスターにまとめられたら終わり、とより包括的なクラスターへ、より上位のクラスターへとまとめていくのが階
7282 層的クラスタリングという手法です。この方法で分類した例を図 27.2 に示します。階層的クラスタリングは、
7283 階層が上がる時にクラスター化されたもの同士の距離をどう考えるかによって、いくつかの方法があります。
7284 先ほどの例ではまず $\{b, c\}$ が 1 つのクラスターになりましたが、これを A とすると、 a と A の距離をどう考え
7285 るか、がいろいろあるわけです。大文字の A の方は、2 つの要素から合成されたもの ($A = \{b, c\}$) ですか
7286 ら、 a と b の距離を取るのか、 a と c の距離を取るのか、はたまた b, c の平均的な値と a との距離を取るの
7287 か等々、基準が必要なわけです。複数の要素を持つクラスター同士の距離の考え方として、その最大値を取
7288 る**最大距離法**, 最小値を取る**最小距離法**だけでなく、データの重心を取る**セントロイド法**, 平均を取る**平均
7289 法**, クラスター内分散とクラスター間分散の差が最小となるように取る**ワード法 (Ward's method)** など

*1 機械学習の文脈では、こうした基準となる変数を置かず分析する手法のことを**教師なし学習**, と呼びます。

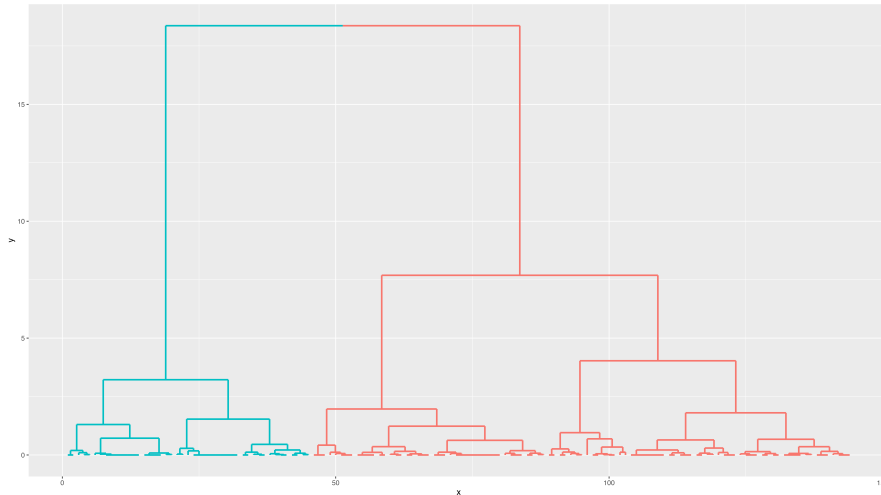


図 27.2 Tigers の選手を年俸で階層的クラスタリング (ウォード法による)

7290 があります*2*3。最後のウォード法は弁別後の結果がわかりやすいため、よく用いられます。

7291 一方、階層的でないクラスター分析で有名なのが **K-平均法 (k-means 法)** です。これは変数間の距離
7292 に k 個の基準点をランダムに置き、その点に近いものをクラスターとします。次に作られたクラスターの重心
7293 を新しい基準点として、再度分類します。その重心を次の基準点とアップデートして再分類... と繰り返して、
7294 分類結果が変わらなくなる前で何度も反復するのです。だいたい数回の反復で分類できること、データの
7295 サイズが大きくても効率よく部類できることからよく利用される方法の 1 つです。最初に幾つのクラスターに
7296 分類するか、明確な基準がないことが欠点でしたが、その後のモデルの展開でクラス数の適合度などを算出
7297 して適切なクラス数を求めるようになりました。

7298 クラスター分析のもう 1 つの分類基準が、クラスターの境界強度によるものです。**ハードクラスタリング**
7299 (**Hard clustering**) では、ある個体がどのクラスターに含まれるかがはっきりしています。この個体はクラ
7300 スター A、この個体はクラスター B、と決められたらそれで決まりです。これに対して**ソフトクラスタリング**
7301 (**Soft clustering**) は、ある個体がクラスター A に含まれる確率が XX%、クラスター B に含まれる確率が
7302 YY%、のように分類が確率的で固定的ではない方法です。

7303 このように、データの特徴から分類をしていくクラスター分析はいろいろな手法があるのですが、今回扱う
7304 混合分布モデルはモデルに基づいたクラスタリング、Model based clustering と呼ばれる手法です。すな
7305 わち、各個体が潜在的な正規分布のどちらに含まれる可能性があるか、その確率を考えながらソフトに分類し
7306 ていく非階層的な方法になります。

*2 厳密に書くと、クラスター P に含まれる要素と重心との距離の二乗和を $E(P)$ とすると、二つのクラスター P,Q を併合して作られる新しいクラスター T を考えた時、 $E(T) - E(P) - E(Q)$ を P,Q の距離と考える方法です。

*3 R では標準関数 `hclust` で階層的クラスター分析ができ、ウォード法を使うときは `method='ward.D2'` とします。この .D2 にはいわれがあって、かつては `ward.D` でウォード法の指定ができたのですが、この関数にバグが含まれていることが判明、修正版を `ward.D2` としたのです。R は無償のフリーソフトウェアですから、バグが含まれているなんてやっばり使い物にならない！ という批判もあるのですが、オープンなソフトウェアであるからこそソースコードが誰にでも見られて指摘することができ、よくなっていくものでもあります。バグがあったこと、それを自浄作用でフィックスしたことを記録するために、2 を残すことになりました。プロプライエタリな商用ソフトの場合、間違いがあってもこっそり修正して、同じ関数名で提供されてしまえば、ユーザは分析が間違っていたことに気づくことすらできません。科学技術に関するソフトウェアはオープンであるべきだ、という R の思想がここに現れていると言えるでしょう。

7307 27.1.2 混合分布モデルの考え方

7308 混合分布モデルの考え方は、アルクスターに含まれるかどうか確率的に決まるというものです。たとえ
7309 ばコイントスをして、表が出たらクラスター A、裏が出たらクラスター B、というように分類するようなもので
7310 す。コイントスはベルヌーイ分布に下がいますから、確率 θ でクラスター A、確率 $1 - \theta$ でクラスター B、と
言っていることと同じです (図 27.3)。

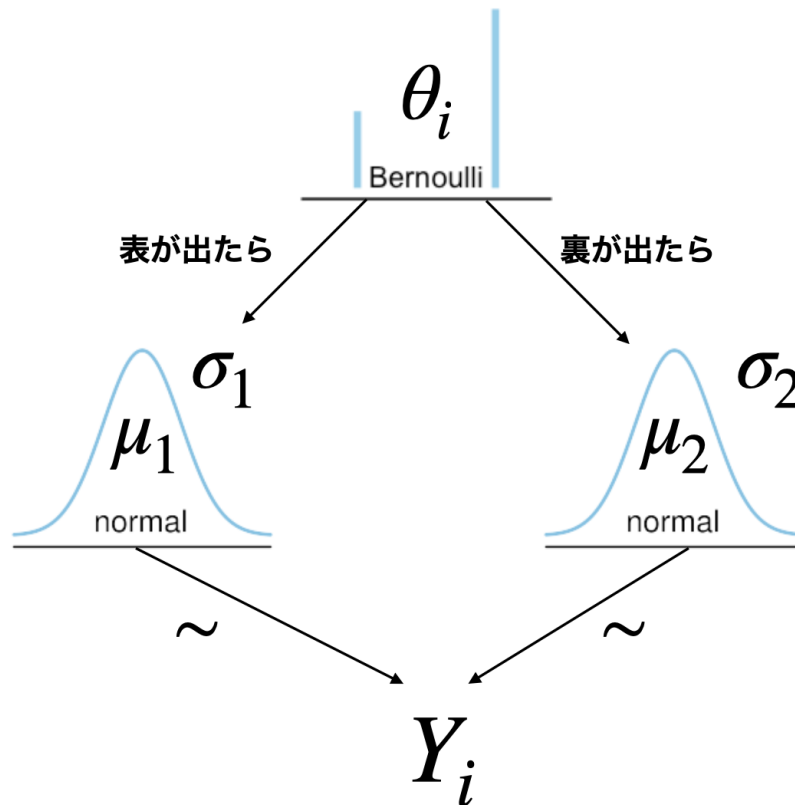


図 27.3 二つの正規分布からデータが出てくるモデル

7311 もちろん 3 つ以上のクラスターに分類することもあるでしょう。この場合はカテゴリーカル分布 (セクション
7312 24.1, Pp.263 を参照) に従う確率変数 z を考えることになります。

7314 データが K 個の正規分布、いずれかから得られていると考えるとしましょう。ここで各正規分布はそれぞ
7315 れの位置母数 μ_k とスケール母数 σ_k を持っているとします。個々のデータ i がどの正規分布からきてい
7316 るかを表す混合確率 λ_{ik} を考えると、 $\lambda_{ik} \geq 0$ で $\sum_{k=1}^K \lambda_{ik} = 1$ になります。つまり、 i ごとの K -simplex
7317 ベクトルを考えることになります。個々のデータ Y_i の出自は、 $z[i] \sim \text{Categorical}(\lambda_i)$ で決まり、データ
7318 $Y_i \sim N(\mu_{z[i]}, \sigma_{z[i]})$ からきている、と考えることになります (図 27.4)。

7319 データ生成メカニズムの設計図は図 27.3, 27.4 の通りです。具体的なプログラミングに進む前に、イメー
7320 ジをしっかり掴んでおきましょう。これらの設計図の一番上にある、どのクラスターからデータが出てくるかを
7321 確率的に決めるというところがポイントです。ですがこのモデルを実装できるようになれば、データ発生の特
7322 条件分岐を確率モデルで表現できるようになるのです。確率 θ で左、 $1 - \theta$ で右のルートを取る、というように、

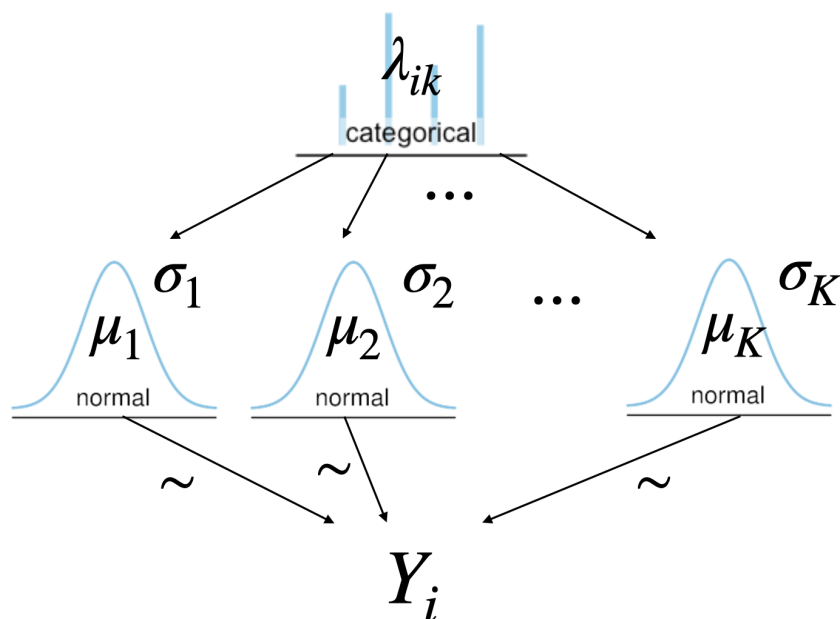


図 27.4 K 個の正規分布からデータが出てくるモデル

7323 データが出てくるメカニズムの背後に潜在的な条件を設定できますから、分析の可能性がグッと広がること
 7324 になります。実践的にはたとえばマーケティングにおいて、顧客がどういうクラスター（主婦層、独身者、大家族
 7325 etc...）に属するかはわかりませんが、買い物の傾向から分類するというようなことができます。あるいは心理
 7326 学的なシーンでも、潜在的に異なる種類（抑うつ傾向があるタイプ、ないタイプなど）が、反応傾向の違いから
 7327 読み取れるようになるかもしれません。

7328 クラスター分析はデータから分類するので、分類後にそれがどう言った変数と関係しているのか、あるいは
 7329 分類されたクラスターにはどういう特徴があるのかを探索することが一般的です。それでもデータだけから分
 7330 類できるのは魅力的な手法であると言っていいでしょう。

7331 27.2 ターゲット記法と周辺化消去

7332 さてこの条件分岐をどのようにコードに書いていくか、というところに話を進めましょう。ここで新しく、2
 7333 つの技術を導入する必要があります。1つはターゲット記法、もう1つは周辺化消去 (marginalizing)
 7334 です。

7335 27.2.1 ターゲット記法

7336 これまで Stan で、確率変数に従う分布の記法は \sim を使って表現していました。すなわち、コード 27.1 の
 7337 ような書き方です。

code : 27.1 sampling 記法

```
7338
7339 1 for(i in 1:N){
7340 2     Y[i] ~ normal(mu, sigma);
7341 3 }
```

7343 これと同じ動きをする、別の表記方法があります。それが**ターゲット記法**というもので、コード 27.2 のように
7344 書きます。

code : 27.2 target 記法

```
7345 1 for(i in 1:N){
7346 2     target += normal_lpdf( Y[i] | mu, sigma);
7347 3 }
```

7350 何やら奇妙な書き方に見えますが、順番に見て行きましょう。まず += の右側、normal が normal_lpdf
7351 になっているところに注目してください。ここで lpdf は log probability density function, すなわち対数
7352 確率密度関数という意味です。normal というのがついてますから、正規分布の対数確率密度関数と
7353 いうことになりますね。ちなみに離散変数の場合は**確率密度 (probability density)**ではなく**確率質量**
7354 **(probability mass)** になりますから、lpmf, すなわち log probability mass function (対数確率質量関
7355 数) になります。具体的には、bernoulli_lpmf とか poisson_lpmf となりますので注意してください。

7356 さて対数確率密度関数ってのがなぜ出てきたのでしょうか。ここで**尤度 (likelihood)** のことを思い出し
7357 てほしいのですが、尤度とはデータが**確率分布**から出てくるときの尤もらしき、出てきやすさのような指標な
7358 のでした。確率分布はパラメータがわかっている時にデータが出てくる確率を記述する関数ですが、同じ関
7359 数でパラメータが未知、データが既知の場合 (多くの研究シーンはこちらですが) に、それを尤度関数とよぶ
7360 のでした。データが得られた時の尤度を計算し、その尤度が最も大きくなる場所をパラメータの推定値とす
7361 るのが**最尤推定法 (Maximum likelihood method)** であり、尤度を使って事後分布を算出、事後分
7362 布の形から推定値を考えるのが**ベイズ推定法 (Bayesian inference method)** なのでしたね。つまり
7363 MCMC で計算するにはデータ全体の尤度を考えているのです。そしてデータ全体尤度は各データ点の尤度
7364 の総積 \prod で考える必要がありますが、確率の掛け算は数字が小さくなるので**対数尤度 (log-likelihood)**
7365 の総和で計算するのです。

7366 具体的に書くと、データ $\mathbf{Y} = Y_1, Y_2, Y_3$ がパラメータ θ のベルヌーイ分布から出てきている場合、データ
7367 の尤度は

$$L(\mathbf{Y}|\theta) = \prod_{i=1}^3 \text{bernoulli}(Y_i|\theta) = \text{bernoulli}(Y_1|\theta) \times \text{bernoulli}(Y_2|\theta) \times \text{bernoulli}(Y_3|\theta)$$

7368 であり、対数尤度は

$$LL(\mathbf{Y}|\theta) = \sum_{i=1}^3 \log(\text{bernoulli}(Y_i|\theta)) = \log(\text{bernoulli}(Y_1|\theta)) + \log(\text{bernoulli}(Y_2|\theta)) + \log(\text{bernoulli}(Y_3|\theta))$$

7369 となるわけです。ここで $\text{bernoulli}(Y|\theta)$ の縦棒 | は、 θ が与えられた時の Y の出る確率、という条件付き確
7370 率の記号になります。

7371 プログラミング的には、 Y_1, Y_2, \dots, Y_n は for(i in 1:N){ ~ } という形で 1 から N までの繰り返し計
7372 算を意味しますから、コード 27.1 や 27.2 でやっていることは、各データ点についての対数尤度を足し合わせ
7373 て行っている作業そのものを意味しています。

7374 ここでターゲット記法にある += ですが、これはプログラミング独特の表記方法で、数学表記では
7375 ありません。プログラミングでは代入を使って、たとえば変数 A に 1 を加えたものを新しい A と
7376 する (上書きする)、ということを $A = A + 1$ と表記できます。この += はこうした上書きを一言で
7377 書く記号で、この A の例だと $A += 1$ とすれば 1 を加えて上書き、と同じ意味になります。これを
7378 踏まえてターゲット記法を見ると、target += normal_lpdf(Y[i] | mu, sigma) ですから、
7379 target = target + normal_lpdf(Y[i] | mu, sigma) としていることと同じです。つまり、対数尤

7380 度を足しあげて行ってね、ということを示明的に書いていることになります。target は足し合わされたもの、
7381 というだけで「目標はこれ」を意味するぐらいの予約語だと思ってください。

7382 サンプル記法とターゲット記法は、書き方が違えど同じ働きをするのですからどちらでもいいじゃない
7383 か、と思うかもしれませんが、ターゲット記法でないと表現できないようなことがあるので、このような代替手
7384 法が用意されているのだと思ってください。ターゲット記法でないと表現できないようなこと、というのはもち
7385 ろん今回の、分布を混ぜ合わせる時がそれです*4。

7386 27.2.2 周辺化消去

7387 さてここまでさまざまな問題を解決してくれた Stan ですが、Stan には離散型のパラメータを扱うことがで
7388 きない、という欠点があります*5。今回は $z[i] \sim \text{Categorical}(\lambda_i)$ と、どの目が出たかというパラメータが離
7389 散的ですから、このままではコード化できません。ではどうするか。仕方がないので、「起こりうる可能性をす
7390 べて数え上げて考える」という方法を取ることにします。最終的な事後分布は対数尤度関数を足し合わせた
7391 ものになりますから、条件分岐したルートすべてについて対数尤度関数を考えて、それを足し合わせてやれば
7392 良いのです。

7393 簡単な例として図 27.3 にあるような、2 つの正規分布モデルを考えてみましょう。コインスで表が出れば
7394 $N(Y_i | \mu_1, \sigma_1)$ 、裏が出れば $N(Y_i | \mu_2, \sigma_2)$ のルートに行けば良いのです。表が出る確率は θ ですから、表
7395 が出るルートを通してデータ Y_i が出てくる確率は $\theta \times N(Y_i | \mu_1, \sigma_1)$ になります。同じく裏が出るルートを通
7396 してデータが出てくる確率は、 $(1 - \theta) \times N(Y_i | \mu_2, \sigma_2)$ ということになります。

7397 ですから、最終的にデータ Y_i が出てくる確率は $\theta \times N(Y_i | \mu_1, \sigma_1) + (1 - \theta) \times N(Y_i | \mu_2, \sigma_2)$ だ、とい
7398 うことができるわけです。

7399 さて、これらは確率をそのまま使った表現になっていますが、Stan の中では対数で考えるのでした。つまり
7400 表が出るルートの対数尤度は

$$\log(\theta) + \log(N(Y_i | \mu_1, \sigma_1))$$

7401 ですし、裏が出るルートもまた同様に \log をとって考える必要があります。対数を取ると積が和になるので、
7402 $\times \rightarrow +$ に変わっていることに注意してください。

7403 ところがこの表ルート、裏ルートは足し合わせないといけません。対数をとったものを直接足し合わせるわ
7404 けにはいきませんから、対数尤度をただの尤度に戻す、つまり \log の逆関数である \exp 関数を通す必要があ
7405 ります。ややこしいですが、

$$\exp(\log(\theta) + \log(N(Y_i | \mu_1, \sigma_1))) + \exp(\log(1 - \theta) + \log(N(Y_i | \mu_2, \sigma_2)))$$

7406 が尤度であり、これの対数をとったものを target に加える必要があるのです。

7407 このように、 \exp を通して、足し合わせ (\sum) て、さらにその \log をとるという作業をひとまとめにした関
7408 数が、Stan にはあります。それが \log_sum_exp 関数です。この関数は、 $\log_sum_exp(x, y)$ とすると
7409 $\log(\exp(x) + \exp(y))$ という操作をします。この関数に渡すものがベクトルであれば、

$$\log_sum_exp \left(\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \right) = \log \left(\sum_{i=1}^n \exp(x_i) \right)$$

*4 そのほかにもターゲット記法は、サンプル記法では省略されている事後対数確率も記録するという特徴があり、ブリッジサンプリングによる事後対数尤度の計算などにはこちらの表記が必要になってきます。もっとも非常に専門的なことなので、ここでは条件分岐の際に必要な記法、と割り切っていただいても構いません。

*5 JAGS という確率的プログラミング言語でしたら、これは可能です。

7410 という操作をすることになります。
 7411 これを使って混合分布モデルの実装をしていきましょう。

7412 27.2.3 混合分布モデルのコード

7413 先ほどの、阪神タイガースにおける選手の年俸分布を 2 つの正規分布から出てきているもの、と考えて
 7414 コード化したものは、Code27.3 のようになります。

code : 27.3 混合正規モデル

```

7415 1 data {
7416 2   int<lower=1> K;
7417 3   int<lower=1> L;
7418 4   array[L] real Y;
7419 5 }
7420 6
7421 7 parameters {
7422 8   array[L] simplex[K] theta;
7423 9   ordered[K] mu;
7424 10  array[K] real<lower=0> sigma;
7425 11 }
7426 12
7427 13 transformed parameters{
7428 14   array[L] vector[K] lp;
7429 15   for (l in 1:L) {
7430 16     for (k in 1:K) {
7431 17       lp[l,k] = log(theta[l,k])+ normal_lpdf(Y[l]|mu[k],sigma[k]);
7432 18     }
7433 19   }
7434 20 }
7435 21
7436 22 model{
7437 23   for(l in 1:L){
7438 24     target+=log_sum_exp(lp[l]);
7439 25   }
7440 26   sigma ~ cauchy(0,5);
7441 27   mu ~ normal(0,10);
7442 28 }
7443 29
7444 30 generated quantities{
7445 31   array[L] vector[K] prob_class;
7446 32   array[L] int<lower=1,upper=K> pred_class;
7447 33   for(l in 1:L){
7448 34     prob_class[l] = softmax(lp[l]);
7449 35     pred_class[l] = categorical_rng(prob_class[l]);
7450 36   }
7451 37 }
7452
7453

```

7454 ■コード解説

7455 data ブロック クラスター数 K, データ長 L, アウトカム変数 Y を宣言しています。

7456 **parameters ブロック** 合計が 1 になる K の長さを持つベクトルである, simplex 型ベクトルを, データ数
 7457 だけ用意します。また, K 個の正規分布からデータが出てきているのですから, 平均と標準偏差も K
 7458 個用意します。ここで大事なポイントとして, 平均 μ が ordered 型としてあるところです。このベクトル
 7459 は, 要素 $\mu_1 > \mu_2 > \dots > \mu_k$ と, 要素が大ききの順に並んでいることです。このような制約をかけて
 7460 おかないと, K 個の正規分布がどの位置にあるのか判然としな苦なるという問題が生じます。たとえ
 7461 ばあるチェーンで $\mu_1 > \mu_2$ となっていて, データ X_k が 1 番目の正規分布から出てきていると考えら
 7462 れたとしましょう。しかし別のチェーンで $\mu_1 < \mu_2$ となっていれば, X_k は平均値が大きい方から出て
 7463 きているのですから, こちらのチェーンでは 2 番目の正規分布から出てきたものとして, MCMC サン
 7464 プルが進むこととなります。複数のチェーンから得られた MCMC サンプルは最終的に統合されます
 7465 から, このようになっていると X_k は正規分布 1 からも 2 からも同じように出てきていることになりま
 7466 す。各チェーンでは正しくサンプリングできているのですが, チェインごとに分布のラベルが違うので,
 7467 統合すると訳がわからなくなる, この問題をとくにラベルスイッチング (Label Switching) 問題と
 7468 いい, 混合分布モデルではよく生じる問題なのです。そこで分布の位置母数に順番の制約を入れて,
 7469 ラベルの変動を止める必要があります。

7470 **transformed parameters ブロック** 長さ K のベクトルをデータの数 L だけ用意した 1p 変数を作りました
 7471 た。これは log-probability の略でこのような名前にしましたが, 任意の名前で結構です。1 行目の
 7472 データがクラスター k に含まれる可能性は, $\theta_{lk} \times N(Y[l] | \mu_k, \sigma_k)$ ですから, その対数をとった
 7473 形で表現しています。

7474 **model ブロック** モデル尤度のところは, 対数 \rightarrow exp 関数で確率の形 \rightarrow ベクトルの要素をすべて足し合
 7475 わせる \rightarrow log をとって対数尤度にする, という計算をまとめて log_sum_exp 関数で行い, それをす
 7476 べて target に追加していくという方法で書いています。事前分布のところは普通のサンプリング記
 7477 法で書きました。ちなみに μ, σ はそれぞれ K 個ずつありますが, Stan は添字の省略を許してくれ
 7478 ます。

7479 **generated quantities ブロック** 久しぶりに生成量ブロックの登場です。ここでは分析結果を使って, 1 行
 7480 目のデータが結局何番目の正規分布から出てきたのかを逆算して出しています。まず, 対数尤度で表
 7481 現されている K 個の要素を持つベクトルを, 確率の形に計算したものを preb_class 変数にしてい
 7482 ます。ここで使われている softmax という関数は,

$$\text{softmax}(x) = \frac{\exp(x)}{\sum_{k=1}^K \exp(x_k)}$$

7483 という計算をするものです。分母はベクトルの総和で分子がその要素ですから, 要するに各クラスに
 7484 入る確率を表すベクトルになっている訳です。これを使って pred_class 変数を作っています。今
 7485 回のデータから予想される所属クラスなので, K 面サイコロを振った出目という形で表現します。
 7486 categorical_rng はカテゴリカル分布からの乱数発生で, その発生確率が先ほどの prob_class
 7487 になっているので, この乱数で出た目が予測されるクラスということです。

7488 これを実行するための R コードをコード 27.4 に用意しました。

code : 27.4 混合分布モデルのコード

```
7489 1 dat <- read_csv("baseballDecade.csv")
7490 2 dat.tmp <- dat %>%
7491 3   filter(position != "投手") %>%
7492 4   filter(Games > 50) %>%
7493 5   filter(team == "Tigers") %>%
```



```

7495 6 mutate(salary.log = log(salary)) %>%
7496 7 dplyr::select(salary.log, Name) %>%
7497 8 rowid_to_column("ID")
7498 9
7499 10 dataSet <- list(K = 2, L = NROW(dat.tmp), Y = dat.tmp$salary.log)
7500 11 model <- cmdstanr::cmdstan_model("latent.stan")
7501 12 fit <- model$sample(data = dataSet, chains = 4, parallel_chains = 4)
7502

```

7503 ■コード解説

7504 1 行目 データファイル baseballDecade.csv を読み込みます。

7505 2-8 行目 分析用のデータに加工しています。

7506 3 行目 フィルターをかけてデータを野手だけのものにします。

7507 4 行目 フィルターをさらにかけて、50 試合場出ている選手に限定します。

7508 5 行目 チームが Tigers の選手だけに絞っています。

7509 6 行目 年俸のデータは歪んでいるので、対数をとって正規分布の形に近似させました*6。

7510 7 行目 使う変数だけに絞っています。

7511 8 行目 行番号を ID という変数名にして追加しました

7512 10 行目 データセットを作っています。

7513 11 行目 cmdstan でコンパイルしています。rstan パッケージを使う場合は rstan のコンパイル関数を使っ
7514 てください。

7515 12 行目 cmdstan のサンプリング関数です。rstan パッケージを使う場合は rstan のサンプリング関数を使
7516 ってください。

7517 推定した結果を見てみましょう。私の環境では出力 10 のようになりました。

MCMC の結果 10

```

# A tibble: 4 × 7
  name          EAP      MED      MAP      SD      L95      U95
<chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu[1]      7.624    7.624    7.638    0.107    7.415    7.826
2 mu[2]      9.702    9.699   10.017    0.284    9.164   10.138
3 sigma[1]    0.656    0.653    0.646    0.063    0.544    0.788
4 sigma[2]    0.762    0.786    0.901    0.237    0.380    1.192

```

7518
7519 対数をとった時の平均値が 7.624 と 9.702 ですから、指数関数を入れて読み直せば $\exp(7.624) =$
7520 $2046.733, \exp(9.702) = 16350.28x$ となります。2000 万円クラスの選手と、1 億 6500 万円クラスの選手
7521 の 2 つのグループに分かれることが見て取れますね。事後予測分布から、所得平均の高い「超一流」選手と
7522 「一般的」な選手とに分けて分布を色分けしてみました (図 27.5)。データから、この選手がどちらのクラスに
7523 所属しているかも想像できる場所がおもしろいですね。

*6 データを変形させずに確率分布の方で対応するのがよい、と常々話してきましたが、ここではデータを変形させています。実は左に歪んだような分布は対数正規分布というのがあり、そこからのデータ生成を考えても良いのですが、対数正規分布はデータの対数が従う正規分布というだけなので、データの方を変えてしまった方が説明しやすいのです。比率や度数など、データの生成メカニズムがそもそも違うばあいは、データを変換せずに使うのが正しいやり方です。

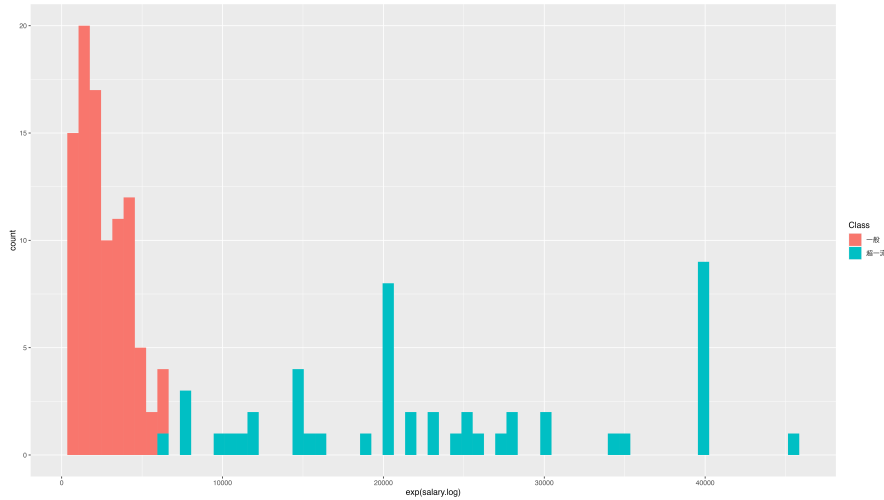


図 27.5 選手のクラスタリング

27.3 ゼロ過剰ポアソン分布モデル

さて、異なる分布の混ぜ合わせ、離散確率分布による条件分岐の例として、もう 1 つ別の例を示しましょう。次の図 27.6 を見てください。

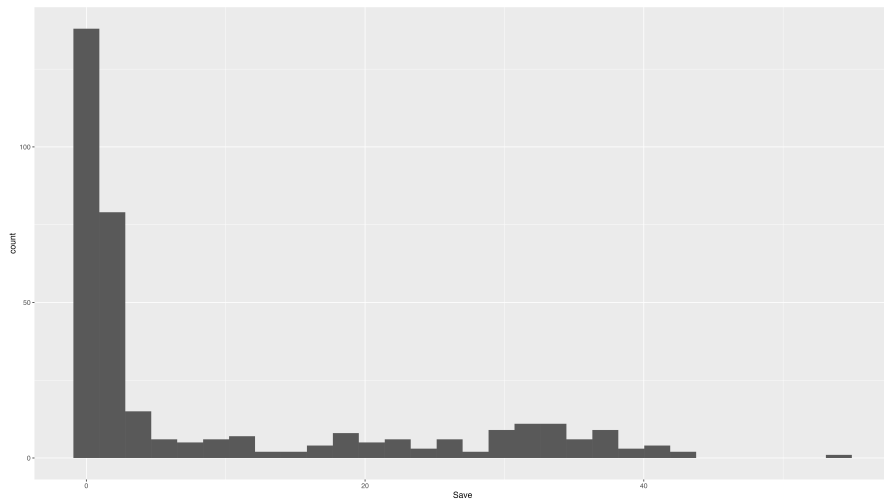


図 27.6 投手のセーブ数

これは投手のセーブ数の度数分布です。野球は 9 回のターン制で、1 回の最初から投げる投手は先発投手と言います。先発投手が 9 回まで投げ切れればいいのですが、疲労も溜まるし戦術によっては投手の打席に代打を投入することもあります。投手が降板した場合は、次の投手が代わって続きから投げますが、「勝っている試合で後半まできたので、最後の 1, 2 回は確実に勝ち逃げできるようにしたい」というときに投入されるのが抑え（クローザー）と呼ばれる投手です。最近のプロ野球は分業が進んでいますから、先発投手が完投することは少なく、100 球程度で交代し、中継ぎ - 抑えと投手リレーしていくのが一般的な戦術です。さて、先発投手は最初から出てきて 6, 7 回まで投げますから（途中でボロボロに打たれたり怪我で降板などが

なければ, ですが), 毎日のように登板することはなく, 一度登板すると 3, 4 日休んで次の試合にでる, ということになります。日本プロ野球は年間 140 試合近くありますが, 通年で 2-30 回も登板することはありません。一方, 抑えの投手は最後の 8, 9 回を投げる程度で, しかも負けている試合などでは出番がありません。必然的に中 4 日の休憩なども入りませんから, 試合数はたくさん出ることができます。そして抑えに成功するとセーブポイントというのがつきますが, このセーブポイントは優秀なクローザーであれば年間数十ポイント取ることができる訳です。

さて, 図 27.6 にあるようにこのセーブポイントを見てみますと, ほとんどの投手がゼロになっています。以下データをみますと, 1 セーブが 53 人, 2 セーブが 26 人...となって, 20~40 セーブある投手が 10 人弱ずつぐらいいる, という分布をしています。これは先ほども言った分業制のせいで, 先発投手にセーブポイントがつくことはありませんし, 先発投手は 1 チームに 10 人弱ぐらいは揃えているでしょうから, 投手全体で見るとセーブポイントは 0 の人が多いのです。

このセーブポイント, 0 以上の整数を取りますからポアソン分布に従うと考えられるのですが, ポアソン分布のパラメータ λ をどう変えても (図 27.7), こんな形にはなりません。

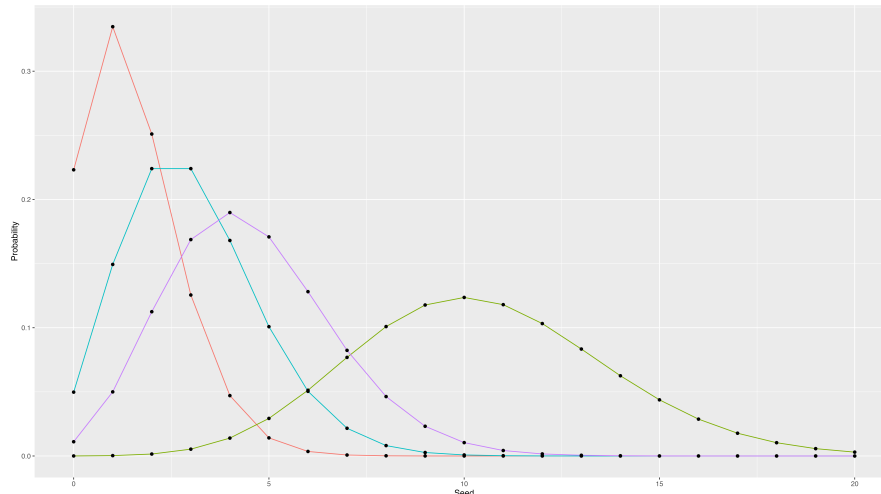


図 27.7 さまざまなポアソン分布。 λ は 1.5, 3, 4.5, 10.5 で描画

そこで混合分布です。まずセーブ数が 0 かどうかを考えます。セーブ数 0 というのは, そもそもクローザーではないか, クローザーなのに 0 セーブポイントのへっぽこか, の二択です。セーブ数が 0 でない限り, ポアソン分布からセーブ数が生成される, と考えるのです。設計図のイメージは図 27.8 の通りです。

このモデルがちょっと特殊なのは, ご覧の通り 0 というデータの値一箇所に極端な偏りがあるところ。この分布のことを**ゼロ過剰ポアソン (Zero-Inflated Poisson)** といいます。

注意してほしいのですが, ポアソン分布から 0 の度数が出てくることもあり得ます。逆に考えると, データ $Y_i = 0$ であれば, これはクローザーでないからそうなっているのか, クローザーなのにへっぽこでセーブ数がつかなかったのかのどちらか, という判断をしなければならぬということです。つまり, データからの分岐を考えることになります。クローザーかどうかは確率 θ で決まるとして, θ ならクローザー, $1 - \theta$ ならクローザーじゃないとすると,

$$\begin{cases} \text{データが 0 ではない} & \rightarrow \theta \times \text{Poisson}(Y_i | \lambda) \\ \text{データが 0 だ} & \rightarrow 1 - \theta \text{ or } \theta \times \text{Poisson}(0 | \lambda) \end{cases} \quad (27.1)$$

このデータが 0 のときに生成メカニズムが条件分岐していますから, ここに `log_sum_exp` 関数を入れる必

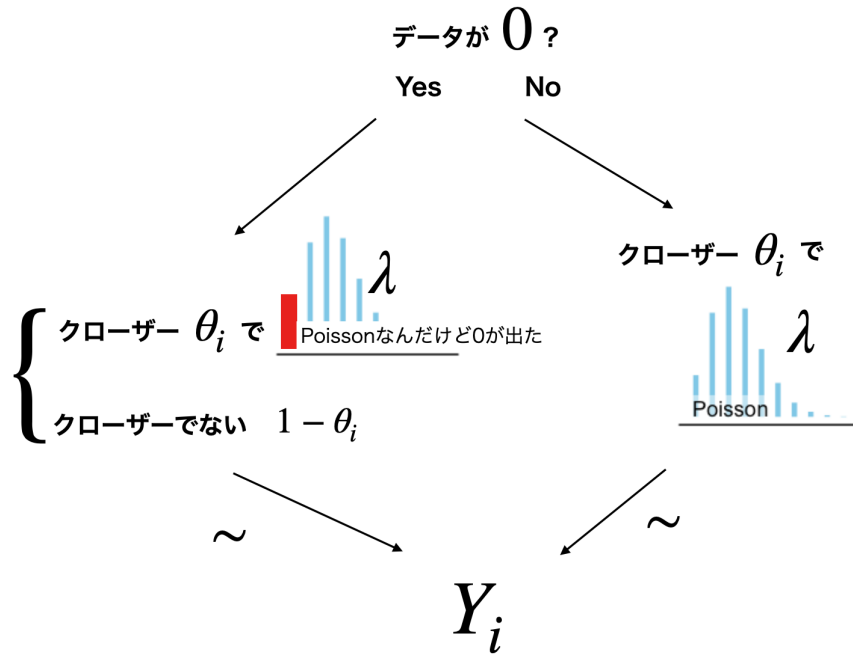


図 27.8 ポアソン分布との混ぜ合わせ

7558 必要があります。

$$\log_sum_exp \left(\begin{array}{c} \log(1 - \theta) \\ \log(\theta) + \text{poisson_lpmf}(0 | \lambda) \end{array} \right)$$

7559 このことを念頭において、コードを書いてみましょう (コード 27.5)。

code : 27.5 ゼロ過剰ポアソン分布

```

7560
7561 1 data{
7562 2   int L;
7563 3   array[L] int Y;
7564 4 }
7565 5
7566 6 parameters{
7567 7   real<lower=0,upper=1> theta;
7568 8   real<lower=0> lambda;
7569 9 }
7570 10
7571 11 model{
7572 12   for(1 in 1:L){
7573 13     if(Y[1]==0){
7574 14       target += log_sum_exp(log(1-theta), log(theta)+poisson_lpmf(0|lambda));
7575 15     }else{
7576 16       target += log(theta) + poisson_lpmf(Y[1]|lambda);
7577 17     }
7578 18   }
7579 19 }
7580

```

7581 ■コード解説

7582 data ブロック データ長 L , アウトカム変数 Y を宣言しています。
 7583 parameters ブロック クローザーかどうかを決める θ と, ポアソン分布のパラメータ λ がパラメータです。
 7584 model ブロック データの中で, プログラミング言語としての if 文による条件分岐が出てきています*7。
 7585 データが 0 であれば, クローザーでないか, ポアソン分布からゼロが出ているかです。そうでなければ,
 7586 確率 θ 経由でのポアソン分布です。
 7587 これを実行するための R コードはコード 27.6 のようになります。

code : 27.6 ゼロ過剰ポアソンモデルのコード

```
7588 1 dat.tmp <- dat %>%
7589 2   dplyr::filter(position == "投手") %>%
7590 3   dplyr::filter(Games > 50) %>%
7591 4   dplyr::select(Save)
7592 5
7593 6 model <- cmdstan_model("ziPoisson.stan")
7594 7 dataSet <- list(L = NROW(dat.tmp), Y = dat.tmp$Save)
7595 8 fit <- model$sample(
7596 9   data = dataSet,
7597 10  chains = 4,
7598 11  parallel_chains = 4
7599 12 )
7600
7601
```

7602 推定した結果を見てみましょう。私の環境では出力 11 のようになりました。

MCMC の結果 11

```
# A tibble: 2 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 lambda  14.666  14.665  14.661  0.261  14.154  15.185
2 theta   0.606   0.606   0.608   0.025   0.556   0.655
```

7603
 7604 これを見ると, 6 割の確率でクローザー判定され, クローザーであれば $\lambda = 14.7$ ぐらいのポアソン分布で
 7605 データが生成されるということですね。事後予測分布のように, この確率モデルで出てくるデータを図 27.9 に
 7606 示してみました。単一の分布では表現できないモデル化ができていないのでしょうか。もっとも, 今回
 7607 は 1, 2 回セーブポイントがついている人のところをうまく表現できませんでした。しかしコードを書き換えるこ
 7608 とで, セーブ数 0 の先発投手と, 本職のクローザーに加え, セーブ数が一桁ぐらいの中継ぎ投手の三種類に
 7609 分割する, というのもできるでしょう。

7610 ここでみた分布の混ぜ方は, たとえば「お店に初来店した人とリピーターの来店数」とか, 「20 歳までの恋愛
 7611 経験の数」のようなデータがあったときに, その特徴から分割してモデルを調整できます。研究のシーンでも,
 7612 たとえば「質問紙で何も考えずにどちらとも言えないと答える人と, ちゃんと考えてどちらとも言えないに丸を
 7613 つけた人」というように回答者のバイアスを除去してその特徴を考えるなど*8, 積極的な分析をできるのです。

7614 さらに言えば, このポアソン分布の λ に線形モデルを入れて, ゼロ過剰ポアソン回帰分析にすることだって
 7615 できますね。たとえばクローザーのセーブ数が年俸によって説明される, というモデルにするなら, コード 27.7

*7 プログラミングでの条件分岐については第 16 講, セクション 16.3.3, Pp.172 でやりましたね。忘れた人は戻って再確認しましょう。

*8 この場合は中点の 3 が多い, というようなケースでしょうから, 3 過剰正規分布とでも言えば良いでしょうか。応用例として清水 (2018) などがあります。

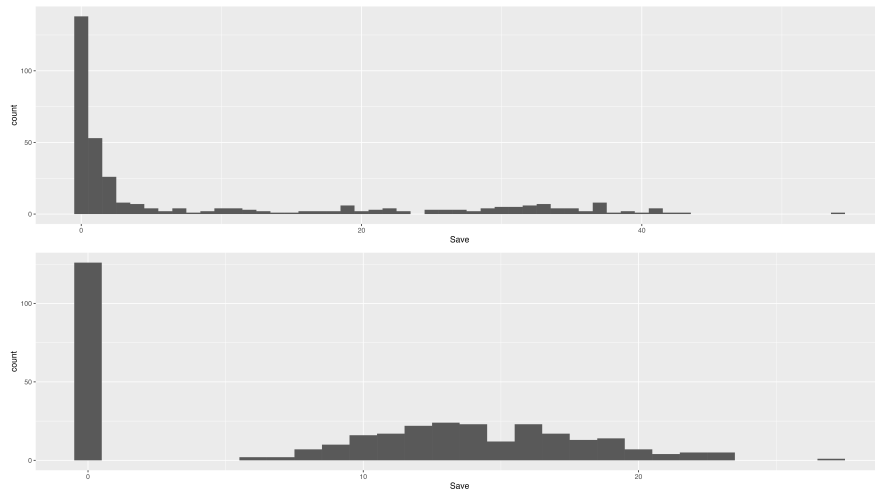


図 27.9 ゼロ過剰ポアソン分布による当てはめ。上がロウデータ, 下が事後予測分布

7616 のようにすれば良いのです。

code : 27.7 ゼロ過剰ポアソン回帰

```

7617
7618 1 data{
7619 2   int L;
7620 3   array[L] int Y;
7621 4   array[L] real X;
7622 5 }
7623 6
7624 7 parameters{
7625 8   real<lower=0,upper=1> theta;
7626 9   real beta0;
7627 10  real beta1;
7628 11 }
7629 12
7630 13 transformed parameters{
7631 14   real<lower=0> lambda[L];
7632 15   for(l in 1:L){
7633 16     lambda[l] = exp(beta0 + beta1 * X[l]);
7634 17   }
7635 18 }
7636 19
7637 20 model{
7638 21   for(l in 1:L){
7639 22     if(Y[l]==0){
7640 23       target += log_sum_exp(log(1-theta),log(theta)+poisson_lpmf(0|lambda[l]));
7641 24     }else{
7642 25       target += log(theta) + poisson_lpmf(Y[l]|lambda[l]);
7643 26     }
7644 27   }
7645 28 }
7646

```

7647 これは一般化線形モデルの応用ですから、ポアソン分布に限らず使える技であることは、すぐにお分かりいた

7648 だけのことと思います。

7649 データの分布や特徴を考えて、それにあった分析モデルを作っていくのは、モデラーとしての階段を一步
7650 登ったこととなります。図 27.10 には統計分析家の成長ステップの模式図を描いてみました*9。最初はデータ
7651 の記述や線形モデルの当てはめなど、得られたデータに振り回されるような形でモデルを作り上げてきました
7652 が、今回のように分布を混ぜ合わせると、さらにその表現力が豊かになることが実感できたのではないでしょ
7653 うか。図にあるように、確率モデルのパラメータが心理事象の(数学的)記述から導出され、それがデータと合
7654 致する可動かで検証するようになれば、さらにもう一段階研究のレベルが上がることになるでしょう。最終
的には、心理学理論の方から新しい確率分布が導出される日が来るかもしれませんね。

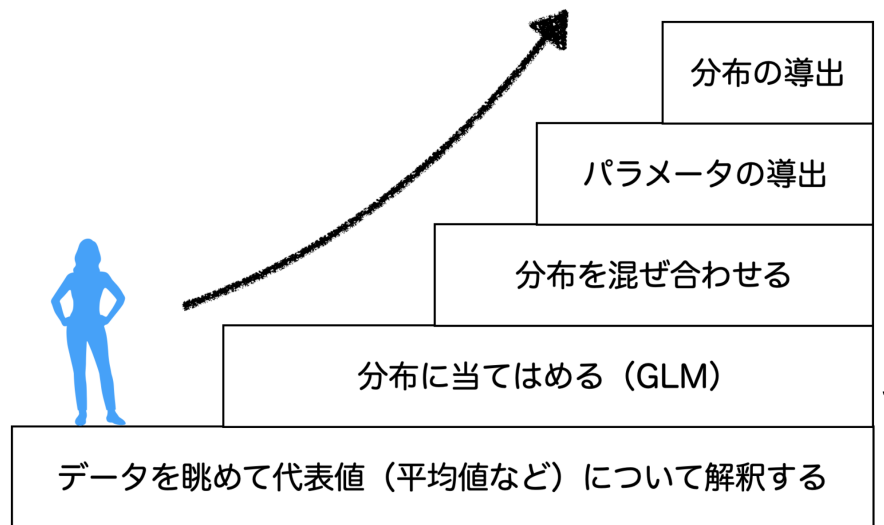


図 27.10 統計モデラーの道。浜田宏先生(東北大学)のアイデアをもとにイラスト化したもの

7655

7656 27.4 課題

7657 以下のモデルを分析する R/Stan コードを提出してください。結果の解釈などを、スクリプトのコメントアウト
7658 や別添ファイルなどで提供してもらえると幸いです。もちろん Rmd ファイルでの提出であれば完璧で
7659 す。なお提出されたコード単体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの
7660 書き方などわからないところがあれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

7661 ■混合分布モデル 今回は阪神 Tigers のデータセットで行ったが、他の球団ではどうだろうか。別の球団
7662 データの年俸(対数をとったもの)のヒストグラムを描き、混合分布モデルが適用できそうなものを見つけた
7663 ら、実際に当てはめて推定してみましょう。

7664 ■ゼロ過剰ポアソン回帰モデル セーブポイントが年俸で説明できるとしたゼロ過剰ポアソン回帰を実行し
7665 てみましょう。年俸データは標準化して単位を整えておくと良いでしょう。

*9 この図は東北大学大学院文学研究科の浜田宏先生が、専修大学社会知性開発研究センター主催の研究会(2019.2.26 開催、
於専修大学サテライトキャンパス)にてご講演いただいた際の、資料をもとに作図したものです。

第 28 章

確率的プログラミング；項目反応理論

さてここまで、線形モデルを中心にモデリングのステップを一段一段登ってきました。この後は、これまでの技術や考え方を応用し、より実践的なプログラミングへと進んでいきます。さまざまなモデルの確率的表現や、そのコーディング技法を学ぶことで、ご自身の研究にも応用できる技術的ヒントが得られるかもしれません。

今回は第 4 講、第 5 講および第 11 講で学んだ項目反応理論を、Stan を使って実装する例を見ていきたいと思います。

28.1 ロジスティックモデルの復習

28.1.1 パラメータモデル側から

ここで簡単に前期の復習をしておきます。

目に見えないものを測定するという意味では、心理尺度や学力検査は同じ技術であるというところから、測定モデルの話を導入しました。とくに学力検査で想定されている測定対象、学力は、正規分布していると考えられること、累積的な能力の蓄積の程度を測定していると考えられることから、正答率が累積正規分布をつかった潜在能力（学力）の関数と考えられるというところからモデル化がすすみました。累積正規分布を直接扱うモデルもありますが、一般にはそれによく近似するロジスティック関数、すなわち

$$f(x) = \frac{1}{1 + \exp(-1.7x)}$$

を使って近似することで、項目の特徴を描写することが考えられたのでした。

x 軸が潜在能力 θ を表していると考えれば、縦軸は能力に応じた通過率になると考えられます。この関数を左右に動かす位置パラメータ、 b_j を加えたモデルが 1 パラメータ・ロジスティックモデル (One Parameter Logistic model) であり、次の式 28.1 で表現されるのでした。

$$p_j(\theta) = \frac{1}{1 + \exp(-1.7(\theta - b_j))} \quad (28.1)$$

ここで $p_j(\theta)$ は項目 j に対する θ の通過率であり、 b_j は困難度 (difficulty) 母数と言われるものです。

さらに項目を特徴づけるために、傾きの大きさを表すパラメータ、 a_j を加えたモデルが 2 パラメータ・ロジスティックモデル (Two Parameter Logistic model) です。モデル式は式 28.2 で表されます。

$$p_j(\theta) = \frac{1}{1 + \exp(-1.7a_j(\theta - b_j))} \quad (28.2)$$

a_j はのことをとくに識別力 (discriminant) と呼ぶのでした。 a_j の値によって傾きがどのように変わる

7689 かを, 図 28.1 の二段目をみて確認しておきましょう。

7690 最後に, 第 4 講では紹介していませんでしたが, 3 つ目のパラメータである c_j , **当て推量母数 (guessing**
 7691 **parameter)** を入れたモデルも紹介しておきます。モデルは式 28.3 のように表されます。

$$p_j(\theta) = c_j + \frac{1 - c_j}{1 + \exp(-1.7a_j(\theta - b_j))} \quad (28.3)$$

7692 ここには切片のように定数 c_j が入っていますから, 関数のベースラインが上に上がるイメージです。これはつ
 7693 まり, 学力 θ にかかわらず一定の正答率を有するというので, 「適当に答えても当たる確率」ということから
 7694 当て推量とよばれています。

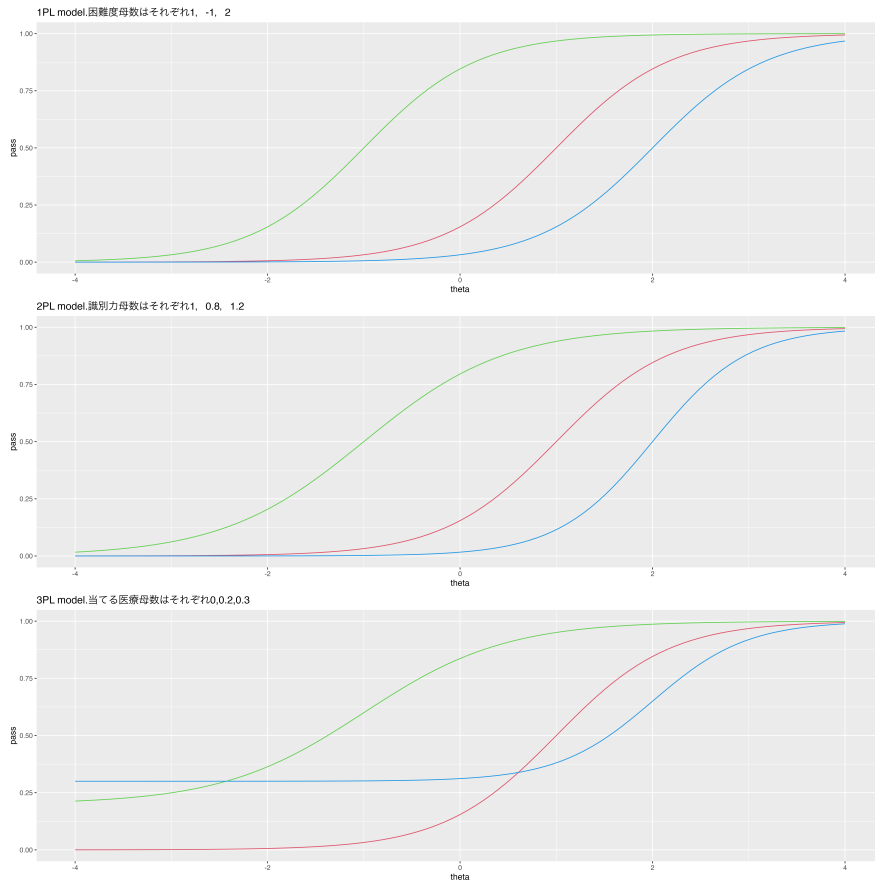


図 28.1 上から順に, 1, 2, 3 パラメータロジスティックモデル。どこがどう変わっていくか確認しておこう。

7695 28.1.2 確率モデル側から

7696 さて今度は同じモデルをデータとの対応の観点から見てみましょう。テストのデータは 0/1 の二値であり,
 7697 ベルヌーイ分布から出てきていると考えることができます。またロジスティックモデルはその言葉の通りロジス
 7698 ティック関数を使って, $-\infty$ から $+\infty$ まであり得る θ の値を, 0 から 1 の範囲に入るように変換している
 7699 わけです。これは一般化線形モデルの文脈でいうところの, **ロジスティック回帰分析**であり, **リンク関数**がロ
 7700 ジット関数, **逆リンク関数**がロジスティック関数になっているモデルだと考えることができます。

7701 データ生成のモデル設計図として考えると, 図 28.2 のようになります。パッケージを使っている分析の時は**最**
 7702 **尤法**による推定結果でしたが, ここで被検者母数に**標準正規分布**を, 項目母数に正規分布を事前分布として

7703 おいだけで、ベイズ推定モデルに置き換えることができました。

これまで利用してきた技術の応用で考えることができますから、すぐにでも実装できそうですね！

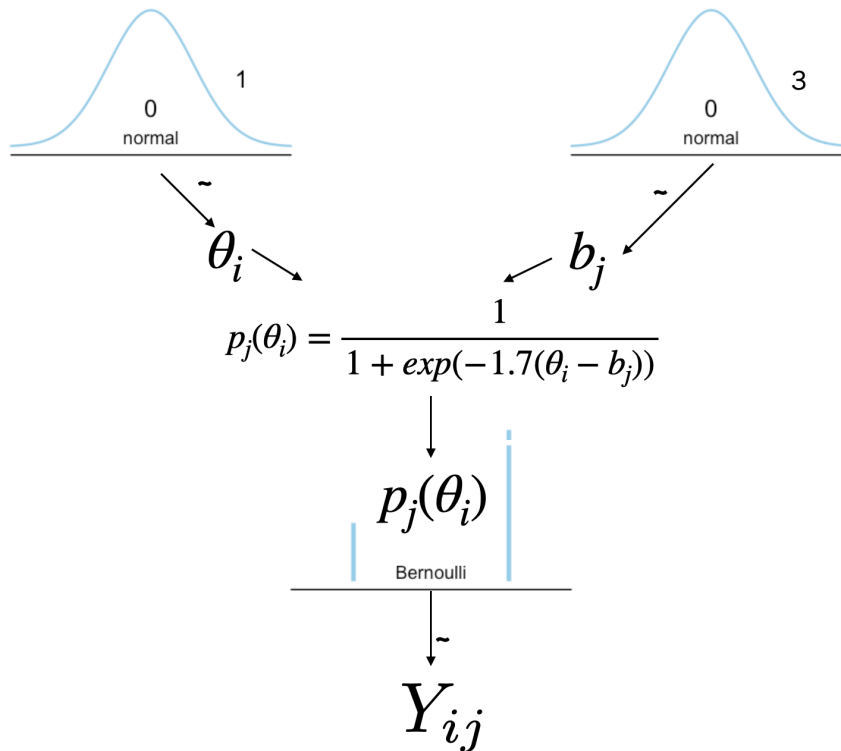


図 28.2 確率モデルとしての設計図。ロジスティック関数とベルヌーイ分布の組み合わせでテスト理論のモデルが表現できる。

7704

28.2 ロジスティックモデルでの実装

7705

7706 それでは設計図をもとにモデルをコードにしていきましょう。

7707

7708 項目数が M 、被検者数が N の $N \times M$ サイズの行列でデータが与えられていると考え、行列の各要素に

7709

7709 対して通過率 p_{ij} を計算し、それをベルヌーイ分布に当てはめるという形で実装したのがコード 28.1 になります。

code : 28.1 1PL モデル

7710

```

7711 1 data{
7712 2   int<lower=0> M;
7713 3   int<lower=0> N;
7714 4   array[N,M] int<lower=0,upper=1> resp;
7715 5 }
7716 6
7717 7 parameters{
7718 8   array[M] real<lower=-5,upper=5> b;
7719 9   array[N] real theta;
7720 10 }
7721 11
7722 12 transformed parameters{

```

```

7723 13   array[N,M] real<lower=0,upper=1> prob;
7724 14   for(n in 1:N){
7725 15     for(m in 1:M){
7726 16       prob[n,m] = inv_logit(1.7*(theta[n]-b[m]));
7727 17     }
7728 18   }
7729 19 }
7730 20
7731 21 model{
7732 22   for(n in 1:N){
7733 23     for(m in 1:M){
7734 24       resp[n,m] ~ bernoulli(prob[n,m]);
7735 25     }
7736 26   }
7737 27   //prior
7738 28   b ~ normal(0,3);
7739 29   theta ~ normal(0,1);
7740 30 }
7741

```

7742 今回はわかりやすくするために、transformed parameters ブロックを使いましたが、直接 model ブロッ
7743 クに書き込んでも構いませんし、bernoulli_logit 関数を使えばさらに高速に安定した推定結果が得ら
7744 れます。

7745 transformed parameters ブロックにパラメータを追加するだけで、2PL,3PL モデルに拡張すること
7746 も容易にできます。2PL モデルに拡張した例がコード 28.2、3PL モデルに拡張した例がコード 28.3 です。

code : 28.2 2PL モデル

```

7747 1  ... (前略)...
7748 2  parameters{
7749 3    array[M] real<lower=0> a;
7750 4    array[M] real<lower=-5,upper=5> b;
7751 5    array[N] real theta;
7752 6  }
7753 7
7754 8  transformed parameters{
7755 9    array[N,M] real<lower=0,upper=1> prob;
7756 10   for(n in 1:N){
7757 11     for(m in 1:M){
7758 12       prob[n,m] = inv_logit(1.7*a[m]*(theta[n]-b[m]));
7759 13     }
7760 14   }
7761 15 }
7762 16 ... (後略)...
7763
7764

```

code : 28.3 3PL モデル

```

7765 1  ... (前略)...
7766 2  parameters{
7767 3    array[M] real<lower=0> a;
7768 4    array[M] real<lower=-5,upper=5> b;
7769 5    array[M] real<lower=0,upper=1> c;
7770 6    array[N] real theta;
7771

```

```

7772 7 }
7773 8
7774 9 transformed parameters{
7775 10   array[N,M] real<lower=0,upper=1> prob;
7776 11   for(n in 1:N){
7777 12     for(m in 1:M){
7778 13       prob[n,m] = c[m] + (1-c[m])*inv_logit(1.7*a[m]*(theta[n]-b[m]));
7779 14     }
7780 15   }
7781 16 }
7782 17 ... (後略)...
7783

```

7784 第 11 講で使ったサンプルコードをつかって、このモデルで推定するためのコードがコード 28.4 です
7785 (cmdstanr での例)。

code : 28.4 IRT のコード

```

7786 1 dat <- read_csv("IRTsample.csv")
7787 2 model_1pl <- cmdstan_model("oneParameter.stan")
7788 3 model_2pl <- cmdstan_model("twoParameters.stan")
7789 4 model_3pl <- cmdstan_model("threeParameters.stan")
7790 5
7791 6 dataSet <- list(N = NROW(dat), M = NCOL(dat), resp = as.matrix(dat))
7792 7 fit1 <- model_1pl$sample(data = dataSet, chains = 4, parallel_chains = 4)
7793 8 fit2 <- model_2pl$sample(data = dataSet, chains = 4, parallel_chains = 4)
7794 9 fit3 <- model_3pl$sample(data = dataSet, chains = 4, parallel_chains = 4)
7795
7796

```

分析結果を見てみましょう。EAP 推定値はそれぞれ表 28.1, 図 28.3 のようになりました。

表 28.1 それぞれのパラメータ推定値

Qid	1PL		2PL		3PL		
	b_j	a_j	b_j	a_j	b_j	c_j	
1	0.727	0.347	1.716	2.865	1.567	0.221	
2	-1.605	0.724	-2.140	1.818	-0.761	0.557	
3	-1.175	1.124	-1.229	2.900	-0.626	0.338	
4	-1.077	0.955	-1.222	2.033	-0.535	0.362	
5	0.577	0.691	0.771	2.451	0.946	0.159	
6	1.271	0.694	1.742	1.100	1.659	0.050	
7	0.564	0.692	0.753	1.027	0.941	0.099	
8	1.590	0.360	3.650	1.581	3.355	0.084	
9	0.514	0.484	0.893	0.836	1.205	0.142	
10	-0.864	1.214	-0.887	2.547	-0.395	0.298	

7797

7798 28.3 整然データでの分析

7799 このように、IRT モデルを簡単に実装できました。

7800 ところで、このモデルをもう少し使いやすくするために、データを整然データ (tidy data) にすることを考

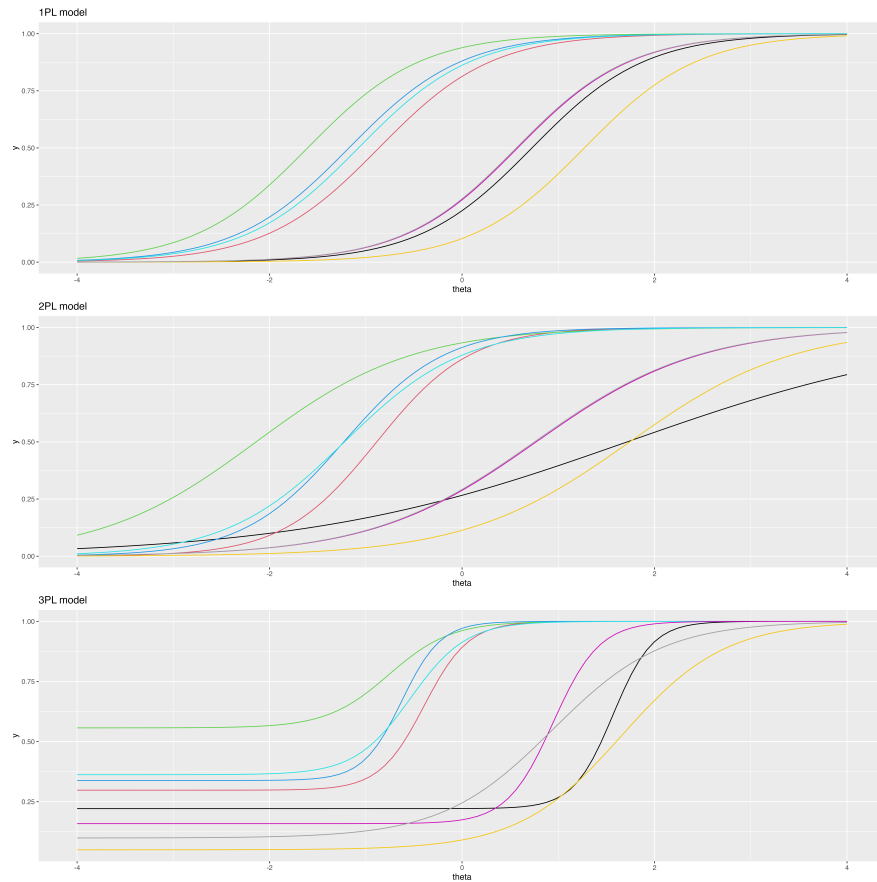


図 28.3 推定値を使って描いた各モデルの ICC

7801 えてみましょう。というのも、今はデータを行列形で渡していますが、これだと欠損値があった場合にうまく機
 7802 能しないからです。Stan はデータに欠損値を取ることができず、与えられるデータに NA が入っていることは
 7803 許されません。

7804 整然データは、1 行に 1 つの数字が入っているようなデータで (セクション 22.3.1, Pp.247 参照), その行
 7805 を見るだけですべての情報が手に入るようなデータです。行列型のデータは、人の ID を行ラベルで、変数の
 7806 ID を列ラベルでみなければなりませんので、参照方向が 1 つに定められていません。具体的には次のような
 7807 形のデータになります (R の出力 28.1)。この形式ですと、一行で誰の (Pid), どの問いに対する (Qid) 反応
 7808 か (value) ということがわかりますね。

R の出力 28.1: 整然データになったテストデータ

```
# A tibble: 20 × 3
  Pid  Qid value
  <int> <dbl> <dbl>
1     1     1     0
2     1     2     0
3     1     3     1
4     1     4     0
5     1     5     0
6     1     6     0
7     1     7     0
8     1     8     0
9     1     9     0
10    1    10     0
11    2     1     1
12    2     2     1
13    2     3     0
```

7809

もし欠損値があればその行を削除してしまえば、完全データになりますから、Stan に欠損値を与えなくて済むことになります。

もちろんこれに対応する形で、Stan のコードを書き換える必要があります。こんどはデータも長くなりますから、何行目に誰のど問いに対する反応が入っているかを、識別変数を使いながら指定することになります。具体的には次のようなコード例になるでしょう (コード 28.5)。

code : 28.5 Tidy Data に対応した 2PL モデル

7815

```
7816 1 data{
7817 2   int<lower=0> L;
7818 3   int<lower=0> N;
7819 4   int<lower=0> M;
7820 5   array[L] int<lower=0,upper=N> Pid;
7821 6   array[L] int<lower=0,upper=M> Qid;
7822 7   array[L] int<lower=0,upper=1> resp;
7823 8 }
7824 9
7825 10 parameters{
7826 11   array[M] real<lower=0> a;
7827 12   array[M] real<lower=-5,upper=5> b;
7828 13   array[N] real theta;
7829 14 }
7830 15
7831 16 transformed parameters{
7832 17   array[N,M] real<lower=0,upper=1> prob;
7833 18   for(n in 1:N){
7834 19     for(m in 1:M){
7835 20       prob[n,m] = inv_logit(1.7*a[m]*(theta[n]-b[m]));
7836 21     }
7837 22   }
7838 23 }
7839 24
```



```

7840 25 model{
7841 26   for(l in 1:L){
7842 27     resp[l] ~ bernoulli(prob[Pid[l],Qid[l]]);
7843 28   }
7844 29   //prior
7845 30   a ~ normal(0,3);
7846 31   b ~ normal(0,3);
7847 32   theta ~ normal(0,1);
7848 33 }
7849

```

7850 ■コード解説

7851 data ブロック データ長 L, 最大被検者数 N, 最大項目数 M, 被検者識別変数 Pid, 項目識別変数 Qid, L
7852 行目の反応 resp としてデータを受け取ります。

7853 parameters ブロック これまでと同じです。

7854 transformed parameters ブロック ここの手を加える必要がありません。

7855 model ブロック L 行目の反応に対して, 識別変数で特定された確率をつかってモデルを当てはめます。

7856 この結果はこれまでのモデルと変わりませんが, 欠損値に対応できたはずですので, 元の完全データに技
7857 と欠損値を与えて改変し, 結果がどう変わるかを確認してみましょう。

code : 28.6 tidy データにして技と欠損値を与える

```

7858 1 dat.tmp <- dat %>%
7859 2   rowid_to_column("Pid") %>%
7860 3   pivot_longer(-Pid) %>%
7861 4   mutate(Qid = str_extract(name, pattern = "\\d+") %>% as.numeric()) %>%
7862 5   dplyr::select(Pid, Qid, value)
7863 6
7864 7 # わざと欠損値を与える
7865 8 dat.tmp$value[1] <- NA
7866 9 dat.tmp$value[11:13] <- NA
7867
7868

```

7869 ■コード解説

7870 1 行目 dat に入っている行列型・完全データを変形していきます。

7871 2 行目 一行ごとに被検者の反応が入っていますので, 行番号を被検者識別変数 Pid として作ります。

7872 3 行目 データを縦長にします。その時のキーとなるのは, 先ほど作った Pid で, この変数は除いて縦長にする関数が pivot_longer です。

7874 4 行目 少し技巧的ですが, ここまでの段階で変数は Pid, name, value という 3 つになっています。なか
7875 でも name 変数には元データの変数名が入っていますので, これを加工して問題番号を取り出しま
7876 す。str_extract 関数の str とは string, すなわち文字列を扱う関数であるという意味です。
7877 str_extract は文字列から条件に合ったものを抜き出す extract というもので, 条件を pattern
7878 で指定しています。ここで \d+ とあるのは正規表現 (regular expression) というもので, 文字列
7879 を一定の規則に従って特定の文字列を表現する方法です。ここでは数字だけを抜き出しています*1。

*1 少し丁寧にいうと, \d は正規表現で数字を意味する記号です。ただ, バックスラッシュ (\) がそのままでは正規表現の記号だと認識されず特殊文字だという宣言をしている (エスケープシーケンス, といいます) だけになってしまいます。そこで, 「特殊文字

7880 また、このようにして取り出せた数字は文字列としての数字ですので、これを `as.numeric` 関数に送
7881 ることで、数値であることを教えています。

7882 5 行目 被検者識別変数, 項目識別変数, 反応の値だけにデータを限定しています。

7883 7 行目 ここでこのデータセットの 1 行目に欠損値 NA を上書きしています。わざと欠損させたのです。

7884 8 行目 同じく 11 行目から 13 行目までも欠損値に上書きしてしまいました。

7885 こうしてできたデータセットは次のようになります (R の出力 28.2)。

R の出力 28.2: 整然データになったテストデータ

```
# A tibble: 20 × 3
  Pid  Qid value
<int> <dbl> <dbl>
1     1     1   NA
2     1     2     0
3     1     3     1
4     1     4     0
5     1     5     0
6     1     6     0
7     1     7     0
8     1     8     0
9     1     9     0
10    1    10     0
11    2     1   NA
12    2     2   NA
13    2     3   NA
14    2     4     0
15    2     5     0
16    2     6     0
17    2     7     0
18    2     8     0
19    2     9     0
20    2    10     1
```

7886

7887 確かに数カ所、欠けているところがありましたね。これをそのまま渡すことができませんから、欠損のあると
7888 ころは削除してデータセットを作り、推定してみましょう。

code : 28.7 データセットを作って推定する

7889

```
7890 1 # 欠損値を消したデータセットにする
7891 2 dat.tmp <- na.omit(dat.tmp)
7892 3
7893 4 dataSet <- list(
7894 5   L = NROW(dat.tmp), N = max(dat.tmp$Pid), M = max(dat.tmp$Qid),
7895 6   Pid = dat.tmp$Pid, Qid = dat.tmp$Qid,
7896 7   resp = dat.tmp$value
7897 8 )
```

のバックスラッシュだよ」を表すために重ねて\\としています。また、+ は正規表現でいうところの「直前の文字が 1 回以上繰り返される」という意味です。今回は 10 という数字が出てくる可能性があり、\\d だけだと 1 しか抜き出せないの、数字が連なっていたら全体を取り出すようにしています。要するに、数字を全部取り出しましょう、が正規表現では\\d+ になるわけです。

```

7898 9
7899 10 model_2pl_ver2 <- cmdstan_model("twoParameters2.stan")
7900 11 fit2.2 <- model_2pl_ver2$sample(
7901 12   data = dataSet,
7902 13   chains = 4,
7903 14   parallel_chains = 4
7904 15 )
7905

```

7906 欠損があっても、データを整形して欠損を取り除いた形で分析し、問題なく推定できたと思います。欠損が含まれているから、そのデータはすべて使い物にならないと考えるのではなく、データの存在するところ・利用できるところは利用しつくすという有効活用ができたと思います。

7907
7908
7909 最後に、推定結果を確認しておきましょう。1 人目の被検者は 1 つ、2 人目の被検者は 3 つの欠損がありました。つまり他の被検者は 10 問分の情報を持っているのに、この 2 人はそれぞれ 9 問、7 問分しか情報が得られなかったこととなります。そのことが結果の推定値にどう変わるのかというと、出力 12 の通りです。

MCMC の結果 12

```

# A tibble: 3 × 7
  name          EAP      MED      MAP      SD      L95      U95
<chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 theta[1]  -1.403  -1.392  -1.438  0.548  -2.530  -0.366
2 theta[2]  -0.774  -0.778  -0.809  0.606  -1.982   0.372
3 theta[3]  -0.656  -0.659  -0.620  0.508  -1.654   0.331

```

7912
7913 これを見ると、10 問分の情報を持っている 3 人目の被検者の能力値が $\theta_3 = -0.656$ と推定されていますが、その SD が 0.508 です。これに比べて、9 問しか情報のない被検者 1 の SD は 0.548、7 問しか情報のない被検者 2 の SD は 0.606 と、情報が少なくなると SD が大きくなっていくことがわかります。推定値の SD が大きいということは、幅が広い、すなわちわからないことがより多くあることを意味します。得られる情報が少なければ、絞り込みが難しくなるというのがデータにも表れていることがわかりますね。

28.4 課題

7918
7919 本講で扱った、1PL, 2PL, 3PL モデルそれぞれを実行する、Stan ファイルや R コードを提出してください。ただし、いずれのモデルも整然データに対応したコードになっている必要があります。データや R コードはシラバスのサイトを通じて提供されています。不明な点がありましたら、TA あるいは小杉まで連絡して指導を受けてください。

第 29 章

確率的プログラミング；変化点と折線 回帰

さて今回は、今までの線形モデルとはちょっと違うモデリングになります。その名も変化点検出、そして折線回帰です。タイトルだけでも面白そうではありませんか？

さらに今回は次のようなデータを扱います (図 29.1)。何を隠そうこのデータは、私の体重の推移のデータなのです。私のモーニング・ルーティンとして、朝目が覚めるとトイレに行きます。そして体重計にのり、体重と体脂肪を測定します。その後顔を洗って*1、計測値を iPhone のアプリに書き込むというがあります*2。このルーチンも早いもので 10 年近く続けていることになります (一時期やめていたこともありましたが、最古の記録は 2012 年 06 月 08 日です)。ともかく、このデータを使って、いろいろ遊んでみようと思います。

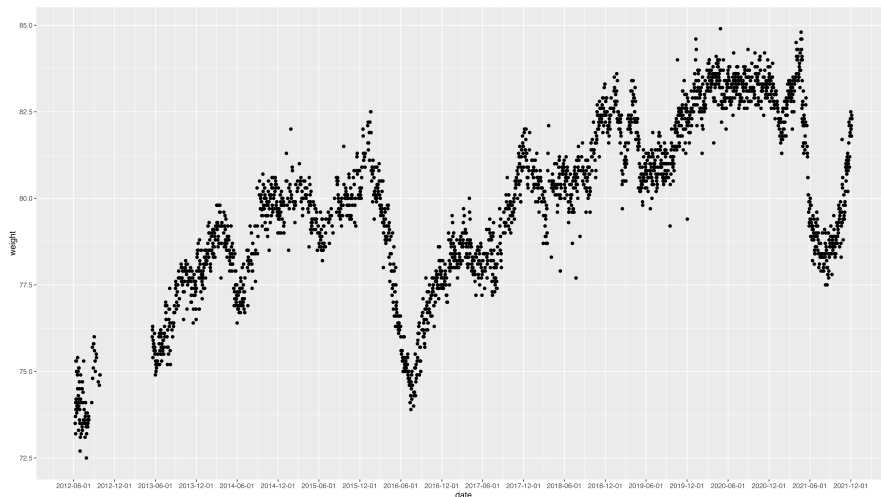


図 29.1 体重の時系列的推移

*1 体の内側の水分を出し、表皮に水分をつけないようにすることで、しっかりと肉の量を測ろうとしているのでこの順番になります。

*2 ついでにいうと、書き込んだ結果はツイートします。フォロワーの中には私の体重の推移を楽しみにしている人がいるのです。おかしな世の中です。

7933

29.1 混合分布モデルの応用

7934 データが 10 年分というのはちょっと多いですから、少し時期を区切ることにします。2019 年から 2020 年の 2 年間に限定したデータが、図 29.2 になります。

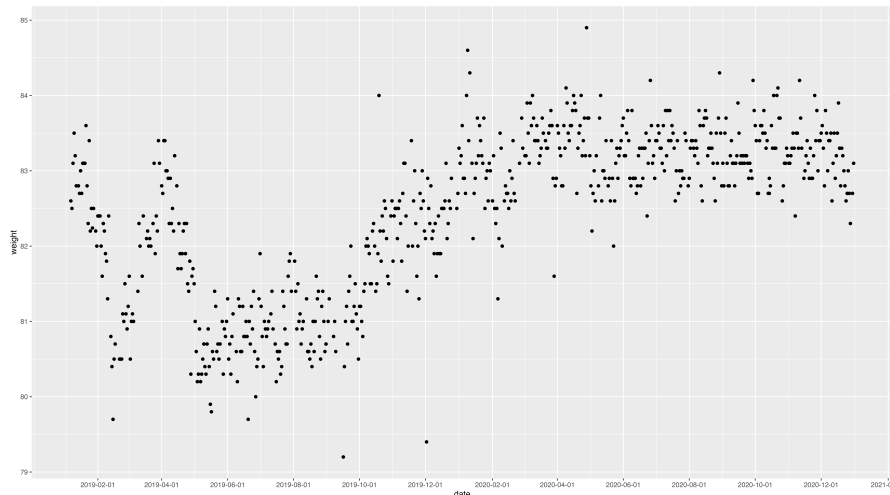


図 29.2 2019-2020 年のデータ

7935 これをみると、うーん残念なことに、2020 年になると体重が増えてしまっているようですね。2019 年は
7936 81kg 台をうろろうしていたようですが、2020 年になると 83kg 台になっているようです。毎回の体重の計測
7937 に偶然的な測定誤差がついているとしても、19 年と 20 年とではそもそも体重が違っているようです。
7938 これは平均値の違う正規分布が混合している、混合分布モデルを考えてみることができそうですね。

7940 第 27 講を参考に、混合分布モデルを考えてみましょう。筆者の体重には 2 つの状態があり、軽い方の状態
7941 なのか重い方の状態なのかは、確率 θ で変わると考えてみます。するとデータは、確率 θ で μ_1 を平均とする
7942 正規分布から、確率 $1 - \theta$ で μ_2 を平均とする正規分布から得られるわけですから、確率モデルは次のよう
7943 になります。

$$p(W_j) = \theta_j \times N(\mu_1, \sigma) + (1 - \theta_j) \times N(\mu_2, \sigma)$$

7944 ここで θ_j としたのは、観測時点 j ごとに θ が変わるという想定です。どちらのモードになるのかは一定の
7945 数字というより、毎回違っているように思えたのでそのようにしてみました。されてこを Stan で実装するには
7946 どうすれば良いのでしょうか。そう、log_sum_exp 関数を使って、2 つの状態をベクトルで表記し、それを足
7947 し合わせる必要があるんでしたね。実際にコードにしてみたのがコード 29.1 です。

code : 29.1 2 つの体重モードモデル

```
7948 1 data{
7949 2   int L;
7950 3   array[L] real W;
7951 4 }
7952 5
7953 6 parameters{
7954 7   array[L] real<lower=0,upper=1> theta;
```

```

7956 8   ordered[2] mu;
7957 9   real<lower=0> sigma;
7958 10  }
7959 11
7960 12 model{
7961 13   for(l in 1:L){
7962 14     target += log_sum_exp(
7963 15       log(theta[l]) + normal_lpdf(W[l]|mu[1],sigma),
7964 16       log1m(theta[l]) + normal_lpdf(W[l]|mu[2],sigma)
7965 17     );
7966 18   }
7967 19
7968 20   mu ~ normal(80,10);
7969 21   sigma ~ cauchy(0,5);
7970 22 }
7971

```

7972 これを使って、体重のデータを分析してみましょう。コード 29.2 を実行し、図 29.3 のような結果を得ます。

code : 29.2 混合分布モデルのコード

```

7973 1 dat1 <- dat %>%
7974 2   dplyr::filter(date > "2019/01/01") %>%
7975 3   dplyr::filter(date < "2021/01/01")
7976 4
7977 5 model <- cmdstanr::cmdstan_model("changePoint1.stan")
7978 6 dataSet <- list(L = NROW(dat1), W = dat1$weight)
7979 7 fit <- model$sample(
7980 8   data = dataSet,
7981 9   chains = 4,
7982 10  parallel_chains = 4,
7983 11  seed = 12345)
7984
7985

```

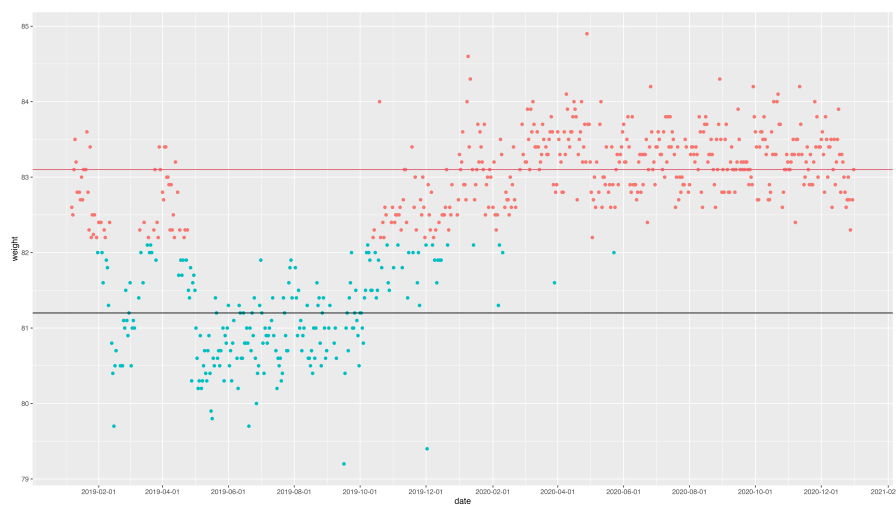


図 29.3 2つの体重のモード

7986 結果の図を見ると、どうやら 2019 年の初期にも 2 回ぐらい平均値が高い、「重いモード」が混在してお

7987 り*3, 2019 年の 10 月ごろからそちらの比率が増えています。もっとも, 2020 年の 1 月にも「軽いモード」に
7988 入っている点はあるようですね。

7989 とまあ, このような分析結果になったわけですが, これをみるとデータが 82kg ぐらいのラインを超えたかど
7990 うかで分割されているな, というのがわかります。当然, 平均値の違う 2 つの分布を混ぜたわけですから, 平
7991 均値の違いによって分かれるわけです。それにしても 2019 年の初期はどうしたんでしょうね。1 月や 4 月は
7992 重いモード, 2-3 月と 5 月以降は軽いモードと, コロコロ入れ替わっています。ここでは 2 つのモードのどちら
7993 からデータが出てきているのか, という事だけを考えているのでこれでいいのですが, 体重というのはそも
7994 そも時系列的な変化をするものですから, どこかで重くなった, どこかで軽くなった, というような時系列的な
7995 つながりについての情報が, うまくモデル化されていないように思えます。そこで, 横軸が時間的な連続であ
7996 るということに注意して, 今度は 2021 年 1 月から 11 月までのデータを見てみましょう。その区間を取り出し
7997 たのが図 29.4 になります。

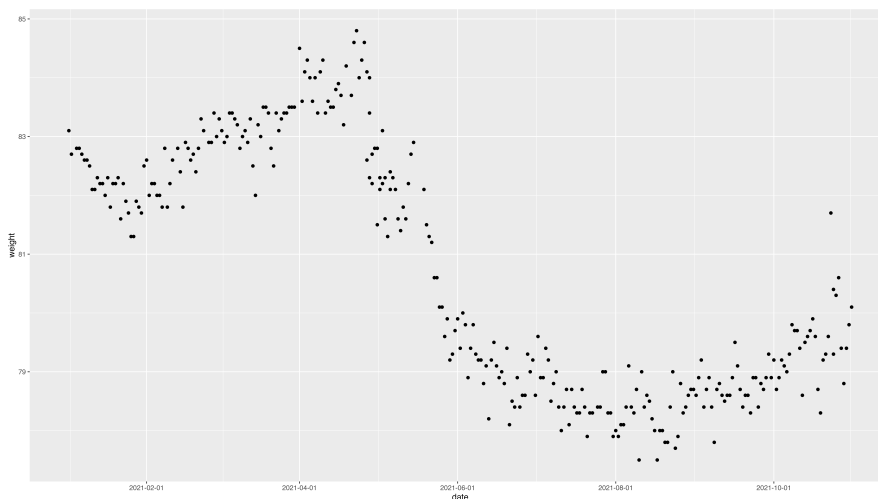


図 29.4 大きく変化したデータ

7998 これは 2 つのモードがあるというより, 途中で大きく変化したというべきではないでしょうか。とくに 2021 年
7999 5 月ごろから 4kg ほど体重がぐぐぐと下がることがあり, 6 月以降は下がった体重のレベルに落ち着いて
8000 いるようです。このように, 横軸が時系列だと考えると, 5 月に何か変化があったのではないかと, ことが
8001 推察されます。これをモデルで表現してみましょう。

8002 29.2 変化点検出

8003 何かのきっかけでデータの様相がガラリと変わってしまった, その変化したところを変化点と呼び, ここで
8004 課題はその変化点を見つけ出す**変化点検出 (Change point detection)**ということになります。たとえ
8005 ばセンサーが検出するデータの針が急に違うレベルに変化すると, 測定している対象の状態がガラリと変わ
8006 たのではないかと, 考えることができますね。

8007 他にもたとえば, 心理学の応用領域として, 科学捜査研究所が担当する**ポリグラフ検査**というのがありま
8008 す*4が, これなども針が大きく触れたことで変化をみることになります。

*3 誰がデブモードやねん

*4 皮膚電気活動や心拍, 呼吸など複数の整理指標を同時に測定するので**ポリグラフ**であり, 熟練の検査官が反応パターンから被
疑者の特別な反応を検出するものです。嘘発見器と呼ばれることがありますが, 正確には**虚偽検出**と読んだほうが良いでしょう。

8009 変化点検出はその名の通り「いつ」変化したのかを検出するものです。今回のデータも、だいたい5月ご
8010 ろに大きな変化があったというのはわかるのですが、いつなのかを特定するのは難しいところ*5。これを
8011 データとモデルから明らかにしようというのです。

8012 ベイズ統計はわからないことを確率で表現し、データでその確率情報をアップデートしていくというもので
8013 す。今回はいつ変化したのがわかりませんから、その変化した時期を τ とし、データ区間の中にその日がある
8014 と考えて、その区間の一様分布を事前分布とします。その時期 τ を境に、データが出てくる分布の位置パラ
メータがずれると考えるのです。設計図は次のようになります (図 29.5)

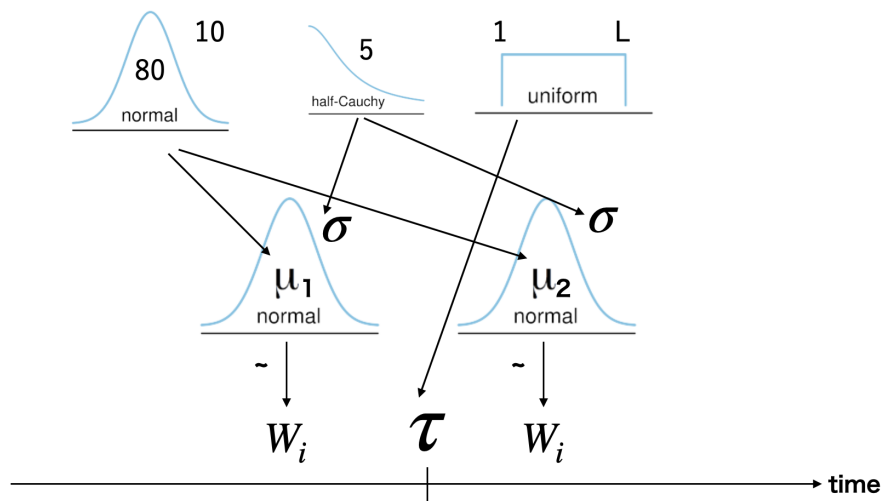


図 29.5 変化点検出のモデル設計図

8015 設計図をもとに、コードにしてみたのがコード 29.3 になります。

code : 29.3 変化点検出のモデルコード

```

8017 1 data{
8018 2   int L;
8019 3   array[L] real W;
8020 4 }
8021 5
8022 6 parameters{
8023 7   real<lower=1,upper=L> tau;
8024 8   ordered[2] mu;
8025 9   real<lower=0> sigma;
8026 10 }
8027 11
8028 12 model{
8029 13   for(l in 1:L){
8030 14     if(l < tau){
8031 15       W[l] ~ normal(mu[2],sigma);
8032 16     }
8033 17   }

```

最近では隠匿情報検査 (concealed information test) ということもあります。ちなみに科学捜査研究所、通称科捜研は、心理の他にも物理、化学、法医、文書 (筆跡鑑定) などの専門領域があり、各都道府県に1つずつ設置されています。

*5 とはいえこのデータは筆者自身の体重変化ですので、何があったのかは実はわかっています。5月11日、筆者が帰宅途中に自転車で転倒する事故を起こし、右肩の鎖骨を骨折しました。17日に手術した後、利き手が使いにくいものですから食事の量が減り、結果的に体重 (筋肉?) が落ちたというのが真相です。

```

8033 16     }else{
8034 17         W[1] ~ normal(mu[1],sigma);
8035 18     }
8036 19 }
8037 20
8038 21     tau ~ uniform(1,L);
8039 22     mu ~ normal(80,10);
8040 23     sigma ~ cauchy(0,5);
8041 24 }
8042

```

8043 日々のデータが順に 1 から L 行目まで並んでいるとします。パラメータ tau があり、1 行目のデータが変
8044 化点 tau より前にあるときは μ_2, σ の正規分布から、tau より後になれば、 μ_1, σ の正規分布からデータが
8045 出てくると考えるのです。今回は後半の平均値が小さいことが明らかですから、 μ のベクトルを ordered 型
8046 で宣言してあります。ベクトル要素の小さい順に並びますから、前半が μ_2 、後半が μ_1 としています。
8047 これを実行して、変化点がいつになるのかを推定してみましょう。推定結果は出力 13 のようになりました。

MCMC の結果 13

```

# A tibble: 4 × 7
  name      EAP      MED      MAP      SD      L95      U95
  <chr> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!> <num:.3!>
1 mu[1]  82.770    82.769    82.762    0.062    82.648    82.893
2 mu[2]  78.891    78.892    78.896    0.060    78.775    79.008
3 sigma  0.750     0.749     0.747    0.031     0.693     0.812
4 tau   144.642   144.578   144.431    0.646   143.321   146.190

```

8048
8049 データの 145 行目がちょうどその変化点だといって、ほぼ間違いないようですね。145 行目はといいます
8050 と、データから 5 月 23 日だったことがわかります。

R の出力 29.1: いつでしょう

```

> dat2[145, ]
# A tibble: 1 × 3
  date          weight bodyFat
  <dtm>          <dbl>  <dbl>
1 2021-05-23 06:50:57  80.6    26

```

8051
8052 違いがわかるようにデータとモデルの推定値を合わせたのが、図 29.6 です。このようにして、ここで様相
8053 が変わったのだなということを、データから見出すことができました。

29.3 折線回帰

8054
8055 先ほどは、変化点を機に平均値が変わる、というモデルでした。では次のような分布の場合はどうなるで
8056 しょうか。今度は 2016 年に注目してみました (図 29.7)。この年はダイエットを志した時期で、最初の半年ぐ
8057 らいでぐんぐん体重が落ちていってるのがわかります。そして残念なことに、8 月ごろでしょうか、ダイエットが
8058 終わりリバウンドが始まったのがよくみて取れますね*6。

*6 このときのダイエット法は、1. ラーメンのスープは飲まない、2. 夕食と晩酌は分ける (食べながら飲むのではなく、食べ終わって、お酒だけ飲む)、3. お酒を飲むときにおつまみは食べない、という 3 か条を守るというものでした。お酒が飲みたいので食事の量

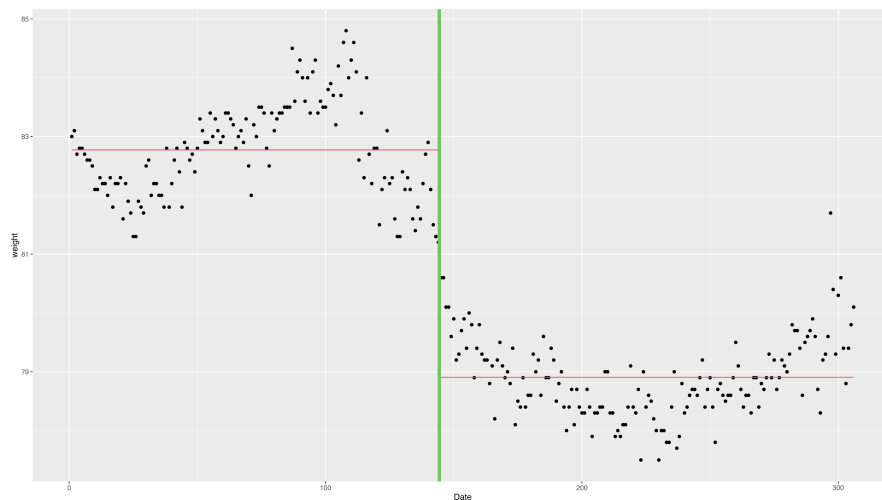


図 29.6 検出された変化点と平均値の違い

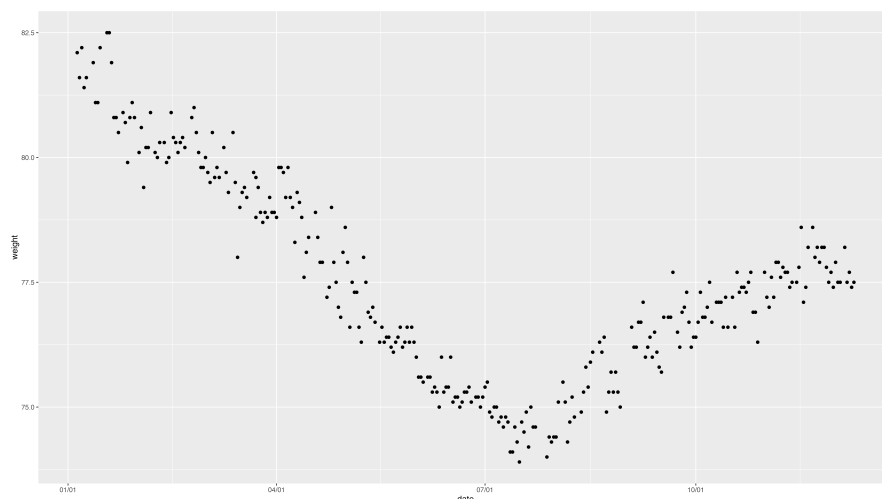


図 29.7 2016 年の変化

8059 さてこのデータに対して、先ほどの変化点検出もできるのですが、変化点の前後で平均値が違う、というよ
 8060 うな簡単な話ではなさそうです。変化点の前はどんどん数字が減っていき、変化点のあとはどんどん数字が
 8061 大きくなっていく、というのが実態です。回帰分析をしたら、変化点前は負の傾きが、変化点後は正の傾きが
 8062 推定されそうな、そんなデータになっていますね。ということで、そのイメージができたのであれば、そのままモ
 8063 デルを描いてしましましょう。設計図が少しごちゃっとしてしまいましたが、ポイントは傾きの係数を前半は負、
 8064 後半は正に限定しているところでしょうか (図 29.8)。また、変化点は大体このあたり・・・ということでデータの
 8065 100 行目 (2016/04/29) から 250 行目 (2016/10/15) の間に置いてみました。

8066 設計図をもとに、コードにしてみたのがコード 29.4 になります。

code : 29.4 折線回帰のモデルコード

8067

を減らし、すぐに飲むことにしたので体重が減り、食事の前に飲んでしまえばいいんじゃないかという裏技を見つけてダイエットが崩壊したのを覚えています。

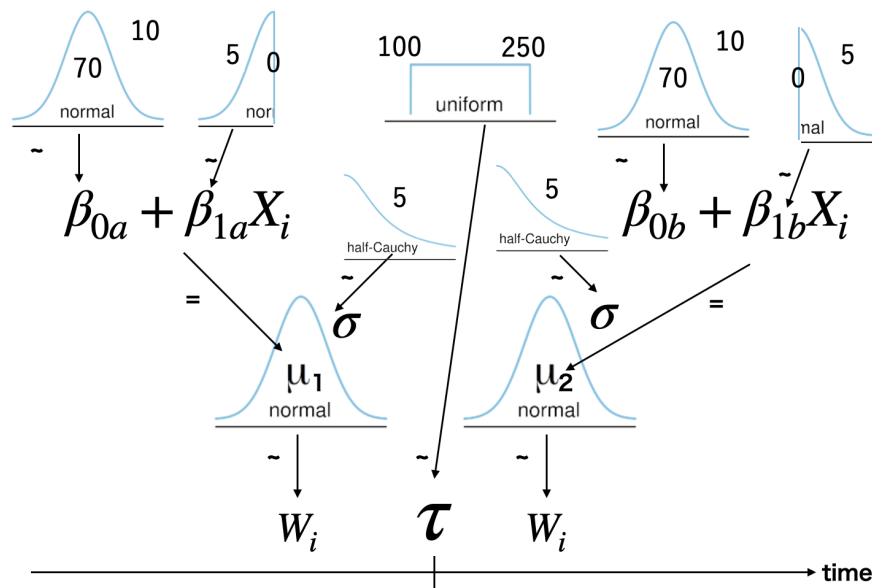


図 29.8 折線回帰モデル設計図

```

8068 1 data{
8069 2   int L; // data length
8070 3   array[L] real W;
8071 4   array[L] real X;
8072 5 }
8073 6
8074 7 parameters{
8075 8   real<lower=100,upper=250> tau;
8076 9   array[2] real beta0;
8077 10  real<upper=0> beta1a;
8078 11  real<lower=0> beta1b;
8079 12  real<lower=0> sigma;
8080 13 }
8081 14
8082 15
8083 16 model{
8084 17   for(l in 1:L){
8085 18     if( l < tau ){
8086 19       W[l] ~ normal( beta0[1] + (beta1a * X[l]),sigma);
8087 20     }else{
8088 21       W[l] ~ normal( beta0[2] + (beta1b * X[l]),sigma);
8089 22     }
8090 23   }
8091 24
8092 25   beta0 ~ normal(70,10);
8093 26   beta1a ~ normal(0,5);
8094 27   beta1b ~ normal(0,5);
8095 28   sigma ~ cauchy(0,5);
8096 29 }

```

8097

8098 回帰係数を正・負に限定するところは、パラメータの宣言で<lower=0>や<upper=0>としていることで表現し
 8099 ています。あとはデータが出てくる正規分布の平均に、線形モデルが入っているだけです。これを実行する
 8100 と、出力 14 のような結果が得られます。データの 187 行目 (2016/08/01) あたりが変化点ですかね。95%
 8101 区間で言うと 185 行目 (2016/07/30) から 188 行目 (2016/08/02) の間に変化点があると言えそうです。

MCMC の結果 14

A tibble: 6 × 7

name	EAP	MED	MAP	SD	L95	U95
<chr>	<num: .3!>	<num: .3!>	<num: .3!>	<num: .3!>	<num: .3!>	<num: .3!>
1 beta0[1]	81.871	81.872	81.874	0.075	81.721	82.014
2 beta0[2]	70.548	70.551	70.549	0.380	69.793	71.274
3 beta1a	-0.043	-0.043	-0.043	0.001	-0.044	-0.041
4 beta1b	0.026	0.026	0.026	0.002	0.023	0.029
5 sigma	0.499	0.498	0.497	0.021	0.461	0.543
6 tau	187.135	187.326	187.460	0.832	185.057	188.215

8102

8103 そして回帰係数をプロットしてみました (図 29.9)。なかなかデータにフィットしていそうです。

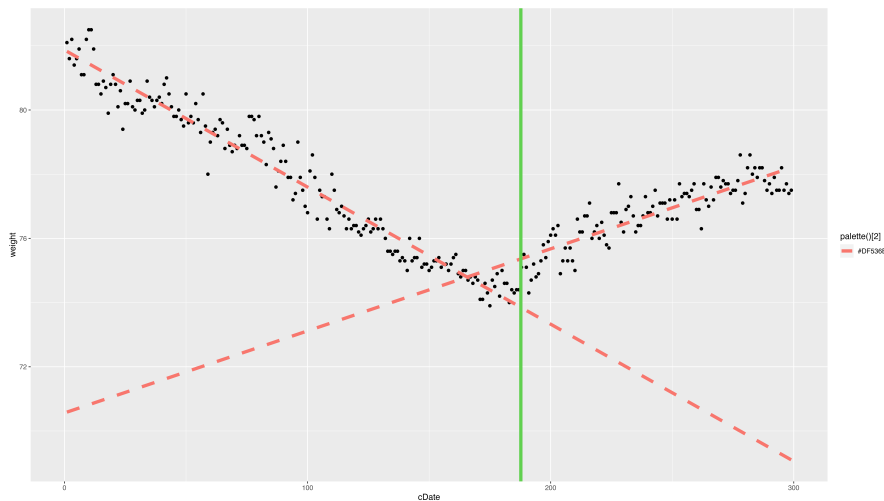


図 29.9 折線回帰モデル推定結果

8104 このままでもよいのですが、変化点と折れたポイントが合致しないのはなんだか気持ち悪いですね。これを
 8105 合わせることを考えてみたいと思います。変化点 τ のあるところで 2 つの回帰線は交わる、つまり同じ値にな
 8106 るはずですから、 $\beta_{0a} + \beta_{1a}\tau = \beta_{0b} + \beta_{1b}\tau$ という式が成り立つはず。ここから逆算して、

$$\beta_{0a} + \beta_{1a}\tau - \beta_{1b}\tau = \beta_{0b}$$

8107 と考えることができますから、推定するパラメータを 1 つ減らすことができます。

code : 29.5 折線回帰のモデルコード 2

8108

8109 1 ... (前略) ...

8110 2 parameters{

8111 3 real<lower=100,upper=250> tau;

```

8112 4   real beta0a;
8113 5   real<upper=0> beta1a;
8114 6   real<lower=0> beta1b;
8115 7   real<lower=0> sigma;
8116 8   }
8117 9
8118 10  transformed parameters{
8119 11   real beta0b;
8120 12   beta0b = beta0a + ((beta1a-beta1b) * tau);
8121 13  }
8122 14  ... (後略)...
8123

```

8124 このコード 29.4 にあるように、beta0b を計算式から出すようにして、改めて推定をした結果が図 29.10 に
 8125 なります。これだと変化点が前にずれて、データの 169 行目 (2016/07/10) で心がポッキリ折れていることが
 わかります。

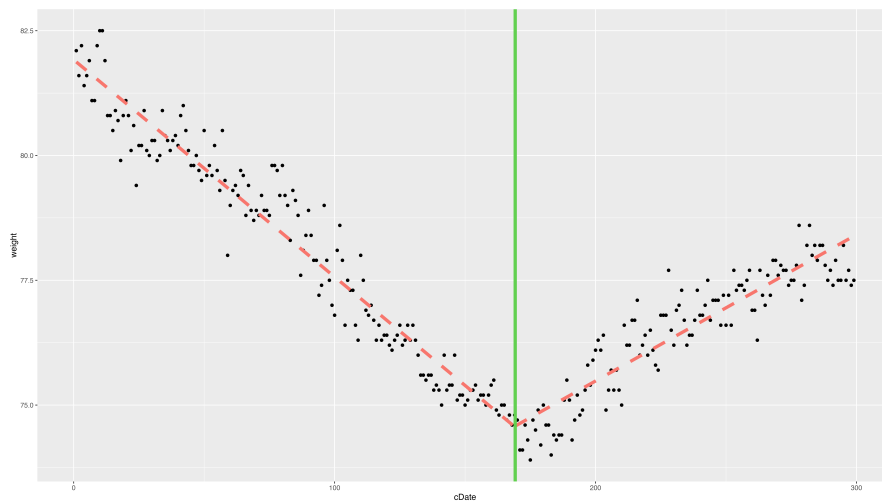


図 29.10 変化点を合わせた折線回帰モデル

8126
 8127 このように、モデリングの力を借りるとわからなかった変化点がどこにあるのか、そしてその変化点の前後
 8128 で予測モデルを変えるとといったようなことが簡単に表現できます。ただの回帰分析であっても、データにあっ
 8129 たモデルを考えることでさまざまな応用可能性が出てくるのではないのでしょうか。

8130 ところで、今回は時系列的なデータに対して回帰分析を行いました。この時の説明変数 X_i は日付、あるい
 8131 は「何日目か」というデータだったわけです。しかし回帰分析には重要な仮定として、各データ点は (同じ分布
 8132 から) 独立に得られていると言うものがありました。たとえば身長と体重の回帰分析、という話をするとき、
 8133 各データ点にあたる個々人の身長・体重は、互いに影響し合わない独立だったはずであり、だからこそ尤度の
 8134 計算の時に各データの尤度を掛け合わせていくことができたのです。それに対し、体重のデータは明らかに
 8135 違います。すなわち、今日の体重が明日の体重に影響している、それどころか今日の体重は明日の体重と大
 8136 に関わっているはずですし、明日明後日とその後の日々にも少なからず影響しているはずなのです。

8137 このように考えると、独立していないデータ点、時系列的なデータに対しては異なる分析をする必要があり
 8138 そうです。次回はこの時系列的なデータを扱うモデルを紹介していくことになります。

8139 29.4 課題

8140 2015 年の 1 月 1 日から 2015 年の 12 月 1 日の間にも、体重の減少と増加が切り替わったところがあり
8141 そうです。変化点を合わせた折線回帰を実行し、何月何日に変化点があったのか推定するモデルを分析する
8142 R/Stan コードを提出してください。結果の解釈などを、スクリプトのコメントアウトや別添ファイルなどで提供
8143 してもらえると素敵です。もちろん Rmd ファイルでの提出であれば完璧です。なお提出されたコード単
8144 体でバグがなく動くことが確認できないものは、未提出扱いになります。コードの書き方などわからないところ
8145 があれば、曜日別 TA か小杉までメールで連絡し、指導を受けてください。

第 30 章

確率的プログラミング；状態空間モデル

前回は体重の変化データを例に、データの平均点が変わるモデル、変化点を検出したり変化点を境に傾きが変わる線形モデルをみてきました。

同じように、今回も時系列的なデータを扱います。前回の最後に案内したように、時系列データに普通の回帰直線をそのまま当てはめることは適切ではありません。どのような注意点が必要なのかについて、少し考えてみましょう。

30.1 時系列データの特徴

時系列的なデータはどのようなときに得られるでしょうか。前回導入した時のように、誰か 1 人が日々の生活を記録し続けている、それも立派な時系列データです。心理学には**日誌法**と呼ばれるデータ収集法があります。一言で行ってしまえば日記なのですが、日々多くの出来事や感情が到来するのを後から振り返って考えるのでは、正しく思い出せなかったり記憶が歪んでしまうこともあります。日々の状態を細かく記録し積み重ねることは、数字になっていなくても貴重なデータなのです。また最近では、**経験サンプリング**と呼ばれる調査法もあります。これは決まった間隔で質問紙がスマートフォンなどに送られてきて、定期的に心情を心理尺度に反映させていく、というものです。こうしたことができるようになったのには、私たちの周りに電子的なデバイスがたくさんあることが理由の 1 つです。さらに最近ではウェアラブル端末といって、身につける電子端末がありますね。Apple Watch など携帯電話と同じような機能を持ったこれらのデバイスは、GPS 機能がついていたり、万歩計による歩数の計測、血圧・脈拍のリアルタイムな計測が行われていたりします^{*1}。これを蓄積したデータとすると、かなりの時系列的なデータが手に入ることになります。

あるいはまた、社会的なデータも時系列的に得られるものが多いです。一番わかりやすいのは株価の指標などでしょうか。日経平均株価など、時事刻々と変化する株価がグラフになって表されているのをみなさんも見たことがあると思います。また、Twitter などの SNS で、今どのようなキーワードが使われているかといったことが、リアルタイムに計測されます。これらの指標は個々人の特徴というより、社会全体のうねりのようなものを表現していることになりませんが、これを研究することも立派な社会心理学的テーマになり得ます。時系列データは、人文社会科学領域ではこれまで経済学やマーケティングなどが専門的に扱ってきたところがありますが、これからは心理学の領域でもこうしたデータを分析することが流行してくるかもしれません^{*2}。

さて時系列的なデータはどういった特徴があるかと言うと、端的に言えば**自己相関 (auto correlation)**がある、ということでしょう。すなわち、ある時点のデータはその前の時点と相関している、ということです。体重のデータの場合でもそうですが、明日の体重は突然ランダムな値になるのではなく、今日の体重に応じ

*1 こうしたデータは**ライフログ (Life-log)**と呼ばれることもあります。

*2 もちろん生理心理学や動物心理学では、古くからこうしたデータを扱ってきています。

8175 て変化するはずだからです。また、データの変動に周期的なパターンが存在することもあります。脳波などの
 8176 データは基本的に波ですから、大きな波、小さな波がパターンを描いています。株価など社会的なデータも、
 8177 季節による一定の変動など全体的な揺らぎがある上で、目的とする意味のある変化を見つけ出さなければな
 8178 らないという課題に取り組んでいます。たとえば周波数のデータの平均はゼロになりますから、基本的な記述
 8179 統計では太刀打ちできないところがあります。回帰分析もその前提として、各データが独立に得られていると
 8180 言うものがありますから、自己相関するデータは当然この仮定に違反していることとなります。こうしたデータ
 8181 を分析するためには、周波数の大きさをスペクトル解析するとか、多次元の行列であるテンソルを使ってさま
 8182 ざまな要素を分解する、といった特殊な応用数学を駆使する必要があります。

8183 今回はこうしたさまざまなモデルの中でも、**状態空間モデル (State Space Model)** を紹介します。

8184 30.2 状態空間モデル

8185 時系列的なデータは、自己相関すなわち前の時点の状態が、次の時点に影響しているという特徴があり、
 8186 これをうまく表現するのが状態空間モデルです。状態空間モデルはこれまでの**モデリング**の技法を使うと
 8187 簡単に実装できます。というか、これまでの時系列的なデータ解析は、自己相関を減らすために一時点前の
 8188 データとの差分をとってそれをモデルにする、時点ごとの変化をスムージングするなど、さまざまな工夫の上に
 8189 成り立っていました。状態空間モデルはそれを非常にすっきり表現してくれているので、それまでの分析ノウ
 8190 ハウがなくとも誰でも簡単に利用できるモデルだといえるかもしれません。

8191 状態空間モデルは、まず観測されたデータと、その背後にある**状態**を区別します。体重変化のデータの例
 8192 で考えてみましょう。計測された体重は、当然体の重さを反映したのですが、測定の際に多少の誤差が生じ
 8193 るでしょう。例に使っている筆者のデータでも、毎朝全裸で計測しているわけではありませんから、季節の変
 8194 わり目で薄いパジャマからジャージにしたり、冬用の厚手のパジャマに変えたりすると、それだけで計測され
 8195 た数字には違いができます。本来知りたいのは、体の重さ、肉の重さそのものなはずですね。ここでいう体重
 8196 そのもの、計測誤差のない肉の重さが、ここでいう状態のことになります。

8197 そして t 時点の状態 μ_t は、次の時点の状態に影響します。服の重さなどを取り除いた t 時点の肉の状態
 8198 に、食事や運動といった変化が加わって明日 $t + 1$ 時点の肉の状態になるわけですから、 $\mu_{t+1} = f(\mu_t)$ と
 8199 考えることができるわけです。

8200 もし翌日の体重が、今日の体重にちょっとした誤差がつくような変動しかしない、と言うのであれば
 8201 $\mu_{t+1} \sim N(\mu_t, \tau)$ と表すことができるでしょう。あるいは、運動量 P_t の関数だというのであれば、 $\mu_{t+1} \sim$
 8202 $N(\mu_t + \beta_1 P_t, \tau)$ のように回帰分析の確率モデルのように表現すれば良いでしょう。もちろん説明変数が増
 8203 えたり、一次関数ではないものを考えても構いません。

8204 一方、計測値 W_t は、この μ_t に誤差 σ がついて得られる、すなわち $W_t \sim N(\mu_t, \sigma)$ という関係になりま
 8205 す。これらの関係を表現したのが、図 30.1 になります。このようにして、状態が時間 $t, t + 1, t + 2, \dots$ を通
 8206 じてつながっており、状態に応じて観測値が得られると考えます。観測値を手に入れる際には偶然誤差 σ が
 8207 生じますし、状態の変化は確率的に変わるのであれば幅 τ で、系統的に変わるのであればそれをモデリング
 8208 してやれば良いこととなります。

8209 この発想に基づいて、設計図を書いてみましょう。ここでは体重の変化に特別な傾向を考えず、ただ幅 τ で
 8210 確率的にノイズが加わる**ホワイトノイズモデル (white-noise model)** で考えてみることにします。少しイ
 8211 メージしにくいかもしれませんが、たとえば次のように表現できます (図 30.2)。

8212 それではこれをコードにしてみましょう。少し長くなりましたが、ホワイトノイズのモデルはコード 30.1 のよう
 8213 になります。

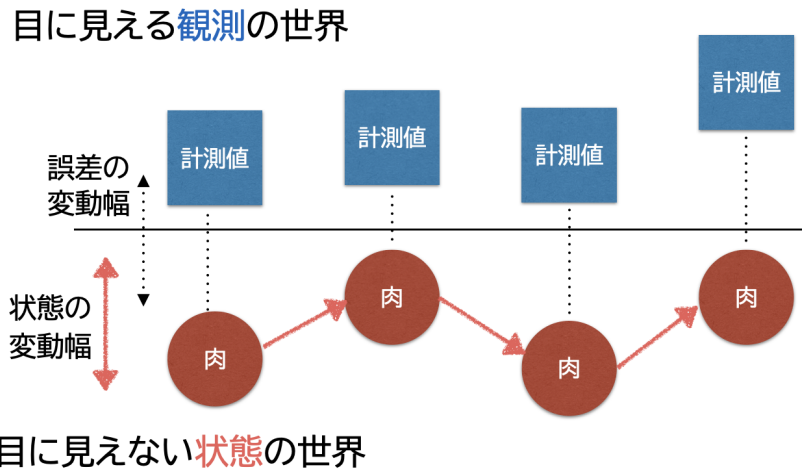


図 30.1 状態空間モデルのイメージ (体重変化の例)

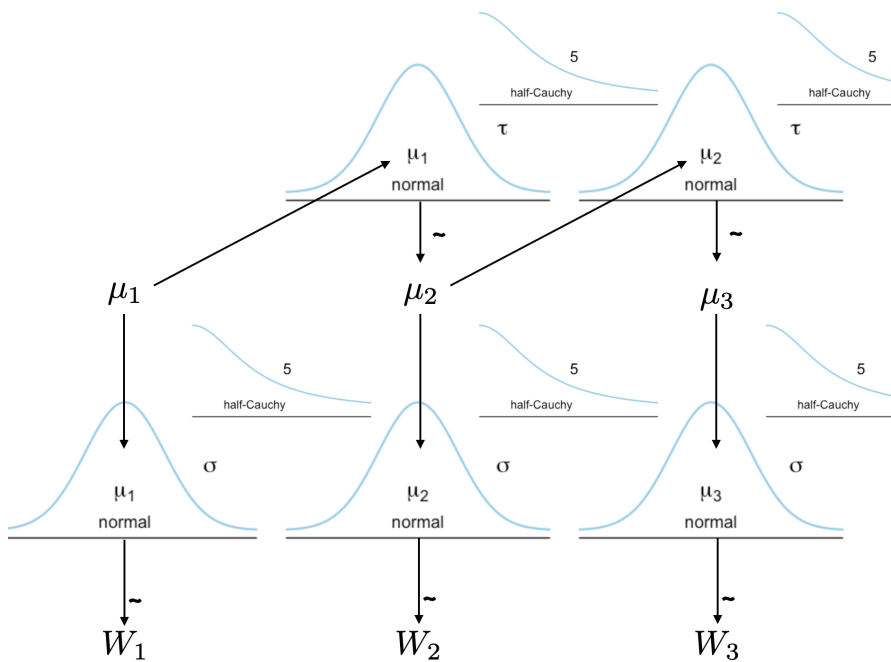


図 30.2 状態空間モデルの設計図

code : 30.1 ホワイトノイズモデル

```

8214 1 data{
8215 2   int L;
8216 3   array[L] real W;
8217 4 }
8218 5
8219 6 parameters{
8220 7   real muZero;
8221 8   array[L] real mu;

```

```

8223 9   real<lower=0> sig;
8224 10  real<lower=0> tau;
8225 11  }
8226 12
8227 13  model{
8228 14    mu[1] ~ normal(muZero, tau);
8229 15
8230 16    for(l in 1:L){
8231 17      W[l] ~ normal(mu[l], sig);
8232 18    }
8233 19
8234 20    for(i in 2:L){
8235 21      mu[i] ~ normal(mu[i-1], tau);
8236 22    }
8237 23
8238 24    muZero ~ normal(80, 10);
8239 25    sig ~ cauchy(0, 5);
8240 26    tau ~ cauchy(0, 5);
8241 27  }
8242

```

8243 ■コード解説

8244 **data ブロック** データ長 L , 各時点の体重 W がデータになります。

8245 **parameters ブロック** データ長と同じだけの状態を表す μ と、状態の変動幅 τ , 状態に付加される計測誤差の変動幅 sig を宣言しています。それに加えて、最初の状態 μ_0 をパラメータとしました。これは 2 時点目の状態 μ_2 は $\mu_2 \sim N(\mu_1, \tau)$ で、3 時点目は $\mu_3 \sim N(\mu_2, \tau)$ で...と表現できるのに対し、最初の点は $\mu_1 \sim N(\mu_0, \tau)$ となって、計測される前の状態を考えなければならないからです。わからないものは確率で表現する、というベイズの流儀に則って、これはパラメータとして推定してやることにします。

8251 **model ブロック** まず 1 時点目の状態は、0 時点目の状態から出てくるものですので、 $\mu_1 \sim N(\mu_0, \tau)$ をモデル化しました。次に計測された値 W_l は状態 μ_l から出てくるものですから、1 時点目からデータ長 L まで順に $W_l \sim N(\mu_l, \tau)$ として尤度を書きます。次に状態のモデルですが、2 時点目以降は前の時点からの影響で表現できるので、 $\mu_i \sim N(\mu_{i-1}, \tau)$ としています。最後に事前分布として、これまでの経験から $\mu_0 \sim N(80, 10)$ とし、変化の幅はいずれも SD ですから半コーシー分布で表現しました。

8257 これを使って、体重のデータを分析してみましょう。データは 2020 年からのものを使います。図 30.3 をみると、2021 年まではあまり変化がなく、21 年初頭に少し体重が下がって、また上がって、その後 5 月から謎の下降^{*3}、そしてその反動^{*4}、とおそらく線形モデルでは表現できないような複雑な動きをしています。

8260 さて、コード 30.2 を実行し、図 30.4 のような結果を得ます。複雑な動きについても、モデルがしっかりと予測しているのがみて取れますね。状態空間のパワフルな表現力をお楽しみいただけましたでしょうか。

code : 30.2 状態空間モデルのコード

```

8262 1  dat1 <- dat %>%
8263 2  filter(date > "2020/01/01") %>%
8264

```

*3 自転車事故による鎖骨骨折と、それに伴ってお箸持ちにくくて食事が減る、運動も減るといったのが原因でしょう。

*4 悔しいです。

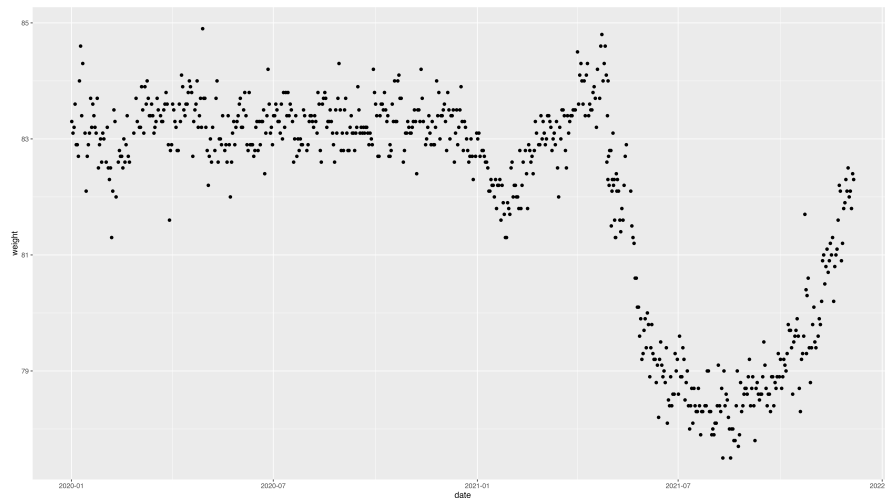


図 30.3 2020 年の体重変化

```

8265 3   mutate(date = as.Date(date))
8266 4
8267 5   model <- cmdstan_model("StateSpace.stan")
8268 6   dataSet <- list(L = NROW(dat1), W = dat1$weight)
8269 7   fit1 <- model$sample(
8270 8     data = dataSet,
8271 9     chains = 4,
8272 10    parallel_chains = 4
8273 11 )
8274 12
8275 13 fit1.stanfit <- fit1$output_files() %>% rstan::read_stan_csv()
8276 14 fit1.df <- fit1.stanfit %>% as.data.frame() %>%
8277 15   as_tibble() %>%
8278 16   rowid_to_column("iter") %>%
8279 17   pivot_longer(-iter, names_to = "Varname") %>%
8280 18   group_by(Varname) %>%
8281 19   summarise(
8282 20     EAP = mean(value),
8283 21     MED = median(value),
8284 22     MAP = map_estimation(value),
8285 23     SD = sd(value),
8286 24     L95 = quantile(value, probs = 0.025),
8287 25     L50 = quantile(value, probs = 0.25),
8288 26     U50 = quantile(value, probs = 0.75),
8289 27     U95 = quantile(value, probs = 0.975)
8290 28 )
8291 29
8292 30 Est1 <- fit1.df %>%
8293 31   dplyr::filter(str_detect(Varname, "mu")) %>%
8294 32   dplyr::mutate(ID = str_extract(Varname, pattern = "\\d+") %>%
8295 33     as.numeric()) %>%
8296 34   arrange(ID)
8297 35

```

```

8298 36 g <- dat1 %>%
8299 37   rowid_to_column("ID") %>%
8300 38   left_join(Est1, by = "ID") %>%
8301 39   ggplot(aes(x = ID, y = weight, ymin = U95, ymax = L95)) +
8302 40   geom_point() +
8303 41   geom_point(aes(x = ID, y = MAP), color = palette()[2]) +
8304 42   geom_ribbon(fill = palette()[3], alpha = 0.2)
8305 43 plot(g)
8306

```

8307 ■コード解説

8308 1-3 行目 データファイルから該当する日程だけ抜き出します。日付を Date 形式に変換しています。

8309 5-11 行目 cmdstanr でコンパイルしてサンプリングするところです。

8310 13 行目 cmdstanr で作ったファイルを stanfit オブジェクトに変換しています。rstan パッケージを使っ
8311 ている人はこの行を実行する必要がありません。

8312 14-28 行目 MCMC サンプルを集計し、EAP, MED, MAP 推定値, SD, 確信区間などを出してい
8313 ます。

8314 30-34 行目 推定された状態 μ_i を取り出しています。

8315 36-43 行目 もとのデータと推定値を合体させ、プロットしています。



図 30.4 2020 年の体重変化をモデルで追跡してみた結果

8316 30.3 欠損値の補間

8317 ところで、データをよくみてみると、実は毎日のデータになっていないところがあることに気づきます。いや、
8318 気づかないかもしれないですね。次のコード 30.3 を実行してみてください。

code : 30.3 日付は連続かな?

```

8319 1 dat1 %>%
8320 2   mutate(lag = lag(date)) %>%
8321 3   mutate(date = as.Date(date), lag = as.Date(lag)) %>%

```



```
8323 4 mutate(FLG = date - lag) %>%
8324 5 dplyr::filter(FLG > 1)
```

8326 ■コード解説

- 8327 1 行目 先ほどの推定に使ったデータファイルです。
- 8328 2 行目 lag 関数で、データを 1 行ずらした列を作ります。
- 8329 3 行目 体重を記録した日の変数 date と、先ほど作った一行ずらした変数 lag はいずれも日付に関する変数ですので、日付型に変更します。
- 8330
- 8331 4 行目 日付変数の引き算をして、一行下のデータと何日ずれていたのか算出しています。
- 8332 5 行目 1 日以上ずれている日をフィルタリングで抜き出しています。
- 8333 これをみると、ちよくちよく抜けていて、時には 4 日も空いていたりすることがわかります*5。ともかく、これでは
- 8334 正確にデータを分析できているとは言えません。

R の出力 30.1: 抜けている日

```
# A tibble: 21 × 5
  date      weight bodyFat lag      FLG
<date>    <dbl>  <dbl> <date>  <drtn>
1 2020-01-13  83.1   26.8 2020-01-11 2 days
2 2020-02-01  82.6   26.7 2020-01-30 2 days
3 2020-02-12  82.6   26.7 2020-02-10 2 days
4 2020-02-26  83.1   26.8 2020-02-22 4 days
5 2020-02-28  83.7    27   2020-02-26 2 days
6 2020-03-02  83.2   26.9 2020-02-29 2 days
7 2020-03-21  83.4   26.9 2020-03-19 2 days
8 2020-05-09  82.8   26.8 2020-05-07 2 days
9 2020-05-18  82.8   26.7 2020-05-16 2 days
10 2020-07-31  82.8   27.7 2020-07-29 2 days
# ... with 11 more rows
```

8335

8336 しかし計測はできていなくても、出張中にも状態の変化は続きいているはずで、次の計測時点はその日の

8337 状態の関数になっているはずで (図 30.5)。

8338 この欠損した状態のデータをどのように考えれば良いでしょうか。実はこのことのヒントはすでに、最初の

8339 コードの中に入っています。状態空間モデルの 1 時点目、 μ_1 は μ_0 から影響されているはずですが、データ

8340 から μ_0 はわかりません。しかしわからないことは確率で表現するのがベイジアンやり方です。今回も、欠損

8341 している測定値はわからないものとして確率で表現し、推定してやれば良いのです！このようにして間を埋

8342 めることをとくに**補間 (interpolating)** と言います。

8343 補間をするためには、R コードの方でも Stan コードの方でもすこし工夫が必要です。まず R のほうで、

8344 データを加工するところから見てみましょう。Stan には欠損値 NA を与えることはできませんから、欠損して

8345 いるところには特別なあり得ない数字、そうですね、999 という数字でも入れておきましょう。

code : 30.4 連続したデータを作り欠損値に借りの値を代入する

*5 2020 年 2 月 22 日から 24 日は岡山県にバイズ塾合宿 (出張) でした。コロナがまだ豪華客船の中だけに抑え込まれている時期で、ギリギリ出張できた時期です。このあと大学の卒業式がなくなったり、前期オンライン授業になったりと言う、波乱の 2020 年度が始まったのでした。

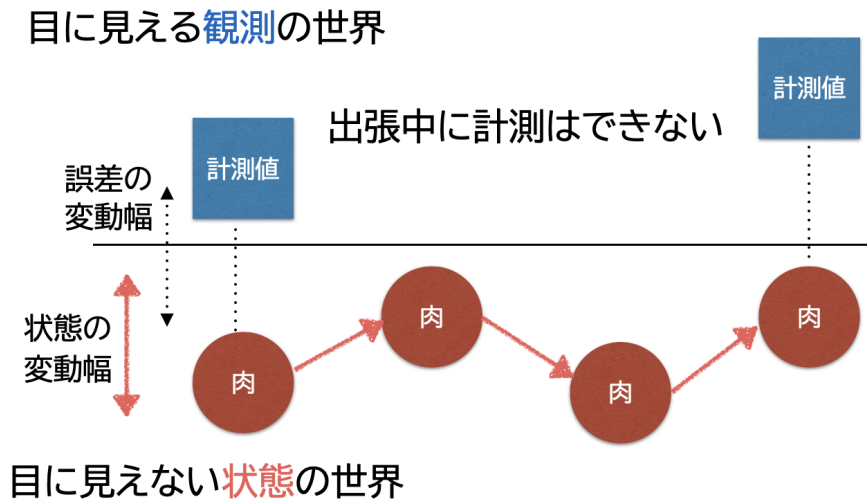


図 30.5 計測できない時点も変化は続く

```

8347 1 fullDays <-
8348 2   data.frame(date = as.Date("2020/01/01"):as.Date("2021/12/01")) %>%
8349 3   mutate(date = as.Date(date, origin = "1970-01-01")) %>%
8350 4   left_join(dat1, by = "date") %>%
8351 5   tidyr::replace_na(list(weight = 999, bodyFat = 999))
8352

```

8353 ■コード解説

8354 2行目 2020年1月1日から2021年12月1日までにデータを限定します。この期間、連続した日付を
8355 date変数に追加します。

8356 3行目 日付変数を日付型に変更しています。

8357 4行目 もとのデータを日付変数に結合しています。left_join関数は左にもとのデータセットを置き、右
8358 からdat1を引っ付けていきます。結合する際は、変数dateをキーにし、キー変数が同じものは同じ
8359 行に合わせるという湯やり方です。このやり方だと、該当する行がない場合(計測データセットの日付
8360 がない場合)、もとのデータセットは残して、体重や体脂肪変数はNAになります。

8361 5行目 replace_na関数で、欠損のところに999という数字を入れています。

8362 このようにすることで、次のようなデータセットができました(出力 30.2)。

R の出力 30.2: 抜けている日を無くした完全データセット

```

> fullDays
      date weight bodyFat
1  2020-01-01   83.3   26.70
2  2020-01-02   83.1   26.80
3  2020-01-03   83.2   26.90
4  2020-01-04   83.6   27.00
5  2020-01-05   82.9   28.70
6  2020-01-06   82.9   26.80
7  2020-01-07   82.7   26.70
8  2020-01-08   84.0   27.10
9  2020-01-09   84.6   27.30
10 2020-01-10   83.4   26.90
11 2020-01-11   84.3   27.20
12 2020-01-12  999.0  999.00
13 2020-01-13   83.1   26.80
14 2020-01-14   82.1   65.00
15 2020-01-15   82.7   26.70

```

8363

8364 これをみると、2020年1月12日は欠損だったのですが、データ上は999という数字が入って完全データ
8365 のように見えます*6。このデータを Stan に与えてやり、体重が999というあり得ない数字の場合は別の処理
8366 をする、という0 過剰ポアソンの技術を応用します (セクション 27.3, Pp.313)。

code : 30.5 欠損値補間のコード

```

8367
8368 1 data{
8369 2   int L;
8370 3   array[L] real W;
8371 4   int<lower=0> Nmiss;
8372 5 }
8373 6
8374 7 parameters{
8375 8   real muZero;
8376 9   array[L] real mu;
8377 10  array[Nmiss] real<lower=0> Miss_W;
8378 11  real<lower=0> sig;
8379 12  real<lower=0> tau;
8380 13 }
8381 14
8382 15 model{
8383 16  mu[1] ~ normal(muZero, tau);
8384 17
8385 18  {
8386 19    int j = 0;
8387 20    for(l in 1:L){
8388 21      if( W[l] != 999){
8389 22        // こっちは尤度
8390 23        W[l] ~ normal(mu[l], sig);

```

*6 ちなみに2020/01/11から12にかけては、犬会という研究会があったため関西に出張していました。宿に泊まったので、12日の朝のルーティンができなかったのです。

```

8391 24     }else{
8392 25         j = j + 1;
8393 26         // こっちはパラメータ
8394 27         Miss_W[j] ~ normal(mu[l], sig);
8395 28     }
8396 29 }
8397 30 }
8398 31
8399 32 for(i in 2:L){
8400 33     mu[i] ~ normal(mu[i-1], tau);
8401 34 }
8402 35
8403 36 muZero ~ normal(80,10);
8404 37 sig ~ cauchy(0,5);
8405 38 tau ~ cauchy(0,5);
8406 39 }
8407

```

8408 data ブロック データ長 L, 各時点の体重 W がデータです。加えて, 推定すべき欠損値の数もデータとして受け取ります。

8410 parameters ブロック データ長と同じだけの状態を表す mu と, 状態の変動幅 tau, 状態に付加される計測誤差の変動幅 sig を宣言しています。また最初の状態 μ_0 を表すパラメータと, 欠損の数だけある推定用パラメータ Miss_W を用意しました。

8413 model ブロック 先ほど同様, 0 時点目の状態から出てくるものですので, $\mu_1 \sim N(\mu_0, \tau)$ をモデル化しています。次に各回のデータについて, W_l が 999 でなければ (同じでない, という論理式は != で表現します), 普通の状態空間モデルのように尤度として計算します。そうではない, すなわちデータ $W_l = 999$ である, つまり欠損値であるはずということなら, $W_{miss} \sim N(\mu_l, \sigma)$ で推定してやります。ここで特殊な変数 j が出てきています。これは「何番目の欠損値か」を表す変数です。ブロックの中の一部だけで変数を宣言するため, まず for 文全体を中括弧で括り, はじめに $j = 0$ という数字を宣言しています。ここで欠損値推定のシーンになると, j のカウンターを 1 つ増加させ, j 番目の欠損値の推定をさせる, というやり方をしています。欠損値の数と行番号が合致しないので, こうした工夫が必要なのですね。その後のコードは先ほどと同じです。

8422 このようにして, 推定してみましょう。Stan に与えるデータセットは次のようにして作ります。

code : 30.6 Stan に与えるデータセット

```

8423
8424 1 dataSet <- list(L = NROW(fullDays), W = fullDays$weight,
8425 2               Nmiss = sum(fullDays$weight == 999))
8426

```

8427 欠損値の数は, オブジェクト fullDays の weight 変数が 999 かどうかを判定させ (== でイコールかどうかという論理判断になります), 条件が合致した数を数え上げるという方法をとっています。

8429 コードが走り出すと, 推定は順調に進むと思います。が, 結果の出方が少しややこしいです。状態 μ_i はすべての日程について推定されますが, これに加えて欠損があったデータも推定されていきますので, 描画する際は「実測値なのか推定値なのか」, 「補間した欠損の推定値は何番目の欠損値だったか」を判断しながら結合していき, データにする必要があるからです。Stan の中で作ったカウント変数 j を R の方でも考えてやらないといけないわけですね。

code : 30.7 プロット用にデータを整形する関数とプロットのコード

8434

```

8435 1 Est2 <- fit2.df %>%
8436 2   dplyr::filter(str_detect(Varname, "mu")) %>%
8437 3   dplyr::filter(!str_detect(Varname, "muZero")) %>%
8438 4   dplyr::mutate(ID = str_extract(Varname, pattern = "\\d+") %>%
8439 5     as.numeric()) %>%
8440 6   arrange(ID) %>%
8441 7   select(ID, MAP, U95, L95)
8442 8
8443 9 Est2miss <- fit2.df %>%
8444 10  dplyr::filter(str_detect(Varname, "Miss_W")) %>%
8445 11  dplyr::mutate(ID = str_extract(Varname, pattern = "\\d+") %>%
8446 12    as.numeric()) %>%
8447 13  arrange(ID) %>%
8448 14  select(ID, MAP, U95, L95)
8449 15
8450 16 ### plot用の関数を準備
8451 17 plotFunction <- function(fullDays, Est, MissEst) {
8452 18   tmp <- fullDays %>%
8453 19     rowid_to_column("ID") %>%
8454 20     left_join(Est, by = "ID") %>%
8455 21     rowwise() %>%
8456 22     mutate(FLG = if (weight != 999) {1} else {2})
8457 23   misJ <- 1
8458 24   tmp$weight2 <- NA
8459 25   tmp$weight2U <- NA
8460 26   tmp$weight2L <- NA
8461 27   for (i in 1:NROW(tmp)) {
8462 28     if (tmp$FLG[i] == 2) {
8463 29       tmp$weight2[i] <- MissEst$MAP[misJ]
8464 30       tmp$weight2U[i] <- MissEst$U95[misJ]
8465 31       tmp$weight2L[i] <- MissEst$L95[misJ]
8466 32       misJ <- misJ + 1
8467 33     } else {
8468 34       tmp$weight2[i] <- tmp$weight[i]
8469 35       tmp$weight2U[i] <- tmp$weight[i]
8470 36       tmp$weight2L[i] <- tmp$weight[i]
8471 37     }
8472 38   }
8473 39   return(tmp)
8474 40 }
8475 41
8476 42 plot.tmp <- plotFunction(fullDays, Est2, Est2miss)
8477 43 g <- ggplot(data = plot.tmp) +
8478 44   geom_point(aes(x = date, y = weight2)) +
8479 45   geom_errorbar(aes(x = date, y = weight2, ymin = weight2L,
8480 46     ymax = weight2U, color = palette()[2])) +
8481 47   geom_point(aes(x = date, y = MAP, color = palette()[3])) +
8482 48   geom_errorbar(aes(x = date, y = MAP, ymin = L95,
8483 49     ymax = U95, color = palette()[4])) +
8484 50   scale_x_date(date_breaks = "1_month",
8485 51     limits = as.Date(c("2020-01-01", "2020-05-01"))) +

```

```
8486 52 theme(legend.position = "none")
8487
```

8488 ■プログラム解説

8489 Est2 を作るブロック MCMC サンプルをデータフレーム化したものの中から、状態 μ に関するデータだけ
8490 を抜き出したものを作る。

8491 Est2miss を作るブロック 欠損値を補間するための推定値を、MCMC サンプルから取り出した欠損値
8492 補間だけのデータセットを作る。

8493 plot 用の関数 プロットはこの後にも行いますので、もう関数を作ってまとめてしまうことにします。もとの
8494 データセット fullDays と、状態の推定値データセット Est, 欠損値の推定値データセット MissEst
8495 を引数にとる関数です。

8496 関数 3 行目 まずもとデータに行番号を ID として変数化します。

8497 関数 4 行目 状態の推定値は、すべての行について推定しているはずですので、行番号をキーに結合
8498 left_join します。

8499 関数 5 行目 もとの観測値が実測値なのか欠損値なのかを表現する変数 FLG を用意しておきます。
8500 これらはすべて、一時的なオブジェクト tmp に納められます。

8501 関数 6 行目 欠損値のインデックスを表す変数 misJ をつくり、カウントを 1 に設定しておきます。

8502 関数 7-9 行目 先ほどの一次オブジェクトにいまから体重の推定値を入れていきますが、この推定値
8503 は実測値で得られている時は不要で NA になるはずですので、いったんすべてに NA を代入して
8504 います。

8505 関数 10-21 行目 オブジェクト tmp を一行ずつ見ていって、もし FLG 変数が欠損値であることを教
8506 えてくれたら、misJ 番目の推定値を代入します。代入が終わったらカウントを 1 つ追加しておき
8507 ます。欠損値でなければ、もとの体重をそのまま移し入れます。ちなみに U や L がついているの
8508 は 95% 区間の上限と下限なのですが、実測値の場合は推定ではありませんのですべて同じ数字
8509 になります。

8510 関数 22 行目 作ったオブジェクト tmp を戻り値として返します。

8511 関数の結果を受けとる plot.tmp は描画用の一時データオブジェクトです。ここに全期間のデータと、今
8512 回の状態推定値、欠損推定値を入れておきます。

8513 描画する x 軸は全日程です。まずは体重 W_j を geom_point と geom_errorbar でプロットします。欠
8514 損の場合は推定値が入っています。次のポイントとエラーバーの geom は、状態 μ のものになります。

8515 さあどうでしょう。コードはともかく、図 30.5 を見てみましょう。実際に観測しているのは黒点、状態の推定
8516 値は緑点と赤いエラーバーで表示されています。青いバーは欠損した時の推定値で、値がないのですから当
8517 然触れ幅は大きいですが、おそらくこの体重が測定できなかった日はこれぐらいだったんじゃないか、という
8518 値が推定されているわけです。数字で見ると、たとえば 2020 年 1 月 12 日の MAP 推定値は 83.4kg, 95%
8519 区間でいうなら [82.6,84.2] です。前日が 84.3kg で翌日が 83.1kg ですから、そんなもんかなあという気も
8520 します。

8521 ところで、測定していないはずの日も補間できるよ、ということを考えてみると、これはすごいことだと思
8522 いませんか。そうです、未来の値もちろん「測定していない日」ですから、このモデルを使うと未来予測が
8523 できるのです！すごい！

8524 早速同じコードで、日程を未来に飛ばしてみましょう。データは 2021 年 12 月 06 日までしかありません*7

*7 この原稿は 2021 年 12 月 10 日 17 時に執筆しています。



図 30.6 欠損していても推定できる

8525 から、2021 年 12 月 31 日までデータを伸ばし、推定してみようではありませんか。コードは省略しますが、実行した結果は次のような図 30.7 として得られます。



図 30.7 未来もある意味欠損値

8526

8527 これをみると、データが得られなくなった次の日から、確信区間が日を追うごとに広がっていく様が見取
8528 れます。当然そうですね。 μ_{t+1} は μ_t にプラスかマイナスか、いずれとも言えない振幅幅で揺れ動くのですか
8529 ら、先に行けば行くほどその可能性は不確かなものになっていくのです。ちなみに描画用のデータから 2021
8530 年最後の数日間の予測値を見てみると、次のようになっています。

code : 30.8 予想される 2021 年末の体重 (一部略)

8531

8532 1 > plot.tmp %>% tail

8533 2 # A tibble: 6 × 11

8534 3 # Rowwise:

8535 4 ID date weight bodyFat MAP U95 L95 FLG weight2

8536 5 <dbl> <date> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>

8537	6	738	2021-12-26	999	NA	82.2	83.8	80.5	2	82.1
8538	7	739	2021-12-27	999	NA	82.3	83.9	80.4	2	82.1
8539	8	740	2021-12-28	999	NA	82.2	83.9	80.4	2	82.2
8540	9	741	2021-12-29	999	NA	82.2	84.0	80.4	2	82.0
8541	10	742	2021-12-30	999	NA	82.3	84.0	80.3	2	82.3
8542 8543	11	743	2021-12-31	999	NA	82.3	84.0	80.3	2	82.2d

8544 未来の状態 μ は 82.2 のあたりでほぼ固定, それに伴って想定される体重 `weight2` も 82.2 程度ですが, 幅
8545 がどんどんと広がっていきます。可能性として, 80kg に落ちているかもしれませんし, 84kg になっている
8546 かもしれない, というのが現時点での予測になります。

8547 予測ができるとは言え, なんだかほとんど変わり映えしない数字で面白くないですね。それは当然, 「明日
8548 の体重は今日の体重 + 誤差だろう」という, とくに何らメカニズムや仮定を入れなかったものですから, ほぼ
8549 横一線の推定にしかならないわけです。

8550 では少し欲張って, 構造を入れみましょうか。ここまでは, μ_t は μ_{t-1} にのみ影響される, というモデルに
8551 なっていましたが, もう 1 日前の影響が入っていることを考えます。たとえば体重が増加傾向にあるとか, 減
8552 少傾向にあるといった, 変化を考えてみることにしましょう。すなわち, $\mu_t - \mu_{t-1} = \mu_{t-1} - \mu_{t-2} + \varepsilon$ です。
8553 これは昨日から今日への変化 ($\mu_t - \mu_{t-1}$) は, 一昨日から昨日への変化 ($\mu_{t-1} - \mu_{t-2}$) と同じようなものだ
8554 ろう (誤差 ε はあるけど), と考えていることになります。このようなモデルを **2 階差分のトレンド**と呼びます。

8555 この式を変形すると,

$$\begin{aligned}\mu_t &= \mu_{t-1} + \mu_{t-1} - \mu_{t-2} + \varepsilon \\ &= 2\mu_{t-1} - \mu_{t-2} + \varepsilon\end{aligned}$$

8557 となることはすぐにわかりますね。これを確率モデルで表現するなら,

$$\mu_t \sim N(2\mu_{t-1} - \mu_{t-2}, \tau)$$

8558 と考えていることと同じですから, そのように Stan のコードも変形します (コード 30.9)。

code : 30.9 2 階差分トレンドのコード

```

8559 1  ... (前略)...
8560 2  model{
8561 3    mu[1] ~ normal(muZero, tau);
8562 4    mu[2] ~ normal(mu[1], tau);
8563 5
8564 6    {
8565 7      int j = 0;
8566 8      for(l in 1:L){
8567 9        if( W[l] != 999){
8568 10           // こっちは尤度
8569 11           W[l] ~ normal(mu[l], sig);
8570 12        }else{
8571 13           j = j + 1;
8572 14           // こっちはパラメータ
8573 15           Miss_W[j] ~ normal(2*mu[l-1]-mu[l-2], sig);
8574 16        }
8575 17      }
8576 18    }
8577 19
8578 20  for(i in 3:L){

```

```

8580 21 //2階差分
8581 22 mu[i] ~ normal(2*mu[i-1]-mu[i-2],tau);
8582 23 }
8583 24
8584 25 muZero ~ normal(80,10);
8585 26 sig ~ cauchy(0,5);
8586 27 tau ~ cauchy(0,5);
8587 28 }
8588

```

8589 コードは2階差分になっているので、状態は3時点目からしか推定できません。最初の2行で μ_1, μ_2 をメカ
8590 ニズムに沿って推定し、変化が取れるようになってから、は2階差分のコードで推定します。for文が3から
8591 始まっていることに注意してください。

このコードと年末までのデータ、そしてプロット関数を使って推定した結果が次の図 30.8 になります。

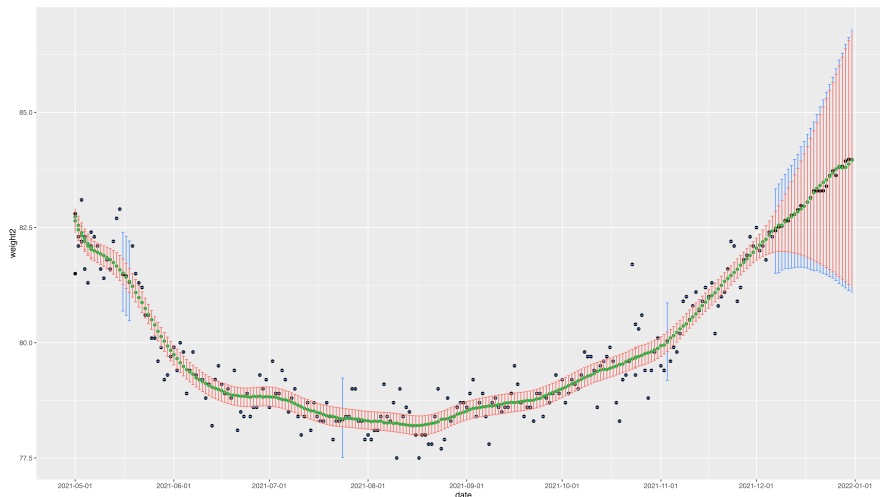


図 30.8 2階差分トレンドの予測

8592

8593 これを見ると、それまでの増加傾向を反映して、年末には 84kg、最悪の場合 86.6kg まで増えている可能
8594 性があることがわかります*8*9。

code : 30.10 2階差分トレンドで予想される 2021 年末の体重 (一部略)

```

8595 1 > plot.tmp %>% tail()
8596 2 # A tibble: 6 × 11
8597 3 # Rowwise:
8598 4   ID date      weight bodyFat  MAP  U95  L95  FLG weight2
8599 5 <dbl> <date>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
8600 6   738 2021-12-26  999    NA  83.9 85.6 81.5  2  83.6
8601 7   739 2021-12-27  999    NA  83.9 85.8 81.5  2  83.8
8602 8   740 2021-12-28  999    NA  84.0 85.9 81.4  2  83.8
8603 9   741 2021-12-29  999    NA  84.1 86.1 81.4  2  83.9
8604 10  742 2021-12-30  999    NA  84.2 86.3 81.3  2  84.1
8605

```

*8 なんてこった、まだ執筆中は答えが出ていませんが、この予測通りになったらとんでもないことですよ。だって私は 20 代 66kg だったんです。あの頃から 20kg も増えている自分の未来なんて想像したくないですよ。またダイエット始めるかなあ。

*9 2021/01/15 追記。2021/12/31,6:07AM の記録は 82.3kg, 体脂肪率 26.6% でした。ちゃんと確信区間の中に入ってますね！

8606 8607	11	743	2021-12-31	999	NA	84.2	86.5	81.2	2	84.0
--------------	----	-----	------------	-----	----	------	------	------	---	------

8608 今回は 2 階差分のトレンドを入れての予測となりましたが、この他にも周期的に影響してくる季節項をいれ
8609 るとか、たとえば「週末は体重が増加するだろう」と考えて週末効果の項を入れるとか、さまざまな工夫を思い
8610 つくことができるのではないのでしょうか。

8611 30.4 状態空間モデルの展開

8612 最後になりましたが、状態空間モデルは表面に現れている値の背後に潜在変数を仮定し、それを滑らかに
8613 繋ぎ合わせたスムージング (smoothing) の技術であるとも考えることができます。時系列を通じて、その
8614 前後の「状態」を、柔らかなバネで結合し、データ全体を通じて変化する関数として相互に調整させあっている
8615 る、と考えることもできるでしょう。

8616 また、時間というのは過去から未来へ進む、1 次元の道のりです。ところで私たちは 2 次元の地図を見たり
8617 三次元の空間を見たりしていますよね。たとえばこの状態空間が、時系列的次元ではなく、X 軸と Y 軸に
8618 広がる網のような二次元結合していることを考えると、空間的な統計分析が可能になります。たとえば地価と
8619 というのは、駅からの距離や大型商業施設、コンビニエンスストアなどの利便性に応じて変化しますが、おそらく
8620 それらは連綿と繋がっていて、徐々に変化していくものであるはずで、であれば徐々に変化する要素をつ
8621 なぎ合わせて、面の状態空間モデルを作ってやることもできるわけです。さらにその地価の変動を考える、つ
8622 まり二次元平面 × 時系列に潜在変数をつなげてスムージングしたモデルを作れば、時系列的な地価の変動
8623 を考えたり予測したりできるかもしれません。

8624 心理学はこれまで、調査法で一時点を切り出してその構造を把握する、というアプローチが多くありまし
8625 た。時系列を扱うとしても、プレ・ポストの変化を見るときか、あまり長期にわたって追いかけるようなことはして
8626 きませんでした。しかし冒頭でお話したように、これからの心理学データは時間的にも空間的にも広がるも
8627 のを計測し、利用できるようになるかもしれません。そのためにはより深い人間の観察や推察によるモデルの
8628 作り込みと、状態空間モデルとベイズ推定のような数理モデルと強力な推定法のタッグが、必要になってくる
8629 でしょう。

8630 30.5 課題

8631 体重変化のデータには、体脂肪の記録もあります。体重 × 体脂肪率から、筆者の筋肉量を計算できます。
8632 2020 年以降のデータを使って、筋肉量の変化を状態空間モデルで推定してみてください。2 階差分トレンド
8633 モデルで、未来のデータも予測できればなお結構です。これらのモデルを分析する R/Stan コードを提出し
8634 てください。結果の解釈などを、スクリプトのコメントアウトや別添ファイルなどで提供してもらえると
8635 素晴らしいです。もちろん Rmd ファイルでの提出であれば完璧です。なお提出されたコード単体でバグがなく動くこと
8636 が確認できないものは、未提出扱いになります。コードの書き方などわからないところがあれば、曜日別 TA
8637 か小杉までメールで連絡し、指導を受けてください。

第 31 章

モデル比較

さて、今回でこの授業も最終回となりました。今回は授業全体のまとめとして、これまでの流れを俯瞰的に捉え直すとともに、今後みなさんがこの講義を通じて学んだことが、どのように利用できるのかに言及していきたいと思います。

本書の前半、心理学データ解析応用 1 として前期にお話ししてきた内容は、心理学で利用されるデータ、とくに心理尺度を用いたものが、どういう根拠で「心」を測定していると言えるのかについて解説してきました。後半、心理学データ解析応用 2 の方は、プログラミングやベイズ統計といった目新しい研究手法の方が目についたかもしれませんが、その背後にあったメッセージは、どういうメカニズムでデータが生成されてきたのかという、**データ生成メカニズム**というアイデアをもつことでした。その意味で、本書全体のメッセージは一貫しています。すなわち、心理学者がデータを取るという時の、データに込められた意味や想定されるメカニズムをどのように表現し抽出し解釈するかという、心理学的営みそのものをお伝えしたかったのです。ここに無知・無自覚なまま、心理尺度で何かが測定できると考えたり、検定の結果から意味を解釈したりできるはずがないのですから。

そのための準備として、確率統計に関する理論や、コンピュータ、プログラミングに関する技能などの習得が必要だったわけです。心理学を学ぼうと思って入ったはずの大学で、どうして数学やプログラミングをしているのか、こんなはずじゃなかったと思った人もいるかもしれませんが、逆にどうして既知の知識や技術だけで、この摩訶不思議な心、人間、社会のありようが理解できているのかと言いたいほどです。もちろん他のアプローチや、他にも学ぶべきことがたくさんありますし、時間は有限ですから苦手な手法は後回しにしたくなる気持ちはわかりますが、実は他の教養を身につけるよりもこちらの方が直接的で近道なルートだといえるかもしれません。

31.1 ベイジアンモデリング

Stan を使った統計的アプローチは、**ベイジアンモデリング (Bayesian Modeling)** と呼ばれることがあります。Stan が事後分布からの乱数発生機であり、事前分布と尤度を記述するだけで事後分布を得ることができるというところが、ベイズの定理に根拠を持ったベイズ統計学ですから、「ベイジアン」という冠がつくのは当然です。しかし後半の「**モデリング**」については、なにもベイズ統計学に限った話ではありません。モデルを立てて考えるということだけを考えれば、広い意味では心理学だって「心」というモデルを立てて考えているのですから。もっと言うと、学問というのは一般的に抽象化 (あるいは単純化) した理想の世界での論理的構造を扱うものですから、すべてモデリングアプローチだと言えるかもしれません。

心理学はその研究方法として、調査、実験、観察、介入といったものがあります。これらを通じて心とはなにかということを考えていくわけです。とくに実験や介入においては、条件の制御や統制を行って、それに伴って

8669 生じた結果の違いを**効果**と考えるのでした。そのことと**帰無仮説検定**は非常に相性が良かった事は、想像に
8670 難くないでしょう。帰無仮説検定は、出てきたデータと変数の関係について、機械的に判断を下すことができ
8671 る方法論でした。帰無仮説検定が農学から生まれてきたことからわかるように、その方法はどの研究分野
8672 に対しても適用できる、共通のルール・御作法なのです。心理学者は実験計画に心理学的要素を埋め込み、
8673 結果を統計的に判断することで、心理学的な結果の解釈を進めるてきました。いいかえれば、実験の計画を
8674 立てることそのものが心理学のエッセンスであったのです。

8675 しかし皆さんはすでに、実験計画が統計的に見れば**線形モデル**にすぎないことを知っていますね。そして
8676 ごく単純な線形モデルによってしか、心理学のエッセンスを表現する方法がなかった時代にあっては、想定し
8677 う心のメカニズムが平均値の差として出てくるように置き換えるほかなかったのです。方法論に自由がない
8678 分、その他のところで工夫するしかなかったとも言えます。そうした工夫の中には芸術的なまでに作り込ま
8679 れているものがあつたりしますから、そこに痺れる懂れる、という人も少なからずいるのは当然ともいえるで
8680 しょう。

8681 線形計画の結果は、操作/介入の効果があつたかなかつたかという判断になりがちです。結果がそれしか
8682 示しておらず、その結果を生むための工夫はすでになされており、得られるデータの精度も平均値差を示すこ
8683 とができれば良い程度にしかない、ともいえるかもしれません。しかし測定のツールは時事刻々と進歩してき
8684 ましたので、より精緻の違いを見て、よりリッチな解釈を許す意味合い豊かなデータも得られるようになってき
8685 ました。そのような時、伝統的な方法論だけで立ち向かうのはもったいないのではないのでしょうか。方法論の
8686 呪縛から自由になり、さまざまな表現ができるようになれば、心理学ももっと発展させることができるでしょう。

8687 **モデリング**は、変数間関係を数式で表現することでもあります。「心」や「気持ち」といった表現を始め、「幸
8688 せ」とか「恋愛感情」といった日常用語で考えるべきところを、 $X_i, Y_{ij}, \alpha, \beta, \dots$ といった記号で置き換えて表
8689 現するには慣れが必要です。ひょっとしたら抵抗を感じる人もいるかもしれません。しかし、数式で表現す
8690 る事は、他の解釈を許さない一意的表現をすることでもあります。人間って色々だよね、当てはまることもあ
8691 れば当てはまらないこともあるよね、という日常的な理解で留まりたいのであれば – そしてそこで止まってい
8692 られるのは幸せなことでもあるのですが – それでいいのですが、学問を進める以上はより明確に、より誤解
8693 のない表現をしなければなりません。個別のケースにしか当てはまらないこと、その程度を確率で表現して、
8694 理論的に全体的な傾向を「こうなっているに違いない」と厳密に表現するには、日常用語ではなく数式用語が
8695 必要なのです。モデリングを進める上で、そうした表現ができるようにトレーニングする必要はありますが、使
8696 えるようになると折線回帰や状態空間モデルのように、さまざまな表現ができることも見てきた通りです。

8697 この授業を受けてきたからと言って、さあ今すぐ自分で心理学的モデルを作れと言われても難しい、と思う
8698 かもしれません。実際これまで私が教えてきた中で、よく質問されるのは「やり方はわかるけれども、自分でや
8699 れる気がしない。どうやってモデルを思いついたらいいの？」ということです。これに対する答え方は2つあつ
8700 て、1つは「色々なパターンを見て模倣する」です。「学ぶ」は「真似ぶ」からきているとも言われたりしますが、
8701 どのような表現の可能性があるのかについては、実際の応用パターンをさまざま見るのが早いでしょう。馴染
8702 みのない料理や調理法があれば、レシピ本を買って応用パターンをいろいろ仕入れますよね。そうすること
8703 で、「こんな工夫ができるのか。じゃあ自分はこうしてみよう、ここをちょっと変えてみよう」と前に進むことがで
8704 けるのですから。こうしたモデルがたくさん載っているものとして、[Lee and Wagenmakers \(2013\)](#) や [豊田 \(2017, 2018, 2019\)](#) などが
8705 ありますので、パラパラとみながら応用例を広く眺めてみるといいでしょう。

8706 もう1つの答えは、「紙とペンをつかって設計図を書く」というものです。この方法は、この講義の中で再三
8707 お伝えしてきたやり方になります。データを見た時に、どういうメカニズムでデータが得られたかを考えるとこ
8708 ろにこそ、心理学者の本当の興味関心があるはずで。このデータは心のこういう仕組みを反映しているの
8709 ではないか。データの背後には変数同士のこういう関係が潜んでいるのではないか、ということを考えていく
8710 わけです。わからないところは確率で表現しつつ、変数間関係を関数で記述する。そのための設計図です。設

8711 計図と Stan があれば答え (事後分布) は得られるのですから、データ生成の仕組みを作り込むことに我々
8712 は全力を傾けることができます。複雑な関数関係をいきなり思いつくのが難しい、ということもあるでしょう。
8713 その場合はデータを可視化し、データに合うような関数はどんな形なのかを考える、あるいはいったん、単純
8714 な線形モデルを当てはめて、それを徐々に複雑にしていく、というやり方が良いでしょう。Kruschke (2014)
8715 は心理学者が書いたベイズ統計の本ですが、その後半は線形モデルです。一般化線形モデルや階層モデル
8716 について、とくに複雑な前置きなく「こういう仕組みになってるんだから、とりあえずそれを反映して、関係は線
8717 形でも大体うまくいくでしょ」という感じで話が進んでいきます。簡単なモデルから徐々に複雑にしていく、と
8718 いうのがモデリングの王道でもあります。完成品だけをみるととても複雑に思えますが、設計図片手に紐解い
8719 ていけば、意外とわかりやすいかもしれません。

8720 **モデリング**は自由に設計図を書くことができるので、**分散分析**や**t 検定**など平均値差だけにこだわってい
8721 る方法をみると、そちらがいかに窮屈そうにしていると思えるかもしれません。講義を通じて得た知識や技
8722 術で、自由にデータを調理していただければと思います。

8723 31.2 帰無仮説検定の代案

8724 ベイジアンモデリングは、自由に書いた設計図から**確率的プログラミング言語 (stochastic program-**
8725 **ming language)** をつかって答えを得る方法です。プログラミング言語が使えるようになるの大変です。こ
8726 れまでの一般的な統計ソフトでは、ボタン 1 つで答えを出してくれたのに、と思う人もいるかもしれません。し
8727 かしそれは、従来の分析方法がさまざまな仮定や前提に基づいた、マニュアル的アプローチを許す範囲のも
8728 のに限定されていたからですし、「簡単、楽ちん」ということが「知らなくていい」ということだとすると、その考
8729 え方がさまざまな問題を引き起こしてきたことを知らなければなりません。これは**心理学における再現性問題**
8730 (池田・平石, 2016) として取り沙汰されており、心理学そのものに疑問符を投げかけるような事態を招いて
8731 います。その原因の 1 つが、統計的技術の誤用・悪用にあります。

8732 心理学ではながらく、帰無仮説検定によって理論を積み重ねてきました。帰無仮説検定では p 値が 5% より
8733 小さければよし、という「ここだけ見ておけば良い」というようなマニュアル的基準があり、差があればすご
8734 い効果が発見されたぞ! と喧伝するというで進んできた側面があります。すでに述べたように、この**有意**
8735 **差**を出すために要因計画の方にこそ心理学者は注力してきたのですが、 $p < 0.05$ になりさえすれば良いと
8736 いう表面的な理解しかしていないマニュアル人間ができあがってしまうことが、問題を大きくしたのかもしれ
8737 ません。

8738 ベイズ統計によるアプローチの場合、データがどのように出てきたのか、事前分布はどのようなのか、といった
8739 ことを考えないわけにはいきません。自由に設計図を描くことができるというのは、裏を返せば、最初は白紙
8740 でしかないということであり、きちんと動く統計ツールを作るためには細部まですべて記述しなければならない
8741 のです。帰無仮説検定では、極端な場合「データ生成には正規分布が仮定されている」ということを知らな
8742 くても、 p 値がどうなるかだけをみる事はできてしまうのに、です。このように、ベイズ統計を使う利点の 1 つ
8743 は、前提や仮定についてすべて自覚的に、明示的に記述しなければならないというところにあり、これが誤用・
8744 誤用を生みにくくする安全装置となっている面があります。

8745 最近では JASP (JASP Team, 2021) という、GUI でベイズ分析もできるソフトウェアが出てきました。帰
8746 無仮説検定と同じような分析を、ベイズでやったらどうなるか、すぐに試すことができます。この便利なアプリ
8747 では Stan や JAGS などとは違って、どのような前提があるかを自覚しなくても、分析結果を得る事はできま
8748 す。ですから、今挙げた「安全装置」としてのベイズ統計の利点は、簡単に外れてしまう日が来るかもしれませ
8749 んし、その方がユーザにとってはいいかもしれません。この講義では第 20 講から第 23 講まで、帰無仮説検
8750 定の代わりにモデリングアプローチをする例を紹介してきました。さて、ではこうした帰無仮説検定の代わりに

8751 して、ベイズ統計を使う利点はどこにあるでしょうか。

8752 もちろんいくつか考えられますので、箇条書き的に解説してみたいと思います。

8753 ■**幾重にも組み合わせる仮定や補正からは無縁** たとえば二群の平均値の差を検定する場合、データが
8754 正規分布に従っているかどうかを検定し、また二群の分散が等しいかどうかを検定して、問題なければ t 検
8755 定に進みます。もし分散が等しくないということになれば、Welch の補正をしなければならない、とマニュアル
8756 にはあります。またたとえば、Within の ANOVA をするときも、球面性の検定を行なってデータが分析の前
8757 提となる仮定にあっていいるかどうかを判定しなければなりません。主効果が見られた、交互作用が見られたと
8758 なれば、今度はどこに差があるのかをみるために - 検定の繰り返しを避けるために、下位検定を行います。

8759 このように、仮説検定をする場合は条件にあっていいるかどうかを検定し、あつていなければ補正をするなど
8760 の手続きが必要です。同時に、仮説検定を同一のデータに何度も繰り返す事は、5% 水準という**タイプ 1 エ**
8761 **ラー**を生じる確率が適切に制御できなくなるため、本来なら実験単位ではなく「一連の分析単位」でこれを調
8762 整しなければなりません。正直なところ、ここまで厳密な調整が、あらゆる心理学研究できちんと行われてい
8763 るかと言われれば、その答えは No なのですが。

ベイズ統計の場合、こうした問題は非常に軽微なものに変わります。たとえば二群の分散が等しいかどうか
については、「等しくないモデル」を考えれば良いのです。具体的にいうと、二群の分散が等しいモデルは、

$$X_i \sim N(\mu_1, \sigma), Y_i \sim N(\mu_2, \sigma)$$

であり、等しくないモデルは

$$X_i \sim N(\mu_1, \sigma_1), Y_i \sim N(\mu_2, \sigma_2)$$

8764 とするだけです (σ に添字がある、すなわち違う大きさとして推定する)。球面性の検定についても分散共
8765 分散行列の要素それぞれを推定するだけであり、その数値に事前に定められた制約はありません。

8766 またたとえば、タイプ 1 エラーの制御についても、そもそもタイプ 1 エラーという発想がありませんから存在
8767 しないのです。データがあつて、平均構造を線形モデルで表現したのであれば、そこから出てくる事後分布は
8768 1 つしかありません。その単一のパラメータ同時事後分布をどのように切り分けようとも、ある判断のエラーに
8769 確率が伴う事はないのです。

8770 帰無仮説検定を厳密に行うには、下位検定を事前に計画してどこにどのような差があるか考えるかも準備
8771 しておく必要があります。帰無仮説検定は勝敗を決める手法ですから、ゲームのルールは事前に決めなけれ
8772 ばならないのです。ゲームの設定如何によっては、判定結果が変わってくることもあるからです。これに対し
8773 てベイズ統計的アプローチであれば、結果 (事後分布) は決まっているので、下位検定を事前に計画する必
8774 要はなく、事後分布をどう切り分けようと分析者の自由です。

8775 ■**停止規則やサンプルサイズの問題が生じない** 同様のことは、サンプルサイズの設計についても言えま
8776 す。帰無仮説検定の場合は、ゲームのルールを事前に決める必要があると言いました。このルールの中には、
8777 サンプルサイズも含まれています。サンプルサイズが大きくなればなるほど、有意であると判定される可能性
8778 は上がっていきます。帰無仮説検定における、最もやってはいけない研究実践のひとつは、データを取りなが
8779 ら検定を繰り返し、有意になったらデータを取るのをやめる、というものです。検定を繰り返すと、誤った判断
8780 をしてしまう可能性が増えますし、サンプルサイズが大きくなると有意になりやすくなります。つまり「勝つまで
8781 ゲームをやめない」というのは卑怯な方法なのです。

8782 しかし実際は、数人相手にデータをとってみたけど有意差はでなかった、でも惜しいから頑張ってデータを
8783 増やした、結果としてある日有意な差になったので、サンプルサイズをちゃんと記載して論文として提出した、
8784 ということがあるわけです。努力して真実に辿り着くのはいい話ですが、この例は努力したことが結果の捏造
8785 になっていたかもしれないという、より悲しい話になってしまいます。

8786 帰無仮説検定は事前にサンプルサイズを決めなければいけません。このサイズ、この基準で勝負する
8787 ぞ、というルールがあってからの一発勝負なのです。どこでデータの収集を止めるかというのは、**停止規則**
8788 **(Stopping Rule)** といって、これまた事前に決めておかなければならないことなのです。事後的に、勝負に勝つ
8789 たからこのやり方でよかったんだらう、と考えるのは間違っているのです。

8790 これに対して、ベイズモデルの場合はこうした問題が生じません。最後が確率的な判断にならないので、こ
8791 うした問題が生じないのです。データが蓄積されていくことは、徐々に確信度が高まっていくことでもありま
8792 す。差があるのかないのか、という判断をするにあたって、よりその違いが明確になっていだけなのです。

8793 ■より柔軟な仮説を立てることができる 帰無仮説検定は、帰無仮説と対立仮説という2つの仮説の比
8794 較判断になります。ここで帰無仮説は「差がない」「相関がない」といったものになります。たとえば二群間の母
8795 平均を比較する場合、帰無仮説は $\mu_1 = \mu_2$ になります。しかし実際には、差がないということを主張したい場
8796 合もあるかもしれません。帰無仮説に「差がある」というのをおくことができないのは、「あるというのは、どれ
8797 ぐらいあればいいのか」という問題がすぐに浮上するからです。たとえば母平均の差が $\mu_1 - \mu_2 = 0.1$ なの
8798 か、 $\mu_1 - \mu_2 = 0.11$ なのか、 $\mu_1 - \mu_2 = 0.111$ なのか…言い出すとキリがないですね。p 値は帰無仮説
8799 のもとで計算されるものですから、差があるという仮説は無限に作れてしまうので、検証すべき帰無仮説も無
8800 限に膨れ上がってしまいます。差がないという状態は $\mu_1 - \mu_2 = 0$ と一意に定まるので、帰無仮説を設定す
8801 ることが可能なのです。

8802 では帰無仮説を積極的に支持すればよいではないか、と思われるかもしれませんが、それができません。
8803 判定結果の p 値は帰無仮説のもとで計算していますから、「帰無仮説のもとで考えるとおかしな結果になっ
8804 た」という背理法が成り立たないのです。「帰無仮説のもとでおかしな結果にならなかった」というのは、帰無
8805 仮説以外の可能性を除外できるものではありません。

このように、帰無仮説検定では「効果がなかった」とか「関係がなかった」とは言えず、せいぜい「効果がな
いとは言えない」というにとどまるのでした。これに対して、ベイズ統計では**モデル比較**の観点から検証すること
になります。帰無仮説は $\mu_1 = \mu_2$ というモデル、対立仮説は $\mu_1 \neq \mu_2$ というモデルです。言い換えれば、帰
無仮説は

$$X_i \sim N(\mu_1, \sigma_1), Y_i \sim N(\mu_1, \sigma_2)$$

というモデルであり、対立仮説は

$$X_i \sim N(\mu_1, \sigma_1), Y_i \sim N(\mu_2, \sigma_2)$$

8806 というモデルだとしているに過ぎないからです。

8807 これまでみてきたように、ベイズ統計のアプローチで平均値の差を考える場合は、 $\mu_1 - \mu_2$ のような差の分
8808 布を**生成量**でつくり、その大きさをそのまま吟味できます。どうしても「差があるのか、ないのか」という判断
8809 がしたければ、**実質的に等価な範囲 (Region Of Practical. Equivalence; ROPE)**(→ セクション
8810 20.2, Pp.220) を事前に設定し、その区間に差の分布が入ってくるかどうかを判断するのでした。そもそも差
8811 も分布しているわけですから、どれぐらい重複しているかという幅で考えるほかないわけです。帰無仮説検定
8812 の文脈においても、**効果量 (Effect Size)** とよばれる**標準化された差**の大きさを算出して考えますが、この
8813 ように「どれぐらい違いがあれば差があると言えるか」という程度問題が前面に出てくるのが重要なのです。
8814 5% 水準のように機械的に判断するのではなく、領域固有の知識(ドメイン知識)に基づいて、実質的な差の
8815 大きさをケースバイケースで判断することが重要です。もちろん機械的に手続きを進められないことは不便か
8816 もしませんが、不便だからといって真実を曲げるようでは本末転倒です。

8817 また生成量を使ったアプローチのところで解説したように、パラメータについて考えるだけでなく、そのモデ
8818 ルから生成されるであろうデータについて仮説を立てたり検証したりできるのも、ベイズ統計的アプローチの

8819 利点です。実際どの程度の差があると意味があるのかについて、具体的な数値シミュレーションで考えること
 8820 ができるからです。パラメータが δ 以上の差があれば良いとか、二群のデータが D 以上違ってくる確率は
 8821 どれぐらいか、といった検証の仕方は、仮想空間上の p 値よりも具体的で意味のある情報をもっています。
 8822 パラメータやデータについて、「同じかどうか」以上の情報を検証できることは、ベイズ統計の長所といえるで
 8823 しょう。

8824 31.3 モデル比較

8825 先ほど、**モデル比較**ということに言及しました。このことについて、最後に少し触れておきたいと思います。
 8826 ベイズアンモデリング、すなわちモデルをたててベイズ統計学的に推定するなかで、そのモデルの優劣を考え
 8827 る方法があります。

8828 もう一度ベイズの公式を見てみましょう。ベイズの公式は次のようなものでした。

$$p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)}$$

8829 ここで右辺の分母は $p(D)$ 、つまりデータの確率をあらわすもので**周辺尤度 (marginal likelihood)** と
 8830 呼ばれます。これはどのように計算されるのでしょうか。尤度は $p(D|\theta)$ で表されますが、 $P(D)$ はこのパラ
 8831 メータが無くなったものになります。どのようにしてパラメータをなくすかという、パラメータのありそうな可
 8832 能性をすべて足し合わせることで**周辺化消去**すれば良い、ということになります。

たとえば性別と学年のクロス表があったとして、ランダムに選んだ 1 人の学生が女性である確率は、すべて
 の学年の女性の確率を足し合わせたものになるわけです。記号で考えると、条件付き確率 $P(D|\theta)$ におい
 て、 θ の可能性を全部考えれば $P(D)$ だけになるということです。これはパラメータ θ が離散的な例ですが、
 連続的な場合は

$$p(D) = \int p(D|\theta)d\theta$$

8833 のように積分で考えればよいでしょう。

8834 これはモデルに含まれるあらゆるパラメータのパターンを網羅して、データが得られる確率を計算したこと
 8835 になりますから、モデルが表現できる世界の中でどれぐらいデータを捕まえることができたか、ということを表
 8836 しているとも言えます。ここでモデルが複数ある場合を考えましょう。あるモデルを \mathcal{M} と表現することにして、
 8837 第一のモデルを \mathcal{M}_1 、第二のモデルを \mathcal{M}_2 、と表したとします。すると事後分布も「あるモデルを仮定した時の
 8838 事後分布」ということになりますから、次のようにすべて「モデルのもとでの」という条件付き確率で表現でき
 8839 ます。

$$p(\theta|D, \mathcal{M}_i) = \frac{p(D|\theta, \mathcal{M}_i)p(\theta|\mathcal{M}_i)}{p(D|\mathcal{M}_i)}$$

8840 ここで右辺の分母を見ると、 $p(D|\mathcal{M}_i)$ 、つまり「モデル i のもとでデータが得られる確率」を表していること
 8841 になります。もしこの確率が高ければ、このモデルはこのデータを生む確率が高いわけです。逆にもしこの確
 8842 率が低ければ、このモデルはこのデータを生む確率が低いということになります。この周辺尤度は、モデルが
 8843 データを支持する度合いであり、証拠 (**エビデンス**) の強さとも言えるわけです*1。

8844 その上で、次のような数字を考えます。

*1 これを考えると、どんなデータでも作れる幅広いモデルを考えておけばいいじゃないか、と思うかもしれませんが、モデルを広く考
 えずぎるとそこから出てくるデータはそのごく一部でしかなく、かえってエビデンスとしては弱くなってしまいます。なんでも予測で
 けるモデルは何も予測していないことと同じなのです。

$$BF_{12} = \frac{p(D|M_1)}{p(D|M_2)}$$

この数字は**ベイズファクター (Bayes Factor)** と呼ばれるもので、モデル 2 に比べモデル 1 がどれほどデータに支持されているかを示す数字です。この数字が 1.0 より大きければ、モデル 1 はモデル 2 に比べて相対的に強くデータに支持されているということになります。逆に 1.0 を下回ると、モデル 2 のほうがモデル 1 よりもいいね、ということです。モデル同士の比による表現ですので、 BF_{21} としても構いません。ここでモデル 1 が対立仮説、モデル 2 が帰無仮説だとすると、どちらがよりデータをうまく表しているか、どちらがよりデータに支持されているかがわかるわけです。最近では、帰無仮説検定の文脈でも BF を併記して、どちらのモデルが良いかを示すこともありますし、この方法ですと帰無仮説 (というか差がないというモデル) を積極的に支持する、ということもできます。

この手法はベイズ推定するものであれば一般的に通用する考え方ですので、1 つのデータに対していろいろなモデルを考えついたとしても、ベイズファクターを計算して相互に比較できるわけです。すなわち、ベイズのモデル比較はベイズファクターという共通の基盤を手にしたとも言えます。そんな優れた方法があるのであれば、なぜもっと早く教えてくれなかったのか、という声が聞こえてきそうです。実は理屈上ではこの通りなのですが、この方法には実践上の問題が残されているのです。それは、**周辺尤度の計算が難しい**というものです。すべてのパラメータについて、考えられるすべての可能性を計算して足し上げるわけですから、それはそれは大変な計算になります。有効な近似計算の方法として、全パラメータの全確率空間を計算し尽くすかわりに、確率が高そうなところを効率よくサンプリングして推定値を得る**ブリッジ・サンプリング (Bridge Sampling)** という手法が考えられており、その計算をする R のパッケージ ([Gronau, Singmann and Wagenmakers, 2020](#)) が提供されていますから、今後こうしたモデル比較も増えてくるでしょう。

またモデル評価の方法として、**情報量規準**によるアプローチも考えられています。これは「良いモデルとは良い予測をするモデルである」と考えるところから始まっており、平均的に良い予測をするモデル程度を指標化するものです。ベイズの定理を基本としながら、データを真のモデルから生成されるサンプルと考えたり、確率分布同士の近さを**カルバック・ライブラー情報量 (Kullback-Leibler divergence)** で表現するなど、非常に高度な数学的体系が確立されています。詳しくは[浜田・石田・清水 \(2019\)](#) など丁寧な解説書があり、R パッケージなどで簡単に計算できるようになっています ([Vehtari, Gelman and Gabry, 2017](#))。数学的にやや複雑な内容を含んでいますので、本講義の中では扱いきれませんが、より進んだ研究や利用のために言葉の紹介だけはしておきます。

31.4 おわりに

本講ではここまで、ベイズ統計と従来の統計的分析を比較して議論してきました。ベイズ統計の方がメリットが大きいような解説をしてきましたが、これは筆者の好みもかなり含まれるものです。帰無仮説検定はなんにも悪いものではなく、前提や仮定をしっかり守れば、一般的なソフトウェアで誰にでも結論を出すことができるという点で、とても良いものなのです。前提や仮定が多くて面倒だから、マニュアル人間ができてしまうというのは、あくまでもユーザ側の問題であって方法論の問題ではありません。ベイズ統計だって、モデルが間違っていれば結論は間違ったものになりますし、なによりデータに合わせてモデルを逐一設計図から書き起こすのは、大変コストがかかることでもあるからです。本来、これら 2 つの流儀は相反するものではなく、それぞれ依って立つ理論や考え方があって、必要に応じて選べるようになるのが一番です。

またベイズ流の研究アプローチ、とくにモデル比較については、まだまだ議論が活発に行われている段階であり、とくに心理学研究にこうしたアプローチがどれほど有用かについては、いまだに結論が出ていないところ

8882 でもあります。パラメータの推定値で考えるのが良いのか、モデル比較をして議論を進めるのかについても、
8883 まだまだ定番の方法というのは見つかっていません。

8884 そもそも心理学の知見を数式に落としてモデリングできるようにするためには、心理学そのものの理論的
8885 発展や、人間についてのより深い理解も必要です。人間の行動と予測が完全にできるようなモデルや一貫
8886 した理論体系は、まだまだ遙か先の目標に過ぎないというのが現状でしょう。それでも有用なツールを武器
8887 に、少しずつ知見を積み重ねていくことができれば、心理学をより楽しく学べるのではないかとおもいます。
8888 Enjoy!

8889 付録 A

8890 よくある質問とミスの例

8891 A.1 Frequently Miss and Comments

8892 Rmd でレポートを提出したのに、なんだか中身の問題じゃないのに突き返された、中身を見てくれよ！と
8893 思う人もいるかもしれません。テキストでは R や Rmd での課題を提出するよう求めているところがありますが
8894 が、その際よく見られる学生さんのミスとその対応についてのコメントをここにまとめておきます。

8895 A.1.1 FMC1；そもそもファイルの書式が違う

8896 Rmd で提出してください、R で提出してください、という指示に対して、違うものが提出されてくることが
8897 あります。書式があっていない、というのは些細なことのように思えるかもしれませんが、学術論文は書式に
8898 捉われず内容に集中するためにも、書式は整えられたものである必要があります。学会誌に掲載されている
8899 論文も、みなさんが書く卒業論文も、レポートに至るまで、書式や指示に沿ったものを準備する必要があります
8900 ず。書式があっていない場合は、門前払いになっても文句が言えないのがアカデミックの世界です。

8901 ということで、R や Rmd で提出してください、という指示があれば、R や Rmd で提出してくださ
8902 い。ファイルの種類は拡張子で分類され^{*1}、R ファイルは.R、Rmd ファイルは.Rmd という拡張子になってい
8903 ます。たまたま.Rproj というファイルを提出してくる人がいますが、これは R プロジェクトのファイルで、これに
8904 は R スクリプトも文書も含まれておらず、みなさんの計算環境情報が少し書いてあるだけです。また、.Rmd
8905 だけ入っていれば良いかというところではありません。まれに Filename.Rmd.R というようなファイルを送っ
8906 てくる人がいます。これはファイル名にピリオドが含まれているだけで、ファイルの種類を識別する拡張子は.R
8907 ですから Rmd ファイルではありません^{*2}。そもそもファイル名には 2 バイト文字はもちろん、!や?などの記号
8908 を含めるべきではありません。ピリオドも当然記号の一種ですから、ファイル名にするのは不適切です^{*3}。

8909 ではどうやって適切なファイル形式にするか、ということですが、最も素直な方法としては RStudio で新し
8910 くファイルを開くときに、R markdown 形式 (略して Rmd) を選ぶようにしましょう。もし間違っ、R script
8911 形式 (.R 形式) で開いてしまった、というときは、そのファイルを破棄して新しく作り直すのが一番ですが、エ
8912 ディタペインの右下にあるファイル種類表示 (図 A.1) をクリックして修正することもできます。

*1 拡張子についてはセクション C.5, Pp.386 参照

*2 最近の OS は拡張子を表示しない設定になっているものも少なくないので、このようなミスが生じます。

*3 長いファイル名などの場合、空白を入れるのも適切ではありません。できなくはないのですが、やるべきではないのです。どうしても空白を入れたければ、アンダースコアなどで区切ると良いでしょう。

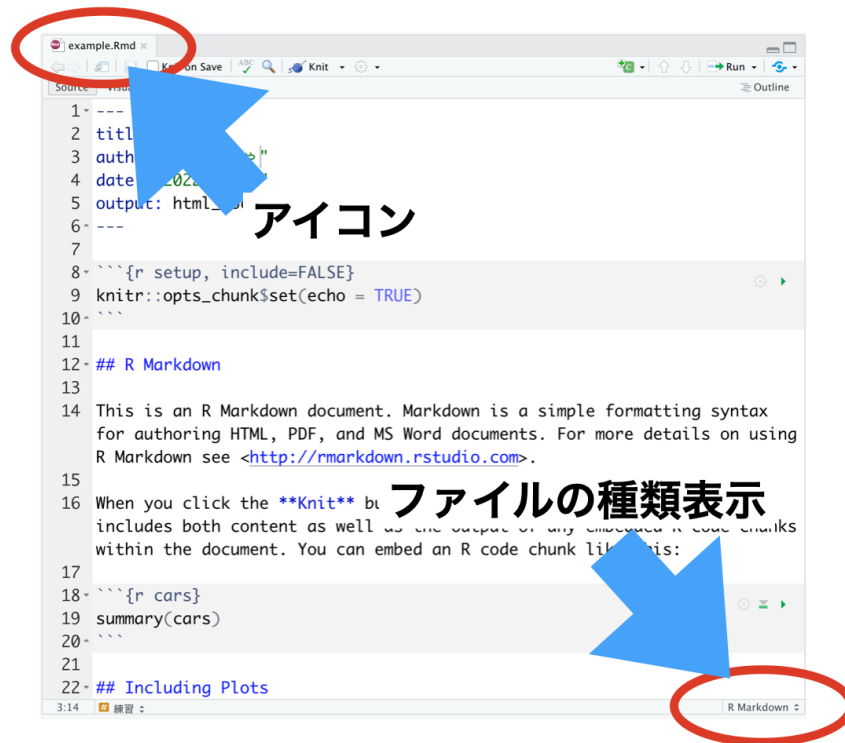


図 A.1 ファイルの種類を判別する方法

8913 A.1.2 FMC2 ; やっぱりファイルの形式が違う

8914 Rmd 形式のファイルは、拡張子だけで決まるものではありません。図 A.2 にあるように、Rmd ファイルは
 8915 冒頭の 6 行 (正確には --- で囲まれた領域) が YAML と呼ばれるところで、文書全体の設定をしています。
 8916 その下に、R のコードを実行する部分 (チャンク (chunk)) や文章の領域があります。文章のところは # 記
 8917 号で見出しを作ったりできます。

8918 この YAML 部分が壊れている、あるいはチャンクが正しく記述されていない場合、拡張子が .Rmd であっ
 8919 ても適切な Rmd ファイルにはなっていません。YAML 部分の書き方はよくわからない、という人も多いと
 8920 思いますので、RStudio で Rmd ファイルを作ったときの状態をなるべく変更しないように注意すると良いで
 8921 しょう*⁴。

8922 チャンクは R のコードを書くところで、バッククォーテーション 3 つでくくるのが決まりです。チャンク領域が
 8923 始まるところに {r} とかいて「ここが R で計算するところです」というのを指定するわけです*⁵。このバック
 8924 クォーテーション 3 つ (```) が全角だったり (` ` `), ダブルクォーテーションだったり (" ") すると、機械は正しく
 8925 チャンクであると認識しません。フォントなどの見せかけ上は、微妙な違いのように見えますが、機械にとって
 8926 違う文字列は違う意味を持ちますので注意してください*⁶。RStudio で編集する場合、チャンクの領域はや

*⁴ Rmd ファイルを新しく開くときに、文書タイトルや著者名、日付、出力ファイル形式などを設定することができるウィンドウが開きますので、そこに必要な情報を書くことで自動的にそれを使った YAML が生成されます。また、本文としていくつかのサンプルコードが最初から含まれていますが、これらに関してはサンプルをそのまま使うことはないので、全て削除してしまっても構いません。

*⁵ ほかに Python や Julia など他の計算言語を混ぜることもできます。

*⁶ 記号の名称や入力については、セクション E, Pp.395 を参照してください。

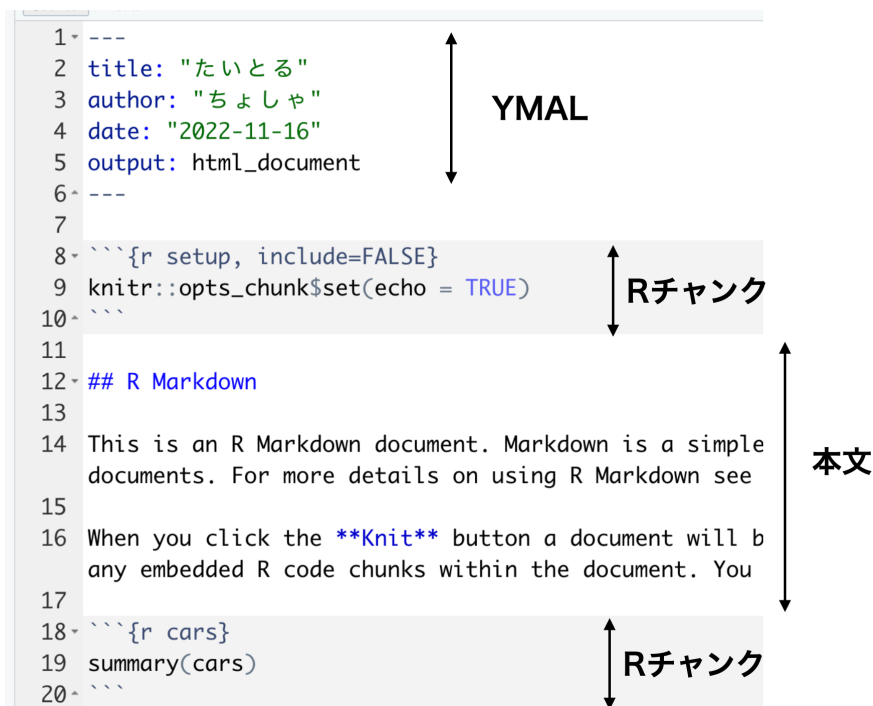


図 A.2 Rmd ファイルの中身における書式

8927 や灰色がかった強調表示がされますので、どこにチャンクがあるかわかると思います。もし強調表示されてい
 8928 ないようであれば、チャンクとして認識されていない可能性を疑った方が良いでしょう。

8929 チャンクはバッククォーテーション 3 つで開き、おわたら同じくバッククォーテーション 3 つで領域を閉じ
 8930 ます。閉じなければずっと R の計算領域が続くと解釈されますし、最後までチャンクが閉じられないと Rmd
 8931 ファイルとして正しくない書式ということになります。チャンクが正しく入力できているかについて、常に注意を
 8932 払っておく必要があります。また、自分でバッククォーテーション 3 つを使ってチャンク領域を開いたり閉じたり
 8933 するのが面倒だ、という人は RStudio のチャンク挿入ボタンから挿入すると間違いないでしょう。

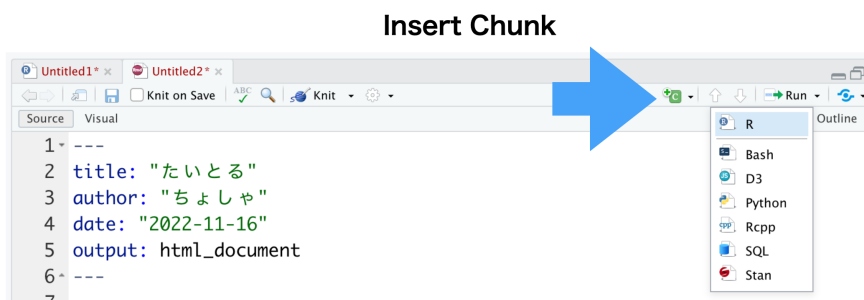


図 A.3 RStudio のチャンク挿入ボタン

8934 A.1.3 FMC3 ; Rmd が knit できない

8935 Rmd は文書作成に加えて、R での計算がセットになったファイル形式であり、文中の数字や分析結果
 8936 は「書き出す」のではなく「その場で生成する」ものです。生成する、というのは Rmd ファイルを変換して、

8937 HTML や DOCX, PDF 形式のファイルにすることを指します*7。このファイル変換をニット (knit) といい
 8938 ます。編み物を編むようなイメージですね。このときに、タイトルをつけ、見出しのサイズを変え、R の計算を
 8939 して結果を埋め込む作業をするわけです (図 A.4)。こうすることで、結果のコピペを避けること、再現性を担
 保することができるようになるわけです。すでに説明したように、チャンクを使って計算に必要な指示を Rmd

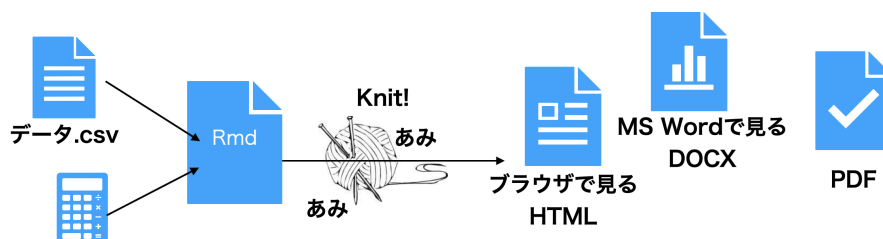


図 A.4 knit してレポートが完成する

8940

8941 ファイルに書きます。この指示はエラーのないコードである必要があります。当然ですが、スペルミスや間違っ
 8942 た R コードは「実行できない」というエラーになるわけです。その場合、knit は中断され結果のファイルが出
 8943 力されません。提出されたレポートは knit して出来上がったもののことを指しますから、knit できない Rmd
 8944 ファイルを提出されてもレポートが提出されたことにはならないわけです。

8945 knit できない Rmd ファイルになっていないか、というのはご自身の RStudio 環境で knit して確認でき
 8946 ることです。提出前に一度、きちんと機能する Rmd ファイルになっているかどうか確認するようにしてくださ
 8947 い。以下で述べるように、自分の環境で knit できても、提出先の (私の) 環境で knit できないということも
 8948 あり得ます。しかし、自分の環境で knit できないのに提出先でできる、ということはありませんので、「knit
 8949 できませんよ」というコメントをつけて返される前に自分で確認するようにしてください。

8950 A.1.4 FMC4 ; 外部環境を参照してしまう

8951 さて、R で行う作業の中には、csv ファイルを読み込むといった「外部のファイルを使う」指示もありえま
 8952 す。例えば次のようなコードです。

code : A.1 Rmd 中の R コードが外部環境を参照してしまう

8953

```
8954 1 dat <- read.csv("social_effects.csv")
8955 2 dat <- read.table("clipboard")
8956
```

8957 このコードでは、上の行は social_effects.csv というファイルを読み込むように、下の行はクリップボードの内容を
 8958 読み込むようになっています。ファイルを読み込む指示は、ファイルがなければ当然エラーになりますから注意
 8959 が必要です。レポート等で Rmd ファイルを提出するとき、Rmd ファイルの他に必要な読み込むべきファイル
 8960 があれば一緒に提出するようにしてください。また、クリップボードの内容は再現できません。クリップボードと
 8961 は、コピーアンドペーストのコピーを行ったとき、PC の内部で一時的にコピー内容を覚えておく場所のこと
 8962 です。つまり、みなさんがみなさんの環境で、クリップボード上にデータを一旦保持している場合は、このコード
 8963 でエラーがしょうじることはありません。しかし、レポート提出先の (私の) 環境で、事前にデータのコピー作業
 8964 を行なっているわけではないのですから、提出先でエラーが発生します。

*7 HTML は Hyper Text Markup Language の略で、ブラウザで開くファイル形式です。DOCX は Microsoft Word の
 ファイル形式です。PDF は Portable Document Format の略で、データを紙に印刷した状態のようにサイズを固定して出
 力したファイル形式であり、OS がもつビューワーや Adobe Acrobat Reader などで見ることができます。

8965 そもそも Rmd ファイルは、作業の再現性を担保するためのファイルになっているわけですから、クリップ
8966 ボードの利用のような「自分の環境だけで可能な記録されない作業」をそのファイルに含めるのは適切ではあ
8967 りません。Rmd ファイルは Rmd ファイルだけで分析作業が完結するように、必要な記録は全て記載されて
8968 いる必要があります。同様に、分析作業に関係のない冗長な指示や無駄な指示は Rmd ファイルに含めるべ
8969 きではありません。

8970 A.1.5 FMC5 ; 提出先の環境を変更してしまう

8971 R はさまざまなパッケージを使って分析環境を拡張することができます。パッケージは **CRAN** を通じて
8972 インターネット経由で配布されますから、ネット環境があれば誰でも最新のパッケージをとってくることができ
8973 ます。みなさんが Rmd ファイルの中で使う R の関数の中には、パッケージの関数も含まれているでしょう。
8974 パッケージがなければ関数が動きませんから、パッケージをインストールする作業も Rmd に書いておきたい
8975 と思うかもしれません。しかし、これは推奨できません。

8976 パッケージのインストールは、実行環境の準備にあたる作業です。Rmd ファイルを使ってコードのやり取り
8977 をするとき、提出先の環境で分析することになりますが、提出先の環境にどのようなパッケージをいつ入れる
8978 かは、提出先の環境の管理者が判断すべき問題です。提出する Rmd ファイルにパッケージをインストールす
8979 る関数、すなわち `install.packages()` 関数が含まれているというのは、相手の環境を勝手に操作してし
8980 まうことと同じであり、セキュリティ的にも適切な発想ではありません。

8981 環境の準備は提出先の管理者が管理すべき問題であり、またすでにパッケージが入っている場合は無駄
8982 な書き作業をさせることにもなります。また `install.package` 関数は CRAN のサーバを参照したりしま
8983 すから、適切な設定がなければ R チャンク実行時にエラーが発生します。いずれにせよ、R チャンクの中に
8984 `install.package` 関数を含めないようにしましょう。

8985 以上が Rmd ファイルや R ファイルでレポートを提出するときの留意点です。その他にも R や RStudio を
8986 使うときによくある質問がありますので、それらも一問一答型で紹介しておきましょう。

8987 A.2 Frequently Asked Questions ; よくある質問と答え

Q. A.2.1: テキストを参考にパッケージをインストールしようとしたところ、エラーが発生しました。

A. 質問ではなくて報告ですね。お返事としては、「わかりました。エラーが発生して大変ですね。」としか言いようがありません。どの環境で、何をしようとして、どのようなエラーが出たのか、明示してください。メールのやり取りで指示が明確になるように、テキストを準備しています。何章、何ページの、どの文章を参考にしたのかも教えてください。

8988

Q. A.2.2: 添付のようなエラーが発生しました (図 A.5)。対応策を教えてください。

A. エラーの発生画面を送ってこられていますが、この画面に写っていない上の方でエラーが発生していますから、これではどのエラーなのかわかりません。スクリーンショットを送ってこられることはよくありますが、ほとんどの場合、適切な箇所が写っていません。複数枚添付してこられる人もいますが、画面を拡大しながら読むのも難しいので、**R コードそのものを送ってください**。そうすると、何行目のどこにエラーがあるかが明確になり、対応に関係ない無駄なやり取り (ex. 「もう少し上を写してください」「違う、もっと上」) が減ります。

8989

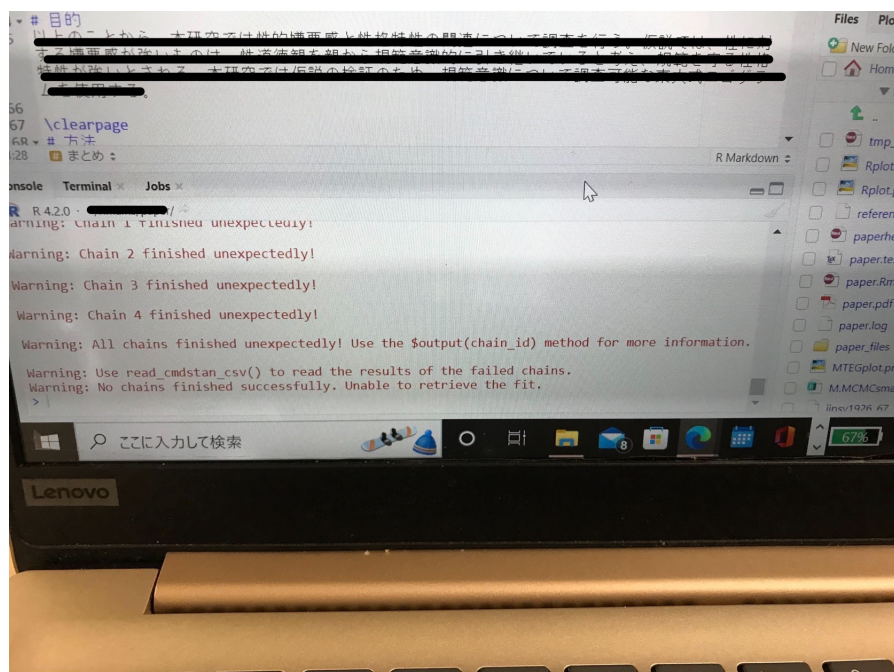


図 A.5 添付されてくるスクリーンショットの一例

Q. A.2.3: ファイルを読み込もうとすると次のようなエラーが出ます。どうすれば良いですか。

'xxxxx' に不正なマルチバイト文字があります

A. 文字化けの原因は、たとえば UTF-8 形式で提供されているファイルを、Windows 標準の CP932 形式で読み込もうとする、という文字コードの不一致です。ですからそれをオプションで指定してあげれば問題解決です。read.csv 関数を使っている場合、オプションの指定は read.csv("foo.csv", fileEncoding="UTF-8") のようにします。tidyverse の read_csv 関数を使っている場合、locale オプションで、locale 関数の encoding を明示的に UTF-8 とします。read_csv("foo.csv", locale=locale(encoding = "UTF-8")) のようにします。

8990

Q. A.2.4: ファイルを読み込もうとすると次のようなエラーが出ます。どうすれば良いですか。

'xxxxx' に不正なマルチバイト文字があります

A. 文字化けの原因は、たとえば UTF-8 形式で提供されているファイルを、Windows 標準の CP932 形式で読み込もうとする、という文字コードの不一致です。ですからそれをオプションで指定してあげれば問題解決です。read.csv 関数を使っている場合、オプションの指定は read.csv("foo.csv",fileEncoding="UTF-8") のようにします。tidyverse の read_csv 関数を使っている場合、locale オプションで、locale 関数の encoding を明示的に UTF-8 とします。read_csv("foo.csv", locale=locale(encoding = "UTF-8")) のようにします。

8991

Q. A.2.5: パッケージを読み込もうとすると次のようなエラーが出ます。どうすれば良いですか。

library(tidyverse) でエラー ; 'tidyverse' というパッケージはありません。

A. パッケージがないので、インストールしてください。

8992

Q. A.2.6: パッケージをインストールしようとする次のようなエラーが出ます。どうすれば良いですか。

Warning in install.packages:

```
'lib = "C:/Program Files/R/R-4.0.5/library"' is not writable
```

A. ユーザがファイルに書き込みをする権限を持っていないので、(パッケージファイルをドライブに) 書き込みできないというエラーです。RStudio を実行する時に、「管理者として実行」を選びましょう。

8993

Q. A.2.7: パッケージをインストールしようとするときのようなエラーが出ます。どうすれば良いですか。

Warning in install.packages:

ディレクトリ 'C:~(任意の文字列)~\????' を作成できません。理由は ``'Invalid argument' です。

A. ユーザ名が全角文字を含んでいるため、文字化けして R 側から操作ができません。解決する方法は二つあります。

ユーザを作り直す方法 ユーザ名に全角文字を含まない、新しいユーザを作成します。設定 > アカウントから新しいアカウントを作りましょう。

インストールフォルダを指定する R のコンソールで `.libPaths()` と実行すると、パッケージをインストールするフォルダが出てきますが、ここを変更する方法です。次の 2 ステップで対応できます。

- まず C ドライブのすぐ下にインストール先のフォルダを作ります。myLib としましょう。
- 次に R のコンソールで `.libPaths("C:/myLib")` 書いて実行します。

これでインストール先が変わりますので、書き込み・インストールができるようになります。老婆心ながら付け加えますと、フォルダ名はなんでも構いませんが、全角文字ではいけません。また場所も C ドライブのすぐ下である必要はありませんが、全角文字や空白を含むフォルダの下に入れてしまってはいけません。OneDrive のようなクラウドを指定するのも良くありません (探しに行った時にオフラインだとまたエラーになります)。

Q. A.2.8: ファイルが読み込めません

```
file(filename, "r", encoding = encoding) でエラー:
```

コネクションを開くことができません

追加情報: 警告メッセージ:

```
file(filename, "r", encoding = encoding) で:
```

ファイル 'foo.csv' を開くことができません: No such file or directory

A. 指定された場所にファイルがないので、読み込むことができないエラーです。確認すべきは、「プロジェクトを開いているか」、「プロジェクトフォルダの中に当該ファイルはあるか」、「ファイル名のミススペルはないか」です。プロジェクトってなんだという人は、RStudio の基本に立ち戻り、プロジェクトでファイルやフォルダを管理するようにしてください。プロジェクト管理については、Pp.87にもその説明があります。

プロジェクトによる管理とは要するに、R が今見ているフォルダの場所を固定する方法です。プロジェクトを開くと、プロジェクトフォルダが「今見ているフォルダ (ワーキングディレクトリ)」になりますので、その中のファイルを参照することになります。プロジェクトフォルダの中に当該ファイルがないと読み込むことができませんので^a、当該ファイルをプロジェクトフォルダ内に移してきてください。

^a フォルダの位置を相対的・絶対的に指定してやれば、どこにおいても読み込むことはできますが、この問題で悩んでいる人は相対・絶対パスの指定というところでさらに疑問が深まることになると思いますので、気にしないでくれて結構です。相対・絶対パスが知りたい人は、付録 C,389 でファイル場所とは何かを再確認してください。

8995

Q. A.2.9: ANOVA 君が読み込めません

```
file(filename, "r", encoding = encoding) でエラー:
```

コネクションを開くことができません

追加情報: 警告メッセージ:

```
file(filename, "r", encoding = encoding) で:
```

ファイル 'http://riseki.php.xdomain.jp/index.php?plugin=attach&refer=ANOVA 君&openfile=anovakun_486.txt' を開くことができません: Invalid argument

A. ファイルを読み込みに行く先が URL、すなわちインターネット上になっています。ANOVA 君のファイルを一度手元の PC のフォルダにダウンロードし、そのローカルのファイルの位置を指定して読み込むようにしてください。

8996

Q. A.2.10: 効果量として Hedges の g を算出してください。という指示について実行すると g ではなく d が出てしまうようなのですが、どうすればよいでしょうか。コードは次の通りです。

```
cohen.d(value ~ condition, data = dat, hedges.crrrection = T)
```

A. Hedges の補正 (correction) のオプションが通っていません。スペルミスです。オプションのスペルが間違っているので無視されたので、関数名通り Cohen の d が算出されます。

8997

Q. A.2.11: 描画の際に次のような注意が出てきます。これはどういう意味ですか。

```
Removed 671 rows containing missing values (geom_point).
```

A. 「671 件の行で欠損値が含まれています」ということです。つまりデータセットの中に欠損値 (観測されていない, 数値が入っていない行) が 671 件あったので、それは表示できませんでしたよ、という意味です。警告が嫌だということであれば、データセットの中で欠けているものを除外する必要があります。R の関数では、`na.omit` で欠損値を除外することができます。

8998

Q. A.2.12: `lm` 関数を実行するコードでオブジェクトがないと言われます。

```
result <- lm( weight - height, data = dat)
```

A. 従属変数と独立変数とをつなげるのはチルダ (`~`) という記号です。そこがハイフン (`-`) になっています。スペルミスの一種です。

8999

Q. A.2.13: `lm` 関数を実行するコードでオブジェクトがないと言われます。

```
result <- lm( weight - height )
```

A. データを与えていないので、`weight` というオブジェクトを探しに行くと、見つからないというエラーです。`data` オプションでデータを与えてあげてください。

9000

Q. A.2.14: 因子分析の Robust 法での RMSEA の p 値って超えてたらまずいですか。Standard(DWLS) の方だけで報告はだいじょうぶでしょうか。

A. 非常に専門的な質問をしておられますが、まず、「因子分析」「Robust 法」「RMSEA」など、個々の用語の意味はわかっているでしょうか。キーワードによる検索で、「ロバスト法とは」「Robust の意味」などは答えが出てくるとは思いますが、これらを分析法の体系的な文脈の中に位置付けないと、答えられない種類の問題です。この例をもとに説明すると、Robust を辞書で引くと「壮健」「たくましいこと」など出てきますが、もちろんそういう意味ではありません。ロバスト法を検索すると、「統計学の分野でロバスト推定法というやり方がある」「観測値に外れ値が含まれている可能性を考え、その影響を抑えることを目的とした手法」などの説明が出てきます。ロバスト推定法は因子分析に限らず、回帰分析など他の手法でも使われる考え方なので、このような一般的な解説になります。ですがここで知りたいのは「因子分析におけるロバスト推定」ですから、因子分析でロバスト推定とはどういうことか、因子分析の観測値における外れ値とは何か、因子分析における推定とは何を推定するのか、そもそも因子分析とは何を目的としているのか、なぜ自分は因子分析方を使うのか、さらに何故因子分析の中のロバスト推定を使いたいと思っているのか、といったことがわかっていないと、この質問に正しく回答することができません。このように、複合的な要素について一回で質問しても、適切な答えに辿り着けないことがあります。聞くことは恥ずかしいことではありません。知らないことを知っているふりをすることが恥ずかしいことであり、知らないまま「みんなそうやっているから」「テキスト/ネットに書いてあったから」と看過してしまうことこそ恥ずべきことなのです。ちなみに、質問に対する回答を間違えることも恥ずかしいことではありませんから、「正しく答えられなければ恥をかく(から質問しない)」というのも同様に恥ずべきことです。こうした恥ずかしさ(保身)から、「専門用語を使ってそれっぽく質問してわかった感じになろう」というのは、かえって遠回りになります。実は回答よりも、質問の仕方です。その人の理解度が明らかになってしまっています。このように書くと、なんでも反射的に「わかりません、教えてください」という人もいますが、それも適切な質問方法ではありません。教員がアドバイスできるのは、「答え」ではなく「理解」が欲しい人に対してだけなのです。質問する場合は、自分は何がわかっていて何がわかっていないのか、何が知りたいのかを明確に言語化して質問するようにしてください。

付録 B

標準正規分布から尺度値を求める計算方法

Likert 法では、態度が標準正規分布すると仮定するのです。標準正規分布をカテゴリの相対度数で分割し、あるカテゴリ c の上限の確率点 z_c 、下限の確率点 z_{c-1} の確率密度の差分を、相対度数 p_c で割ることで、尺度値 Z_c が下の式で得られます。

$$Z_c = \frac{(y_{z_{c-1}} - y_{z_c})}{p_c}$$

ここで y_{z_c} は標準正規分布をカテゴリで区分し、当該カテゴリ c までの累積確率点 z_c における確率密度、 p_c はカテゴリ c の相対頻度です。図 B に記号の対応関係を示しましたので、確認してください。

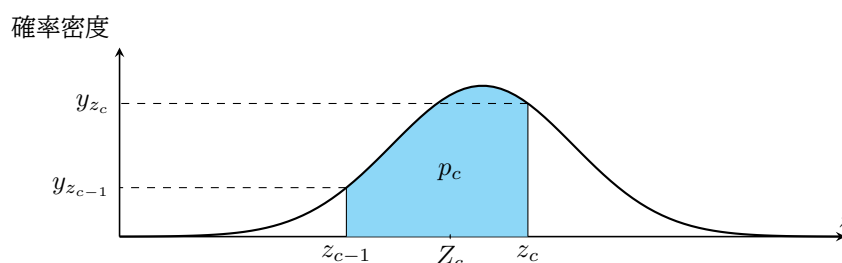


図 B.1 標準正規分布と対応する記号の確認

尺度値 Z_c を求める計算が確率密度 $y_{z_{c-1}}, y_{z_c}$ と、相対度数 p_c で算出されるというのは一見奇妙です。どうしてこのようになるのかを見ていきましょう。

まず標準正規分布の確率密度の式を確認しておきます。確率点 z における確率密度 y は次の式で算出できます^{*1}。

$$y_z = f(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

この関数は確率密度の曲線を表しており、確率はその面積です。 z_{c-1} から z_c までの面積（確率）は、積分を使って

$$p_c = \int_{z_{c-1}}^{z_c} f(z) dz = \int_{z_{c-1}}^{z_c} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

*1 e^x を $\exp(x)$ と書くことも少なくありませんし、この e は自然対数の底で $e = 2.718\dots$ という実数です。この数字は微分しても変わらない、 $(e^x)' = e^x$ という便利な特徴を持っています。

9016 で表されます。

9017 さて、ある確率点 z における密度の高さを $f(z)$ としたとき、 z の微小な増分 Δz を考えると、区間の面積
9018 $S = f(z)\Delta z$ を考えることができますから、逆にある点 z が知りたい時は

$$z = \frac{\sum S_z}{\sum S}$$

9019 とすれば良いこととなります。積分はこの微増分 Δz の極限を

$$\lim_{\Delta z \rightarrow 0} \Delta z = dz$$

9020 と考えることですから、ここで考えたいのは

$$Z_c = \frac{\int_{z_{c-1}}^{z_c} f(z)z dz}{\int_{z_{c-1}}^{z_c} f(z) dz}$$

9021 になります。分母は確率密度関数の積分ですから面積すなわち確率で、ここでは p_c であり、その面積を実
9022 データの相対度数で代えてもいいでしょう。

9023 分子については少し式の変形が必要です。記号を見やすくするために、 $a = z_{c-1}$, $b = z_c$ と書き換えてお
9024 きましょう。

$$\begin{aligned} \int_a^b yz dz &= \int_a^b \frac{1}{\sqrt{2\pi}} e^{(-\frac{z^2}{2})} z dz \\ &= \frac{1}{\sqrt{2\pi}} \int_a^b e^{(-\frac{z^2}{2})} z dz \end{aligned}$$

9025 ここで変数変換をして、 $u = -\frac{z^2}{2}$ とすると

$$\begin{aligned} \frac{du}{dz} &= -z \\ du &= -z dz \end{aligned}$$

9026 となり、積分の下限は $-a^2/2$ 、上限は $-b^2/2$ になりますから、以下のように展開できます。

$$(\text{与式}) = -\frac{1}{\sqrt{2\pi}} \int_{-a^2/2}^{-b^2/2} e^u du$$

$$\int e^u du \text{ は } e^u \text{ なので}$$

$$= -\frac{1}{\sqrt{2\pi}} [e^u]_{-a^2/2}^{-b^2/2}$$

$$= -\frac{1}{\sqrt{2\pi}} e^{(-\frac{b^2}{2})} - \left(-\frac{1}{\sqrt{2\pi}} e^{(-\frac{a^2}{2})} \right)$$

$$= \frac{1}{\sqrt{2\pi}} e^{(-\frac{a^2}{2})} - \frac{1}{\sqrt{2\pi}} e^{(-\frac{b^2}{2})}$$

$$y = \frac{1}{\sqrt{2\pi}} e^{(-\frac{z^2}{2})} \text{ であることから,}$$

$$= y_a - y_b$$

$$= y_{z_{c-1}} - y_{z_c}$$

9027 以上のことから、リッカート法において標準正規分布をもとに尺度得点を決めるには、

$$Z_c = \frac{(y_{z_{c-1}} - y_{z_c})}{\int_{z_{c-1}}^{z_c} f(z)dz} = \frac{(y_{z_{c-1}} - y_{z_c})}{p_c}$$

9028 とすれば良いことになります。

9029 この計算に至る理論的背景は、より専門的には**系列範疇法 (Method of Successive Categories)**

9030 と呼ばれ、順序カテゴリに数値を付与する心理測定論、精神物理学理論からきています。詳しくは [Guilford](#)

9031 (1954) や [西村 \(1977\)](#) も参照してください。

9032 付録 C

9033 電子計算機のイロハ

9034 C.1 前置き

9035 このセクションは心理統計ではなく、コンピュータについての四方山話をダラダラと書いています。そんなの
9036 聞かなくてもわかってるよ、という人もいるかもしれませんが、知らなくてもみなさんはきっとスマートフォンや
9037 タブレットを使っていることと思います。しかし知らずに使うことと、知ってて使うことには大きな違いがありま
9038 すし、今後大学でレポートや論文を書いたり、それに必要な統計処理をするためにも、計算機の基本的な特
9039 徴を知っておくと、トラブルに会った時に「ああこれってひょっとして」というヒントが得られたり、納得できる
9040 ようになるかもしれません。知らなければ「何だかわからないけどパソコンが壊れた」というか、「パソコン運が
9041 悪い」「自分はパソコンが苦手なのだ」と間違えた帰属をしてしまうことになります。

9042 コンピュータ関係の授業で聞いたことがある話、これから聞く話もあると思いますが、もし苦手意識を持っ
9043 ている人がいたらこれを機に再入門するつもりで読んでください。

9044 C.2 コンピュータの基礎

9045 21 世紀に生きる私たちは身の回りをコンピュータに囲まれて生きています*1。それは携帯電話の形をして
9046 いたり、ノートパソコン、デスクトップパソコンの形をしていたりします。また時々テレビなどで報道されますが、
9047 気象予報や飛沫がどのように飛び散るかをシミュレーションする大型計算機「富嶽」などもコンピュータです
9048 ね。これらは形は違いますが、いずれも電子計算機であり、電子計算機には次の 5 つの装置があります。

9049 **入力装置** キーボードやマウス、タッチパネルなどを使って情報を取り込むデバイス (装置)

9050 **出力装置** モニタやプリンタ、タッチパネルなどを經由して情報を出力するデバイス

9051 **演算装置** プログラムの命令に従って計算処理 (四則演算や論理演算) をする装置。一般に電子計算機の
9052 中央で一括して処理するので、Central Processing Unit(CPU) と呼ばれるものです。

9053 **制御装置** 演算結果に従って他の装置に指示を出す装置のこと。演算装置とまとめて CPU に実装されてい
9054 ます。

9055 **記憶装置** 計算結果などの情報をいったん保持しておく装置のこと。コンピュータの内部にあって一時的な

*1 私は 1976 年生まれですが、生まれた頃は周りにコンピュータなんかありませんでした。小学生の頃にマイコン (マイクロコン
ピュータ、小さなコンピュータという意味でもあります、My Computer、私のコンピュータという意味でもあります。つまり個
人単位でコンピュータが使えるようになった、というだけでも大きな出来事だったのです。) という言葉が出てきて、なんかかっこ
いいなと思った記憶があります。私が 10 歳になったころ、ビデオゲーム (テレビゲームでやるゲームが家庭でできるようになり、
それをこのように呼びました。) が身の回りに出てきました。ファミリーコンピュータ、とくに「スーパーマリオブラザーズ」によって日
本中の子供たちが熱狂したのが 11 歳の頃、「ドラゴンクエスト」によって社会問題になったのが 12 歳の頃になります。ともかくこ
の頃は、コンピュータといってもゲーム機のような扱いでした。

9056 計算に使われる一次記憶装置, ハードディスクドライブ (Hard Disk Drive,HDD) やソリッドステート
9057 ドライブ (Solid State Drive,SSD) などの二次記憶装置など。

9058 最後の記憶装置については, 一次記憶装置, 二次記憶装置と種類が分かれていますがこの区別は簡単
9059 で, 電源を落とした時に記憶が消えてしまうのが一次記憶装置, 電源を落としても記憶が消えないのが二次
9060 記憶装置です。たとえば $12 + 38 =$ という計算をするとき, 頭の中で「えーっと一の位が 2 と 8 だから 10 に
9061 なって繰り上がるから・・・」と考えてから, ノートに $= 50$ という答えを書くとします。次の問題に進むと, 先
9062 ほどの「1 繰り上がるから・・・」という情報は忘れてますよね。でもノートに書いた $12 + 38 = 50$ というのは
9063 残っています。このノートがいわば二次記憶装置であり, 頭の中で一時的に保持していた情報が一次記憶装
9064 置ということになります。一次記憶装置は RAM(Random Access Memory) とも呼ばれます*2。

9065 とところで, コンピュータがやっているのは計算だけです。私はこの資料を PC に向かって書いており, キー
9066 ボード (入力装置) を叩きながら, 画面 (出力装置) を見て文字を連ねています。これも「キーが押されたら文
9067 字を表示させ, その文字列を記録する」という処理を機械が淡々とこなしているに過ぎません。あるいはマウ
9068 スやトラックパッドで, アイコンを指し示し, カチリと押す*3ことで選択し, 押し込んだまま移動させ (ドラッグ),
9069 離すことで別の場所に置いたりします (ドロップ)。トラックパッドの場合は 2 本指で同時に押しついたりし
9070 タッチパネルの場合は二本指を広げたり (ピンチアウト), 逆に二本指を狭めたり (ピンチイン), 3 本以上の指
9071 でファサーっと触って場所を広げたりします。たとえば「ファイルを掴んでゴミ箱に捨てる (削除する)」という
9072 が我々にとっての操作ですが, コンピュータの内部では実はこんなことをしていません。ファイルは (二次) 記
9073 憶装置に書き込まれた情報です。記憶装置は原稿用紙のように小さなマス目がたくさんあって, ファイルとは
9074 そのマス目の XXX 番目から YYY 番目までの情報, ということです。このファイルを削除するというのは,
9075 記憶装置のある場所 (アドレス) に「削除されたものなので画面に表示しない」という情報を書き込む, という
9076 操作をしているだけです。じゃあなぜ私たちは「ゴミ箱にドラッグ&ドロップ」なんてするのでしょうか? それは
9077 そのほうがわかりやすいからですよ。「ファイルを削除する」というのは, 「メモリアドレスの XXX 番地に別
9078 の情報を書き込む」という操作だと言われてもピンとこないので, コンピュータが人間にとってわかりやすい表
9079 現を見せて見せてくれているのです。このユーザにとってわかりやすい幻を見せてその気にさせてくれるという
9080 デザインのことを, ユーザーイリュージョンと言います。ともかくこういう「画面で見ながら操作する」ことをグ
9081 ラフィカル・ユーザ・インターフェイス (Graphical User Interface,GUI) といいます, これのおかげでコン
9082 ピュータの操作は随分楽になっています*4。

9083 C.3 コンピュータの歴史

9084 コンピュータの装置, すなわちハードウェアについての解説につづいて, ソフトの側面についても解説を加
9085 えようと思うのですが, そのためには少し歴史的な流れを説明したほうがわかりやすいかもしれません。

9086 コンピュータの発展の歴史は, 小型化の歴史でもあります。最初にできたコンピュータは ENIAC とい
9087 います。1946 年の話です。この ENIAC は 27 トン, 広さにして倉庫 1 つ分 ($167m^2$, 90 畳以上の広さ) が
9088 必要なもので, 真空管を使った計算機でした。軍事的な計画のために開発されたオーダーメイドのもので

*2 実はこのように人間を 1 つのコンピュータに喩えて, そこではどのように計算がされているのか, 人間の記憶装置や演算装置, 入出力装置の特徴はどうなっているのか, というのを研究するのも心理学の仕事です。記憶や演算, 制御については認知心理学や学習心理学, 入出力については知覚心理学や生理心理学が専門的に扱っています。そういう意味でもコンピュータの登場は心理学に大きな影響を与えているのですが, それはまた別の講義で。

*3 押すときの音から, この操作をクリック click と言います。2 回続けて押すことをダブルクリックと言います。

*4 実は人間の意識もこのユーザーイリュージョンのようなもので, 実際の体の動かし方や感覚情報の受け止め方, 処理の仕方は別ですよ。すべての情報を意識に上げるのではなく, 「私は XXX をしている」と身体がそれっぽい幻想を投影して見せてくれているのが意識の正体ではないか, という議論があります。興味のある人は [Norretranders \(2002\)](#) を読んでみてください。

9089 すから、一般人が触れるはずがないものです。時代が下がって真空管が半導体、IC チップになった頃、
9090 やっとサイズが小さくなって、家庭用・個人用のコンピュータというのできるかもという時代が来ました。
9091 Apple コンピュータの創始者、スティーブ・ジョブズとスティーブ・ウォズニアックが最初のパソコン (Personal
9092 Computer,PC), Apple I を発売したのが 1976 年。爆発的に売れた Apple II が発売されたのが 1977 年
9093 です。AppleII はブラウン管表示装置とキーボードを持っていましたので、今の PC の原型とも言えるかもしれ
9094 ません*⁵。PC を作っている会社は Apple だけでなく、IBM や DEC などがありましたが、まだこの頃は
9095 パーソナルなレベルのものよりも大型計算機の開発が進んでいました。IBM が PC を作ったのは 1980 年
9096 で、この頃から小型化が進められていきます。

9097 私事で恐縮ですが、1976 年に生まれた私が初めて PC を手にしたのは 1991 年、高校入学のお祝い
9098 買ってもらった Fujitsu の FM-Towns という機体でした。この頃は「マルチメディア」という名前もなく、「ハ
9099 イパーメディア」と読んで売り出していました。この機体は他の PC(NEC の PC-9801 や Sharp の X68000
9100 シリーズが有名でした)とは違って、CD-ROM ドライブをつけていたことが画期的だった時代です。その後
9101 1994 年に大学生になりましたが、この頃は連絡を取り合うツールはポケベルが主流であり、携帯電話 (や
9102 PHS) のような個人端末は高級品という時代でした。大学に入るとコンピュータを学ぶ授業があり、アカウン
9103 トをもらったりするのですが、それは大学が持っている大型計算機に端末からアクセスするためのものでし
9104 た。いろんな部屋にあるのは「端末」で、それほど機能の優れた PC ではなく、複雑な計算 (統計的な計算な
9105 ど) は大型計算機に仕事を依頼しその返信を待つ、というスタイルでした。関西私立のマンモス校でしたので
9106 学生数は非常に多かったのですが、多くの学生が一度にアクセスしても、大型計算機はものすごくものすご
9107 く計算が早かったので、瞬時に回答をもたらしてくれるものでした*⁶。つまり、まだ「専門的な計算は大型計算
9108 機」という時代であり、パーソナルなコンピュータになるにはもう少し時間が必要でした。

9109 どれぐらいの時間が必要だったかという、実はその次の年なのです。1995 年、Windows 95 という OS
9110 が発売されました。これを機に日本でも PC がどんどん浸透していくことになります。Windows95 は物凄
9111 んだぞ、と発売前からテレビでも散々とりあげられ、発売日には行列ができて真夜中のカウントダウンと同時
9112 に大フィーバー、という売れ行きでした。当時のそれは何が凄かったのでしょうか？ コンピュータにはそれ動か
9113 す基本ソフトが必要です。HDD にデータを書き込み、キーボードからの入力をディスプレイに表示する、と
9114 いったごく基本的な装置を統括し、メモリ番地をファイルという単位で扱うと言った基本的な操作は OS いう
9115 ソフトウェアが担当します (図 C.1)*⁷。この OS、大型計算機は Unix と呼ばれるものを使っていましたが、各
9116 企業が個人向けにコンピュータを売り始めるときにも当然必要で、各社で開発もしていましたが、PC の共通
9117 規格をつくることで OS 部分は共有できるようになりました。そこを提供したのが Microsoft 社のビル・ゲイツ
9118 です。どんなパーツで作られた PC であっても Windows という OS が共通のフィールドを用意してくれるの
9119 で、ユーザは Windows で動くアプリケーションを選ぶだけで良い、ということになったのです。

9120 そして Windows95 は、GUI、つまり「ファイルを掴んでポイ」といった直感的な操作で使えることも大きな
9121 特徴でした。大型計算機で使われている Linux は基本的に Command User Interface,CUI で、黒い画面
9122 にプログラムを書いて実行するといった手法で、初心者には人気がなかったのです。GUI については、その
9123 頃 Apple を追放されていたスティーブ・ジョブズが、NeXT という会社で GUI を備えた OS を開発してい
9124 ました。この NeXT はその後 Apple に買収され、ジョブズは Apple に戻って活躍することになります。その
9125 頃から巷では、Windows の GUI は Apple の OS を真似したものだと言われていたのですが、商業的に
9126 は Windows が大勝利、というわけです。ともかくこれを機に企業などはもちろん一般家庭でも PC を使うよ

*⁵ Mac の歴史については Walter (2012) が読み物としておもしろいですよ。

*⁶ Time Sharing System,TSS, 時分割システムとよばれる機構です。命令を小さな単位に分割し、それを順次捌いていくという方法でした。

*⁷ 物理的な機構と OS との間に Basic Input/Output System,BIOS というのが入りますが。

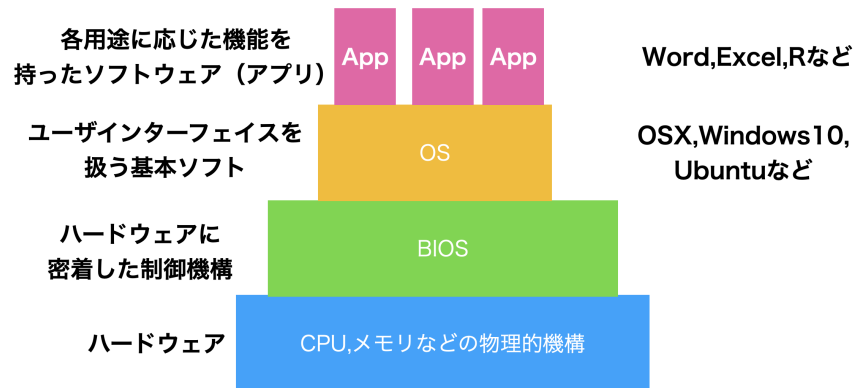


図 C.1 コンピュータのハードウェアとソフトウェアの関係

9127 うな時代になりました。

9128 ちなみにインターネットが広まったのもこの頃です。私は大学 3 年生の時、大学の授業で初めてインター
 9129 ネットを介して世界の情報を得る、という経験をしました^{*8}。その頃から徐々に、大型 PC のアカウントではな
 9130 くインターネットで使えるメールアドレスというのを個々人が持つようになりました。携帯電話も廉価な PHS
 9131 が広まり、メッセージだけでなく徐々に音声、写真、短い動画が遅れるようになっていきます。当初は当然画
 9132 質・音質も今とは比べものにならないのですが、それでもインターネットを経由してつながるという経験は想像
 9133 を超えたものでした。

9134 皆さんはすでに、携帯電話やインターネットがある時代に生まれた世代だと思います。こんな話はすでに過
 9135 去のものであり、不便な頃の話聞かされても困る、と思うかもしれません。私の世代は、幸いにもこのように
 9136 ちょうどコンピュータが使われはじめ、広がり、高性能になっていくにつれて育ってきていますので、そのぶん
 9137 機械の根本的な理解に直結しやすい社会環境にあったのです。みなさんは、便利な時代ではありますが、言
 9138 い換えると「初心者が苦労しないように補助輪をつけておく」「何もしなくてもできているような感覚が得られ
 9139 るようなイリュージョンを見せておく」という状態におかれているので、補助輪が対応できない道に進んだり、
 9140 多くの人とは違う使い方をすると、急にサポートが外れどこで困っているかわからなくなる、ということに
 9141 なりかねません。非常に大雑把な Historical Review ではありますが、理解の一助になればと思います。

9142 ところで、先ほどサポートという言葉を使いましたが、このサポートを商売にしているのが Microsoft や
 9143 Apple という IT 企業です。これらの会社は、ユーザが使いやすいようにソフトウェアを提供してくれますが、
 9144 有償ですし想定外の使用をするユーザに対してはサポートをしてくれません。逆にいうと、「決められた路線を
 9145 走らなければ面倒を見ない」ということでもあり、これはユーザに不自由を強いているとも言えます。また、ど
 9146 のような仕組みで動いているのかを尋ねても、企業秘密と言って答えてくれません。コンピュータは誰でも自
 9147 由に使えるべきものであり、勝手にユーザの情報を盗んだりしていないか、と言ったチェックをするためにも
 9148 オープンであるべきではないか。そういう考え方に基づく、フリーソフトウェアというソフトウェアのあり方があ
 9149 ります。Unix という OS を PC 用にした Linux がそうであり、オフィスソフトの LibreOffice、統計環境の R
 9150 などもこの精神に賛同するものです。フリーソフトウェアは自由であり、無償です。お金と秘密を払ってサポ
 9151 ートを受けるのではなく、ユーザが相互に助け合ってオープンで自由な世界を広げていこうという活動です。急
 9152 に何の話なんだ、と思うかもしれませんが、心理学ひいては科学的活動すべてにおいて重要な問題であるこ

^{*8} 教職関係の科目で、教育技術として今後使われるだろうということで担当教員が実演してくれました。隣の部屋から電話線を延長し、モデムというパソコンの信号を音情報に帰る機器を繋げて、NASA の Web ページを Netscape というブラウザでみたのが初めての体験でした。

9153 とをご理解いただきたいと思います。

9154 C.4 情報の単位

9155 コンピュータを取り巻く世界の話はこれまでに、ソフトの側面についての解説に入りました。

9156 コンピュータは文字、音、絵、動画ファイルいずれについても、すべて 0 か 1 のデータとして管理します。
 9157 0/1 の 1 つの単位を 1bit(ビット)と言います。1bit であれば Yes か No か、という二択の情報しか提
 9158 供できませんが、これが 7 つあれば $2^7 = 128$ ですから、これで 128 種類の状態を表現できます。コン
 9159 ピュータは一般に、8bit で 1 つの単位として計算します。8bit のことを 1byte(バイト)と言います。この
 9160 1byte が 1024 集まったものを 1kb(キロバイト)と言います^{*9}。1kb の次は、1024kb=1Mb(メガバイト)
 9161 です。さらに 1024MB=1GB(ギガバイト)で、1024GB=1TB(テラバイト)、1024TB=1PB(ペタバイト)、
 9162 1024PB=1EB(エクサバイト)と続きます。K,M,G,T,P,E といった名称は 1000 倍ごとに変わる大きさの桁
 9163 をあらわしているのであって、「ギガが減る」というのは本来意味をなさない表現です^{*10}。

9164 このビット・バイトは情報に関する基本単位なのであちこちに使われます。記憶装置について使われるとき、
 9165 一次記憶装置も二次記憶装置も同じ単位なので混同するかもしれませんが、2021 年現在では二次記憶装
 9166 置の単位は GB から TB が使われます。USB フラッシュメモリーや外付け HDD など数 TB の容量が一
 9167 般的でしょう。これに対して、一次記憶装置は 4~64GB ぐらいが相場かと思います。一次記憶装置は暗算の
 9168 途中経過のように一時的に記憶する場所に過ぎないので十数 G でも問題ありませんが、二次記憶装置は結
 9169 果の記録なので大きければ大きいほど余裕が持てますね。たとえば数 TB でもテレビドラマや映画を何本も
 9170 記憶できるのですから、PB や EB なんて使うのかな、と侮ってはいけません^{*11}。すぐにそれぐらいのサイズ
 9171 が必要な時代が来ることでしょう。

9172 実際、それぞれ単位ではどれぐらいの情報が記録できるのでしょうか。1byte は 128 文字表現できる
 9173 ので、英語のアルファベット 26 文字に加え、数字や簡単な記号であれば 1byte で表現できます。たとえ
 9174 ば A という文字は 01000001、B という文字は 01000010、小文字の a は 01100001、と言ったように 0/1
 9175 の文字列 8 個と一対一対応させて考えるのです。日本語は 1byte では足りませんので、一文字あたり
 9176 2byte が割り当てられます。また 1KB(=1024byte) は 500 文字ですから、原稿用紙一枚ぐらいになります。
 9177 1MB(=1024KB) は文字だけだと新聞一紙(朝刊の 40 ページ分)ぐらいで、昔の記録媒体であるフロッ
 9178 ピーディスク一枚に保存できるのがちょうど 1MB でした^{*12}。ちなみに「カメラ映像 + 音声」のオンラインビデ
 9179 オ会議を 1 時間やると 200~300MB ぐらいの容量をやとりしていることになります。ビデオ会議では文字
 9180 (チャット)だけでなく、画像(動画)や音声も送りますね。実は画像や音声も、0/1 に置き換えています。
 9181 音声の場合、1 秒間を短い間隔にくぎります。この区切りのことをサンプリング周波数といい、たとえば 44.1
 9182 kHz という単位は 1 秒間に $44.1 \times 1000 = 44100$ 点のデータの採取をします。この 6 データ点において、音
 9183 の振幅を区切ります。ビットレートと言いますが、たとえば 16bit で区切る場合は $2^{16} = 65,536$ 段階で区
 9184 切り、その高さの音がある (1) かない (0) かで表すわけです。このように時間と音階を細かく区切り、その目

^{*9} キロメートルやキログラムのように、キロは 1000 の単位を表す言葉ですが、コンピュータは 2 進数なので 1000 ちょうどではなく 1024 で 1 つ上の位に上がることになります。

^{*10} 同様に「USB を紛失する」というのもよく聞くおかしな表現です。USB は Universal Serial Bus の略で、データ転送規格のことを指します。なくすことができるのはそれにつなげるメモリースティックなどです。

^{*11} 私事ですが、私が 1991 年に生まれて初めて手した PC、FM-Towns はハードディスクが 40MB、RAM は 2MB で、CD-ROM がついていましたが、それは 360MB の容量でした。購入するときに、「ハードディスクが 40MB もあって何を記録するんです? CD-ROM なんか情報が詰まり過ぎてますよ!」と店員さんに笑われたのを今でも覚えています。数年後、PC の動きが遅くなったので、2 万円で 2M の RAM を追加したのも良い思い出です。今でこそ、2MB なんて USB フラッシュメモリーでも売ってないほど微小なサイズですが。

^{*12} 正確には 1.2MB 入る規格 (2DD) と 1.44MB 入る規格 (2HD) とがあります。

9185 に情報があるかないかを積み重ねてデータとするわけです。テキストよりも圧倒的に情報が多くなるのが分か
9186 りますね。画像も同様に、図面を細かい単位 (ピクセルなど) で分けて、その色合いを色々な段階で区切り
9187 ます。色は R(赤)G(緑)B(青) の組み合わせで表現でき、それぞれを $8\text{bit} = 2^8 = 256$ 段階で表現したりし
9188 ます。一点一点にその情報がありますから、図の情報も非常に多くなるのがわかると思います。動画はその画
9189 像が時系列的に細かく分割されたものと思ってください。このように分解しますので、テキスト、音声、図、動
9190 画の順にデータサイズが大きくなります。通信機器が最初ポケベル=数 byte の情報しか送れなかったもの
9191 から、徐々に絵文字、ショートメッセージ (音声)、写真^{*13}、動画が送れるように発展していきました。今では町
9192 中のあちこちで、誰もが手軽に動画をモバイル端末で見られるようになっています。あらためて、すごい進歩
9193 ですね^{*14*15}。

9194 ところで、1byte は 256 種の情報が記録できるので、英語のアルファベットや数字は 1byte あれば十
9195 分だが、日本語や中国語など、英語以外の言語は文字種が多いので、2byte で一文字を表すという話
9196 をしました。この 2 バイト文字も、たとえば「あ」とか「亜」という文字に 111000111000000110000010 とか
9197 111001001011101010011100 という文字列を割り当ててののですが、言語ごとによってどの数字をどの文字
9198 に割り当てるかという対応表が変わってきます。これを文字コードと言います。日本語はかつて Shift-JIS と
9199 というコードで変換していましたが、今は世界のあらゆる言語に対応している共通企画である、UTF-8 という
9200 文字コードで変換することが一般的です。ところがなぜか、日本の Windows OS だけ Shift-JIS をいまだに
9201 使い続けており、他の PC とファイルをやり取りするときに文字コードの変換エラー問題が起きます。ファイル
9202 を開いて文字化けをしたりとか、プログラムが実行される際に「ファイルにアクセスできない」というエラーが生
9203 じたりするのは、この文字コードの問題が大きいのです^{*16}。受身的な対策法になってしましますが、PC で
9204 つかうユーザ名やファイル名などは半角英数文字を使い、短くした方がこうしたエラーに出くわしにくくなりま
9205 す。逆に、全角文字やスペース (空白) などを含んだファイル名、やたらと長いファイル名を使っていると、こう
9206 した問題に出くわしやすくなるということです。

9207 C.5 ファイルの種類と拡張子

9208 ここまで述べてきたように、計算機というのは基本的に物理的実体 (記録装置、記憶装置) の上で 0/1 の
9209 データをやり取りしているだけです。記録 (記憶) 装置上に置かれている情報のセットは「ファイル」という形
9210 で記録されています。スマートフォンやタブレットは、ユーザの利便性のためにファイルの存在を意識しなくて
9211 も良いようになってはいますが、バックエンドでは実行されるアプリケーションもファイルですし、開かれる音
9212 声や動画もファイルです。とくにパソコンでは、どの媒体、どのアプリで使うどういうファイルかを識別するた
9213 めに、拡張子 (かちようし) と呼ばれる識別記号をファイルの後ろにつけています。拡張子はファイル名の背後
9214 にピリオドで区切って追記されています。ついてないように見えても、OS がそれを表示させない設定にして
9215 あるだけであることに注意してください。代表的な拡張子と、それに対応づけられているアプリケーションは次

*13 できた当時は写真を撮ってメールができることをとくに「写メールする」と言い表したほどです。

*14 実は音でも画像でも、分割してそのままデータにしてしまうと膨大になりすぎるので、人間が気付きにくい周波数や色合いなどは削除して作ります。これを非可逆圧縮処理と言います。一度落としてしまった情報は戻らない、という意味です。ライブや生きている人間が処理している情報は、携帯の画面から得られるものの何億倍もの情報量なんですよ!

*15 デジタル化のすごいところは、こうした文字、音、図版、動画といったものを bit という共通の単位に落とし込んだことです。こうすることですべて一元的 (bit という共通次元) で処理することができるようになったのです。メディアの違いが問題にならなくなり、0/1 の情報であれば複写も簡単にできてしまいます。情報化社会においては情報に特別性はなく、情報があるかないか、それを生み出せるかどうかこそが重要なのです。

*16 Windows だけ世界標準から外れているので、早く修正して欲しいのですが、歴史的な経緯からユーザ数が多くて切り替えられないでいることと、こうした違いがあることをユーザに説明しない (素人は知らなくていい、と馬鹿にされているようなもの) ので、問題が解決される日はまだ先になりそうです。

9216 のようなものがあります。

9217 .docx マイクロソフト社の文書作成アプリケーション、Word で使うファイル

9218 .xlsx マイクロソフト社の表計算アプリケーション、Excel で使うファイル

9219 .pdf Adobe 社が開発した Portable Document Format 形式。OS が違っても同じレイアウトで文書を表示できるのが利点で、PDF 形式を読むことができるアプリケーションは多数。

9221 .txt シンプルな文字だけのテキストファイル。文字の飾りやレイアウトなどの情報がない最もプレーンな形式なので、OS が違っても文字コードさえ合っていれば読むことができる。

9223 .mp3 音楽、音声のファイル。音声データには他の種類もあります。

9224 .jpg 画像のファイル形式の一種。

9225 .png 画像のファイル形式の一種。

9226 .csv comma/character separated variables ファイル。変数をカンマ (,)、あるいはタブ、半角スペースなど文字コードで区切ったファイルという意味で、中身は.txt と同じく装飾のない文字/数字だけであり、文字コードさえ間違えなければ OS を問わずに読み書きできる。データのやり取りはこの形式で行われることが多い。

9230 .zip 圧縮ファイルの一種。1 つまたは複数のファイルをパックして圧縮してあるもの。ファイルの冗長な部分をうまくまとめてコンパクトにまとめ上げるため、ファイルサイズが小さくなるし、複数のファイルもひとまとめにできる。また圧縮の際にパスワードをかけることもできるため、メールなどに添付する場合はこの形式にまとめられることが多い^{*17}。可逆圧縮であり、圧縮されたファイルは展開する (解凍するともいう) ことでパッキングを開封できる。zip ファイルの圧縮/展開は各種 OS が標準的に対応している。

9235 .txt や .csv といった形式は、「装飾のない、文字だけの」ファイルです。こうした種類のことを ASCII ファイルと言います。メモ帳などのエディタと呼ばれるプログラムで読み書きできます。逆に .docx など特定の会社が提供するアプリケーションに対応しているファイルは、アプリの中でのさまざまな操作・装飾を暗号化して保存しており、メモ帳などで読んでも意味がわかりません。対応しないアプリでは開くこともできません。こうした形式は ASCII ファイルに対してバイナリファイルと言います。

9240 OS は拡張子を見てファイルの種類を判別し、そのファイルを開くのに適したアプリケーションを自動的に起動し、開いてくれます^{*18}。 .csv ファイルは Excel などのアプリケーションで開くことも当然できますが、その際文字コードのエラーが生じたり、保存するときに文字コードを変えたりして、形式・内容が気づかずに変わっていることがあります。Windows も良かれと思ってやっていることなのですが、処理が徹底してないのか、かえって不便になってしまっています^{*19}。

^{*17} PPAP というピコ太郎の楽曲を思い出す人もいるかもしれませんが、圧縮ファイルの文脈では「Password つき zip ファイルを送ります。Password は次のメールで送ります」Angoka(暗号化) Protocol の略です。つまりメールでパスワード付きのファイルを送り、そのファイルを開くためのパスワードをまたメールで送るという、日本でよく見られるおかしな風習です。おかしな、というのは、メールがハッキングされていたらパスワードもどうせバレるわけで、同じメールに書いてあるのはもちろん馬鹿馬鹿しいですが、すぐ次のメールに書いてあるのも同じぐらいに馬鹿馬鹿しいことです。情報セキュリティ対策手法のつもりで行われる慣習が広まっていますが、PPAP の標語のもと、馬鹿馬鹿しいのでやめましょうという風潮になってきました。

^{*18} 見たことのない拡張子の場合は、どのアプリケーションで開くべきか尋ねてくるでしょう。

^{*19} 実際、この授業での課題データを UTF-8 形式の csv ファイルで提供しても、Excel で開いたばかりに文字化けして分析できなくなる、という相談がこれまで多く寄せられています。根本的な解決策として、Windows を使うのをやめることをお勧めします。

9245 C.6 クラウドとは

9246 すでに述べたように、計算機は基本的に 0/1 データのやりとりであり、それを保存してあるのがファイルと
9247 よばれるものです。ファイルは HDD や USB メモリ、SSD に保存することができます。

9248 ところで、最近はこうした手元の物理的実体にファイルを置くことに加えて、クラウドに保存することも少な
9249 くありません。クラウドとは雲という意味で、インターネットの向こう側のどこか、ということの意味します。です
9250 が、基本的にはインターネットで繋いだその先にも電子計算機があるのです。たとえばパソコン A とパソコン
9251 B をケーブルで繋ぐと、パソコン A からパソコン B のファイルにアクセスできます*20。このケーブルをどんど
9252 ん伸ばすと、遠く離れていてもこの操作ができます。このケーブル網を世界レベルに広げているのがインター
9253 ネットです*21。

9254 ここで覚えておいて欲しいのは、当たり前のように、ネットといっても基本的には電子計算機と電子計
9255 算機を繋げている実体がどこかにあって、ファイルのやりとりをしているだけだということです。ブラウザはウェ
9256 ブサイトの情報を書いたファイルを取り込んでホームページを見せてくれていますし*22、Youtube は動画の
ファイル、Instagram は画像のファイルへのアクセスをして見せてくれているのです。最近ではクラウドサー

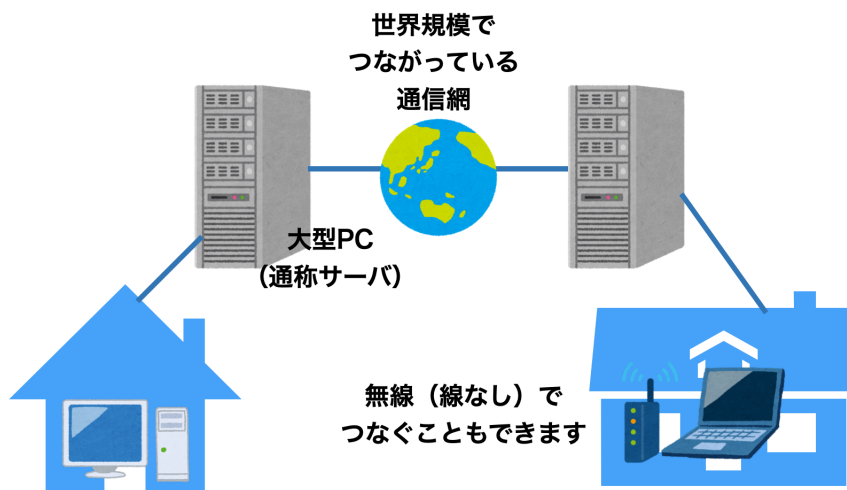


図 C.2 インターネットの世界

9257 ビス、というのがよくあります。中でも Dropbox とか OneDrive, iCloud などが有名です。これらは、ファイ
9258 ルをインターネットを介した別の大型電子計算機に保存 (アップロード) し、インターネットを介して読み込み
9259 (ダウンロード) して使う、というものです。手元の電子計算機の中に記録されているファイルのことを「ローカ
9260 ル」、サーバなど遠隔地にある大型機に記録されているファイルのことを「サーバ」「クラウド」と呼んで区別しま
9261 すが、このローカルとサーバのファイルを常に同じものに同期しておく便利なシステムです。こうしたクラウド
9262 サービスを使うと、たとえば大学で作業して保存したファイルを、USB メモリに保存して自宅のパソコンにコ
9263 ンピュータ

*20 もちろんファイルの情報をどのように信号に変えて送受信するか、アクセス権限はどうかといった細々したことを調整しなけれ
ばなりません、そのあたりの仕事をしてくれるのが OS のありがたいところです。

*21 インターネットは軍事通信網として始まりましたが、それを電話回線や企業内通信網、大学間通信網などと接続しあって世界中
に広がっています。網と網が相互に (inter) 繋がっているため inter net であり、大学内・企業内のネットワークのことはイントラ
(intra) ネットと言います。ちなみに一般用語としてのインターネットは Internet と大文字で書き始めます。

*22 ホームページとは本来ブラウザが最初に開くページのことを指し、企業や個人が情報発信しているページのことはウェブサイトと
いうのが適切です。

9264 ピーする, という操作をしなくても, 大学でクラウドサービス上のフォルダ内に保存しただけで, 自宅のパソコンにそのファイルがコピーされているため^{*23}, 意識せずにそのまま作業を続けられるということです。自動的にバックアップをとってくれているとも言えるので便利です。

9267 もちろん注意すべきこともあります。「他の人に見られたら困るファイル」, 「アプリケーションの実行ファイル」などは通信網の向こうではなく, 手元 (ローカル) に置いておきましょう。個人情報・機密情報などを, クラウドドライブに保存すると, 悪意を持った人が大型計算機に攻撃を仕掛けて情報を盗んでいく可能性があるからです。情報化の怖いところは, 取られても気づかない (コピーすればいいだけで元ファイルになんの影響もない) ところにあります。加えて, 取られたものをばら撒かれる = インターネットを介して誰でもアクセスし保存できるようにされると, すべてを回収できなくなるのも問題です。失言が記録されて拡散されると大変なことになるのは, みなさんもこの時代に生きる人間のマナーとして色々見聞するところだと思います。また, クラウドサービスで自動的に同期されるといっても, アプリケーションの実行ファイルなどはローカルに保存すべきです。アプリケーションは実行に際してさまざまな関連ファイルにアクセスしますので, 1 つでも場所が違っているとエラーになって動かなくなります。同じ OS でも, です。インターネットからとってきたソフトウェアをうっかり OneDrive に保存してしまうと, アクセスできないエラーで起動しない, ということもありますので注意してください^{*24}。

9279 C.7 ファイルの位置の指定

9280 ここでファイルとそのパスについての話をしておきたいと思います。

9281 計算機が情報を 0/1 で管理し, それらがファイルとなってどこかに保存されている, ということでした。私たちは Finder や Exploer などファイルブラウザをつかって, 実行したいファイル, 参照したいファイルを探していきますね。ファイルはまとめてフォルダの中に含まれていますし, フォルダの中にフォルダがあるといった, 階層状態になっていることも少なくありません。ちなみに **フォルダ**と同じ意味で**ディレクトリ**という言葉が使われることもあります。

9286 C.7.1 相対パスと絶対パス

9287 普段 PC を使っているときは気にすることがありませんが, R や RStudio などプログラミング言語をつかっているときは, 「今どこで作業しているか」という現在地が重要になってきます。たとえば, RStudio で `C:\User\kosugitti\Document\kiso1\` というところでプロジェクトを開いているとします。スラッシュ (\) はフォルダ, コロン (:) はドライブを表す記号です。プロジェクトフォルダは, C ドライブの User フォルダの下にある, kosugitti フォルダの下にある, Document フォルダの下にある, kiso1 というフォルダということになります (プロジェクト名が kiso1 だとそうなります。)。この kiso1 フォルダが現在地です。

9293 このフォルダの中で, Rmd ファイルや R スクリプトファイルを使って, 他のファイルを参照するようなコードを書くとしましょう。たとえば `script1.R` というファイルに `read_csv` 関数を書いたとします。読み込みたいファイルは, 同じフォルダの中にある, `sample.csv` とだとします。このとき, `read_csv` の書き方は次のようになるでしょう。

^{*23} これはユーザが特段の指示をしなくても, アプリケーションがユーザの見えないところで (バックグラウンドで) アップロード, ダウンロードの作業を進めているからです。パソコンはシャットダウンしていなければ, 裏でこうした作業を淡々とこなし続けてくれます。

^{*24} マイクロソフト社は, これまたユーザのためを思ってやっているのかもしれませんが, デフォルトで OneDrive に保存させようとしています。それでうまくインストールできなかったという相談も時々寄せられています。抜本的な解決策として, Windows を使わないことをお勧めします。

code : C.1 相対パスで読み込む

```

9297 1 dat <- read_csv("sample.csv")
9298
9299

```

9300 しかし別の書き方もあります。たとえば code:C.2 のような書き方でも問題ありません。

code : C.2 絶対パスで読み込む

```

9301 1 dat <- read_csv("C:\User\kosugitti\Document\kiso1\sample.csv")
9302
9303

```

9304 後者 code:C.2 の書き方は、ファイルの場所を全部書いてありますから、確実にその場所が特定できます。
 9305 それに比べて前者 code:C.1 の書き方は、なぜファイルを書いただけでいいのでしょうか。これは、このコード
 9306 を実行している現在地と同じフォルダの中に sample.csv ファイルがあるからです。プログラムは、命令を受
 9307 けるとファイルを探しにいきますが、現在地と同じフォルダの中を探すことになっているのです。この現在地、
 9308 すなわち現在作業しているフォルダのことを**作業フォルダ (working directory)**と言います。

9309 では作業フォルダと別のフォルダの中にファイルがあれば、アクセスできないのでしょうか。そんなことはあ
 9310 りません。code:C.2 の書き方を使えば、作業フォルダがどこにあっても位置を特定できますから、作業フォ
 9311 ルダを問わずに書くことができます。ちょっと長くて面倒ですが、確実にある場所を指定しているからです。
 9312 この書き方のことを**絶対パス**による指定と言います。一方、code:C.1 の書き方は、今の作業フォルダから
 9313 見た場所、という相対的な書き方になっています。この書き方のことを**相対パス**による指定、と言います。相
 9314 対パス指定で、違うフォルダにアクセスする場合には、次のようにします (code:C.3)。ここでは、Document
 9315 フォルダの中に、kiso1,kiso2 フォルダがあり、kiso1 フォルダの中で作業している時に kiso2 フォルダの
 9316 sample2.csv ファイルを読み込む例を書いています。

code : C.3 絶対パスで読み込む

```

9317 1 # 絶対パス指定
9318 2 dat <- read_csv("C:\User\kosugitti\Document\kiso2\sample2.csv")
9319 3 # 相対パス指定
9320 4 dat <- read_csv("../kiso2/sample2.csv")
9321
9322

```

9323 絶対パスはそのままなのですが、相対パスは..****という記号になっていますね。このピリオドを 2 つ打つ方
 9324 法で、「ひとつ上のレベルの」という意味になります。このように、現在地からの相対的な位置関係で、ファイル
 9325 を指定することもできます。

9326 絶対パスと相対パスのどちらが良いのか、というのは一概には決められません。絶対パスは、PC が
 9327 変わったりフォルダの構造が変わったりすると役に立ちませんから、使い勝手が悪いと言えなくもない
 9328 ですが、確実に指定できる方法です。相対パスは、PC が変わったりしても「現在地から相対的に見て
 9329 どこか」という話ですから、たとえばこの例で kosugitti フォルダごと別の場所に移しても (たとえば
 9330 D:\Univ\Classes\kosugitti\kiso1 のように)、コードはエラーなく動きます。kiso1 フォルダ、kiso2
 9331 フォルダの相対的な位置関係が変わらなければいいのですから。バックアップを取ったり、複数の環境で同
 9332 期しながら作業する場合などは相対パスの方がいいでしょうね。

9333 いずれにせよ、現在どこで作業しているかということ、すなわち作業フォルダの場所を、意外と意識しておか
 9334 なければならないということには注意が必要です。ファイルをどこに置いたか、どんなファイルを置いたか、自
 9335 分はどこにいるのか、これが変わってくると「ファイルが見つかりません」というエラーになるのです。言い方
 9336 を変えると、**ファイルが見つかりませんエラーの原因は、この 3 つのどれかであることがほとんどです。**

C.8 ファイルのバージョン管理

これからみなさんは大学生活の中で、たくさんのファイルを生み出していくことになるでしょう。たとえば 4 年生の時に卒業論文を書くことになりませんが、データファイル、分析ファイル、図を書いたファイル、引用文献リストを書いたファイル、卒論本文などなど、1 つの研究でも複数のファイルが作られることはよくあります。さらに、これらのファイルは日々加工されますから、その度に上書き保存することになります。いわば、ファイルがバージョンアップしていくのです。

卒論などの場合はとくに、「途中で保存しておく」ことが重要です。途中まで書いていた時に、横に置いていたマグカップが倒れて PC にコーヒーがかかり、変な音を立てて PC が壊れてしまった、ということがあるかもしれませんからね。紙に書いていた時代は、その手のハプニングがあってもせいぜい原稿用紙数枚がダメになっただけで、「ちくしょう、やりなおしかぁ」で済んだのですが、電子データの場合は電子の藻屑になると復元させることができません。ですから**バックアップは非常に大事**なのです*25。

バックアップの基本は、「別の場所に」「別のファイル名で」というものです。同じ名前の上書きすると元に戻ることができませんから、面倒でもコツコツと違う名前をつけましょう。そうするとよくあるのが、`soturon1.docx`, `soturon2.docx`, `soturon3.docx`, `soturon3(修正).docx`, `soturon3(最新).docx`, `soturon3(最新)(修正).docx`, `soturon3(最新)(修正)(提出版).docx`, `soturon3(最新)(修正)(提出版2).docx`... というようになっていくやつで、「どれが最新版だっけ...」と書いてる本人でも探すのに苦労することになります。

この問題の解決策として、`soturon1001.docx`, `soturon1005.docx` のように日付を入れるというものがあります。10 月 1 日分、10 月 5 日分、としていけば「いつまで戻れば良いか」もわかるのでいいやり方ですね。日付の数字をファイルの先頭につけておくと、並べ替えも簡単です。この日付をつけて保存するというのを習慣化し、1. 昨日までのファイルを開いて今日の日付で別名保存する、2. 作業を進めて、時々上書き保存、最後にも上書き保存、3. PC に USB メモリをさして、バックアップ保存して作業終了、というルーティンを作っておくと、確実に記録が残って良いでしょう。

ただし、この方法の問題は、ファイルサイズが大きくなりすぎることです。図表などを含めたファイルが数百 MB になることは少なくありません。それを次々複製していくわけですから、大容量の USB メモリでも限界が来るかもしれません。これは「丸ごとコピー」していることが原因で、たとえば昨日は 10 行目まで書いた、今日は 11 行目から 14 行目まで書いた、というのであれば、この増えた 4 行分 (差分) だけを追加保存すればいいのに... と思いませんか。

こうしたバックアップやバージョン管理をやってくれる仕組みとして、Git というものがあります。Git は作業フォルダの中身の差分だけを記録し、必要であれば過去のバージョンに戻すこともできるシステムです。毎回上書き保存 (commit する、といいます) のたびに「どこを変更したか」というメモを付けて保存しておけば、そのメモを見ながら「ここの時点まで戻ろう」という使い方をすることができます。ファイル名は変更する必要なく、同じファイル名で進めていけますから、同じようなファイルがたくさんあって訳がわからなくなるということもありません。さらに保存先をクラウドにした GitHub というものもあり、これを使うとクラウド上に追記していくことができます。この GitHub は IT 企業などでプログラムをチームで進めていく時にも使われている技術で、全体のプログラムに個別の機能を複数人が追加、管理者が必要なものだけ取り入れる、というように使われています。国里ゼミや小杉ゼミでは、卒論を GitHub で管理し、学生が書いた分を commit し、それを

*25 ちなみに私の経験上、レポートが電子の藻屑になったので助けてください！と言われることがよくあり、4 年間の大学生活の間では平均 10-15 名に 1 人の割合で発生することのようです。

9374 hub 上にアップロード (push といいます) する, というようにします。教員の方は学生の進捗が管理できま
9375 すし, どこがどう変わったかが分かりやすく, バージョン管理と同時にバックアップもできるという便利な仕組み
9376 です。GitHub は無料でアカウントを作ることができますから, 興味があれば皆さんもチャレンジしてみたく
9377 さい*26。

9378 さてここで, ひとつ注意をしておきます。卒論の原稿やプログラムは日々変化するものですからバージョン
9379 管理が必要ですが, データファイルはアップデートする必要がありません。いや, アップデートしてはおかし
9380 のです。たとえば 100 人分のデータを取って分析をしていて, 後で「やっぱりこのデータを削ろう」というのは,
9381 研究不正が疑われかねません。自分に都合の良いデータだけで議論し, 都合の悪いデータは削除して統計
9382 的に有意な結果が出るように細工しよう, なんてことがあれば困るのです。データファイルはバージョンアップ
9383 せず, それを加工, 計算するプログラムがバージョンアップしていく。そしてその加工プロセスは誰でも見るこ
9384 とができるように公開されている。少なくとも, 科学的な営みをする上では, そうしたやり方が必要なのです。
9385 自分だけのデータで自分だけの分析方法で, 良い結果だけ示すというのは適切な方法ではありません。

9386 Open Science Framework(<https://osf.io/>) はこうした「オープンな科学」にむけた取り組みの一種で
9387 す。このアカウントは誰でも無料で作ることができ, ここにファイルをアップロードしたり, 分析計画を事前に記
9388 録しておくことができます。何も難しいことではなくて, クラウド上のファイル置き場だ, というぐらいに思ってい
9389 ただければ結構です。ここにおかれたファイルは自動的にバージョン管理され, 同じファイル名のものがアップ
9390 デートされてもその記録 (ログ) が残ります。最近はこの OSF をつかって, 論文文化されたデータやプログラム
9391 を公開するという取り組みも進んでいます。

9392 クラウド, バックアップ, オープンサイエンスといった新しい研究方法は日々生まれてきています。皆さんも
9393 便利な機能はどんどんキャッチアップしていきましょう!

9394 C.9 おわりに

9395 古臭い話をしてしまうのは, 私が歳をとった証拠でしょう。皆さんはこんな話を知らなくても, スマホやタ
9396 ブレットを使いこなしていることと思います。細かいことを知らなくても, ユーザとして利用するだけなら知らな
9397 くて良い話なのかもしれません。私はここにも書いたように, 高校生の頃からコンピュータの発展と一緒に大
9398 人になってきましたから, 学ぶともなく学んできたところがあります。皆さんは生まれた頃からコンピュータや
9399 があったネイティブ・デジタル・カウボーイですから, 苦労なんかする必要なかったわけです。

9400 しかし細かい仕組みを知らないということは, 問題が生じた時に「何か・誰かが, どこかでどうにかなって,
9401 今私が困っている」という状況になります。問題を特定できないと, 解決することもできません。コンピュータ
9402 は文房具に過ぎませんから, それを使いこなせないほうが格好悪いのです。しかも今後ますますコンピュータ
9403 に囲まれた世界になっていくのは自明ですから, ここに学習コストをかけない方が勿体無い。幸い, わからな
9404 いことに対して, 自ら調べて学んだ利する時間と環境が用意されているのが大学という世界なのですから, 今
9405 のうちにしっかり基礎固めをしておきましょう。

9406 このくだらない懐古的エッセイのような文章が, 何かの足しになれば幸いです。

*26 このテキストやシラバスも GitHub で管理していますし, 公開されているサイトも GitHub 上のもので。これからは教科書も
日々成長していくものになるかもしれません。

9407 付録 D

9408 ギリシア文字一覧

ギリシア文字ってかっこいいけど、読み方わからない・・・という人のための一覧。ついでに $\text{T}_{\text{E}}\text{X}$ 表記も。

表 D.1 ギリシア文字とその読み方

読み方	大文字	小文字	英語表記	$\text{T}_{\text{E}}\text{X}$ 表記
アルファ	A	α	alpha	<code>\alpha</code>
ベータ	B	β	beta	<code>\beta</code>
ガンマ	Γ	γ	gamma	<code>\gamma</code>
デルタ	Δ	δ	delta	<code>\delta</code>
エプシロン	E	ε	epsilon	<code>\varepsilon</code>
ゼータ	Z	ζ	zeta	<code>\zeta</code>
イータ	H	η	eta	<code>\eta</code>
シータ	Θ	θ	theta	<code>\theta</code>
イオタ	I	ι	iota	<code>\iota</code>
カッパ	K	κ	kappa	<code>\kappa</code>
ラムダ	Λ	λ	lambda	<code>\lambda</code>
ミュー	M	μ	mu	<code>\mu</code>
ニュー	N	ν	nu	<code>\nu</code>
クサイ	Ξ	ξ	xi	<code>\xi</code>
オミクロン	O	\omicron	omicron	<code>\mathrm{o}</code>
パイ	Π	π	pi	<code>\pi</code>
ロー	R	ρ	rho	<code>\rho</code>
シグマ	Σ	σ	sigma	<code>\sigma</code>
タウ	T	τ	tau	<code>\tau</code>
ウプシロン	U	υ	upsilon	<code>\upsilon</code>
ファイ	Φ	ϕ	phi	<code>\phi</code>
カイ	X	χ	chi	<code>\chi</code>
プサイ	Ψ	ψ	psi	<code>\psi</code>
オメガ	Ω	ω	omega	<code>\omega</code>

9410 付録 E

9411 記号の入力とキーボードの場所

9412 プログラミングのミスでよくあるのが打ち間違い、スペルミスです。X と x, S と s など大文字と小文字で
 9413 形が同じものや, l(エルの小文字) と 1(数字のイチ) の違いなどは, プログラミング用フォントにして違いがわ
 9414 かるようにするとか, 文字列の意味から類推する (Normal とあればノーマルであって, ノーマ・イチではない
 9415 と察する) など工夫が必要かもしれません。

9416 また, 理由はよくわからないのですが頻発するスペルミスは, データ (data) をデート (date) と書いてしま
 9417 うことです。データはラテン語の datum(与えられたもの) の複数形なのですが, 最近のデータサイエンスの
 9418 文脈では data も単数形と考えるようです。ともかく, 日付を表す date とは由来も意味もスペルも全て違う
 9419 ので, 気をつけましょう。

9420 さて, これらはまだ序の口。プログラミングのコードを読んでも, 日本語の五十音に入っていない記号の違
 9421 いがわからない, どこでそれが入力できるかわからない, 質問しようにも読み方がわからないといったものも
 9422 あります。ここではこれらをまとめて解説します。記号の上では微妙な違いですが, 当然形が違うのでプログ
 9423 ラミング上は違う文字として扱われますので, 形の細部までよくみてください。なお, プログラミングでつ変わ
 9424 る時は言語に依存することもありますので, ご注意ください*1。

9425 特に目立つのはハイフンとチルダ, アンダースコアの入力ミスです。それぞれ文字を書く領域における位置
 9426 が違ったり, 形が違ったりするのでよくみてください。

ハイフンは真ん中	A-B
アンダースコアは下	A_B
オーバーバーは上	A [~] B
チルダはニョロ	A~B

9428 これを踏まえて, そのほかの記号の名称や意味を表 E.1 で確認しましょう。

9429 一般的な日本語キー配列の場合, 1つのキーに4つの文字・記号が割り当てられていますが, 英語入力
 9430 モードの場合はキーの左側, 日本語入力モードはキーの右側を見ることになります。キーを押すと下の段の
 9431 文字が入力されますが, シフトキーを押しながらキーを押すことで上の段の文字が入力されることになります
 9432 (図 E.1)。

9433 これを踏まえて, 日本語キーについては E.2, US キーについては図 E.3 に代表的な記号とキー配列の位
 9434 置を示しました。入力に困った場合は一度図を見て確認してください*2。

*1 たとえばコメントアウトは C 言語では \\, R では #, TeX では % など, それぞれ異なります。

*2 US キー配列はキートップに一種類の文字しかなく, 美しい配置なのでおすすめです。

表 E.1 記号と読み方

記号	読み方	解説
:	コロン	英文中では「すなわち」などの意味。セミコロンと間違えないように
;	セミコロン	英文中では文章の区切り, 接続詞のようにつかう
.	ピリオド	英文の終わりを意味する。日本語で言う句点
,	カンマ	英文の区切りを意味する。日本語で言う読点
@	アットマーク	メールアドレスに用いられることで有名
\$	ドルマーク	米国の通貨単位。R では変数名指定のときにもちいる
/	スラッシュ	割り算の記号
*	アスタリスク	掛け算の記号
+	プラス	足し算の記号
-	マイナス	引き算の記号
^	ハット	累乗の計算の記号
=	イコール	等号。プログラミングでは==で一致しているかどうかの判定にも
!	エクスクラメーション	強調。プログラミングでは否定 (NOT) の意味になることも
_	アンダースコア, アンダーバー	位置に注意。ハイフンではなく文字領域の下の線 変数名をつなげる時 (ex. snake_case) に使ったりする R では独立変数と従属変数をつなぐときに使う
~	チルダ	ハイフンやオーバーライン, アンダースコアと間違えられる率高め
-	オーバーライン, オーバーバー	アンダースコアの逆。滅多に使わないが。文字化けを直したときにみられる。
%	パーセント	プログラミングではコメントアウトの時などに使われたりする
&	アンパサンド	プログラミングにおける AND(論理積) の記号など
	縦棒	プログラミングにおける OR(論理和), 条件付き確率の記号にも
\	バックスラッシュ	プログラミングではコメントアウトの時などに使われたりする
"	ダブルクォーテーション	文字列の開始・終了を表す。同じ記号で閉じる
'	シングルクォーテーション	文字列の開始・終了を表す。同じ記号で閉じる
`	バッククォーテーション	文字列の開始・終了を表す。同じ記号で閉じる
[]	大括弧	プログラミングでは配列を意味することがある
{}	中括弧	プログラミングではブロックの開始と終了を意味することがある
()	小括弧	数式のまとまりを表す

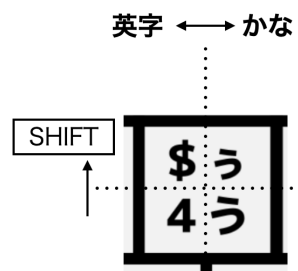


図 E.1 日本語キーで入力する場合

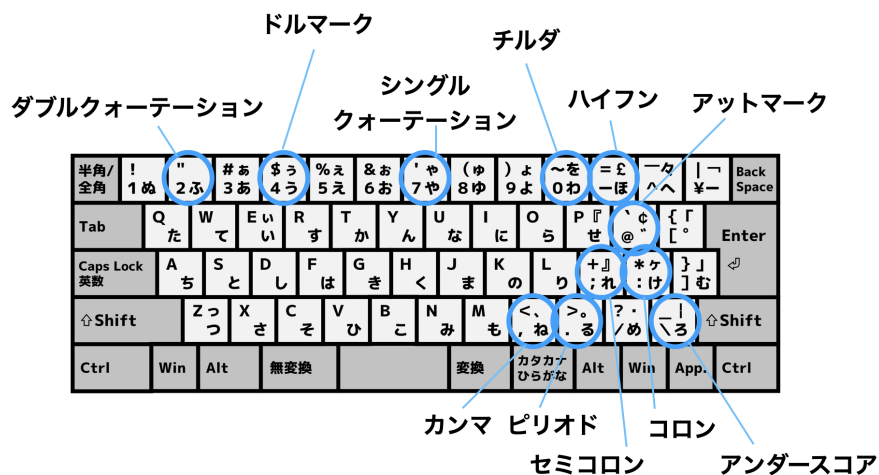


図 E.2 代表的な記号と日本語キー配列

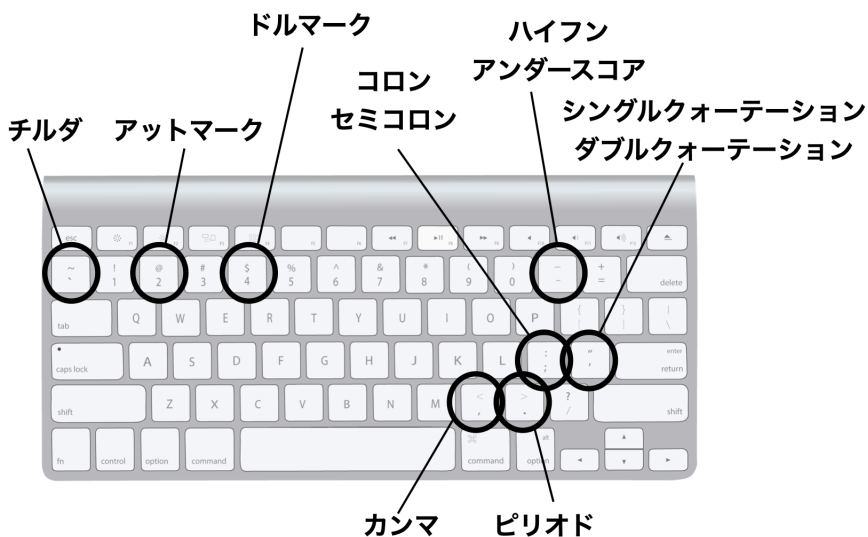


図 E.3 代表的な記号と US キー配列

9435 付録 F

9436 本講義に対応する詳細シラバス

9437 F.1 イントロダクション

9438 F.1.1 授業内容

9439 科目中でのこのコマの位置づけ

9440 この講義の位置付けは、基礎的な心理統計の学習は終わった後の応用的内容となる。基礎的な内容とし
9441 て、確率の基本的な考え方、線形モデル(回帰分析、群間の平均値差の検討)、さまざまな推定法による母数
9442 の推定と検定の考え方を理解しているものとする。これに基づいての応用であるから、扱うデータも単変量で
9443 はなく多変量であるし、数学的には行列表現を用いることになる。これらを使って、回帰分析や因子分析の数
9444 理的理解を目指す。このコマではこの講義によって扱われる領域を外観するとともに、基礎的な内容で扱っ
9445 たものがしっかりと定着しているかどうかを確認することを目的とする。

9446 コマ主題細目

9447 **正規線形モデルの世界** 単変量ではなく多変量を扱う統計の領域に入るので、多変量データとはどのよう
9448 なものであるかに言及した上で、本講義の扱う領域を概観する。心理統計の応用的分野では、正規分
9449 布を仮定した線形モデルがその大半を占めている。正規線形モデルに含まれるさまざまな下位モデル
9450 の名称を紹介するとともに、構造方程式モデリングに統合されることや、非線形なモデルとの違いにつ
9451 いて理解する。加えて正規分布ではない分布を扱うモデルも増えてきている昨今、これらについての
9452 モデリングアプローチの存在についても講義する。

9453 → 正規線形モデルの枠組みについては、[三中 \(2018\)](#) 参照。

9454 **尺度の四水準** 心理統計の基礎で触れたが、データ化として扱う数値はその尺度水準によってどのような計
9455 算が可能かということに違いが生じる。このことは、そのまま分析モデルや名称の違いに繋がるため、
9456 改めて名義、順序、間隔、比率の4水準を確認しておく。

9457 → [Stevens \(1946\)](#) の論文は短く、ネットで読むこともできる。入門書としては[川端・荘島 \(2014\)](#) の
9458 Pp.9-16, あるいは[山田・村井 \(2004\)](#) の Pp.22-25.

9459 **平均と分散** 間隔尺度水準以上の数字であれば、平均値や分散、標準偏差によってその特徴を要約でき
9460 る。ここではこれらの代表値の表記について、数学的記号とともに確認する。加えて、分散式を展開し
9461 て表現したものや、分散がデータから得られる情報の上限であることを確認する。

9462 **共分散と相関係数** 共分散やそれを標準化した相関係数は、複数の変数間関係を表現する最もシンプルな

9463 ものの 1 つである。ここではこれらの複数の変数間に関わる代表値について、数学的記号とともに確
 9464 認する。加えて、この他の関係の表現方法として、距離や共頻度などの共変量について解説し、それら
 9465 の違いに応じて統計モデルが変わりうることを確認する。

9466 → 記述統計量については川端・荘島 (2014) の Pp.26-33 など基礎的な心理統計の教科書を参照
 9467 すると良い。

9468 キーワード

- 9469 • 正規線形モデル
- 9470 • 尺度水準
- 9471 • 平均と分散
- 9472 • 共分散と標準偏差

9473 F.1.2 授業情報

9474 ■コマの展開方法 講義

9475 予習・復習課題

9476 ■予習 心理統計の基礎について、今一度基礎的なテキストを参照しながら、自分の理解度を再確認してお
 9477 くと良い。とくに尺度水準や記述統計量の計算方法などは今後この講義でも頻出するので、確認しておく必
 9478 要がある。

9479 ■復習 数式の展開を踏まえて理解しておくとともに、実際のデータを使って計算しながら確認すると良
 9480 い。とくに分散は二乗のオーダーになるので元の単位に比べて大きな数字になること、標準化のプロセスや相
 9481 関係数の大きさなど、逐一確認しておくべきである。

9482 F.2 心理尺度を作る

9483 F.2.1 授業内容

9484 科目中でのこのコマの位置づけ

9485 目に見えないものを測定するために心理学が洗練してきた手法が、心理尺度である。心理尺度の作成方
 9486 法としてサーストン法、リッカート法、SD 法などがあり、心理学の初頭コースで習うものも少なくないが、その
 9487 本質は反応カテゴリに数値を割り当てる、というところにある。「そう思わない」「ややそう思わない」などとい
 9488 たカテゴリに対する反応が、なぜ 5 や 4 といった数字にすることが許されるのか。カテゴリカルな反応が連
 9489 続的な量として扱うことができる理由などについて、よく知られていない現実がある。この点をしっかり理解し
 9490 ないまま進んだ分析を行うと、結果の解釈はもちろんそもそもの研究が足元から崩壊することにもなりかねな
 9491 い。本講義ではこの点について、作成方法から数値化まで一通り確認し、最後に心理尺度の評価基準である
 9492 信頼性と妥当性について理解する。

9493 コマ主題細目

9494 サーストンの等現間隔法 サーストンと呼ばれる尺度構成法は、態度とよばれる心理学的特性を仮定し
 9495 ている。この態度は対象、符号、強度をもち、正規分布すると仮定されている。個々人の態度を測定す

9496 るために、事前に評定者集団を用意して項目を採点しておく必要がある。そこで評定値に等間隔性を
9497 持たせる工夫をしているため、態度の数値化ができるという原理を理解する。

9498 → サーストン法による尺度作成については末永 (1987) の Pp.149–152 参照。

9499 **リッカートのシグマ法** リッカート法は最もよく使われるスタイルの心理尺度である。カテゴリーに無頓着
9500 に数字を割り振る慣例がみられるが、本来は潜在的態度が正規分布することを想定し、確率分布の確
9501 率点を得点とする方法であった。このような数値化がされているからこそ、順序尺度ではなく間隔尺
9502 度水準と「見なす」ことが許されてきているのである。この原理を理解しておくことは、後のより進んだ
9503 尺度作成法を理解する助けになる。

9504 → リッカート法による尺度作成については、宮谷・坂田・林・坂田・入戸野・森田 (2009) の
9505 Pp.150–153 を参照。ただしシグマ法についての言及はなく、田中 (1977) などの古典を当たらねば
9506 ならない。

9507 **尺度を評価する** 作られた尺度を評価する方法として、IT 相関を求めるものや内的整合性信頼性を求め
9508 るものがある。その背後には信頼性と妥当性の考え方があることを確認する。

9509 → 心理尺度の信頼性については、末永 (1987) の Pp.156–158, 妥当性については Grimm and
9510 Yarnold (2001) の第 4 章も参照。

9511 キーワード

- 9512 • サーストンの等現間隔法
- 9513 • リッカートのシグマ法
- 9514 • 信頼性
- 9515 • 妥当性

9516 F.2.2 授業情報

9517 ■ コマの展開方法 講義

9518 予習・復習課題

9519 ■ **予習** 基礎実習で尺度作成法や尺度の分析をしたことがあれば、その時の資料を再確認しておく。とくに
9520 反応カテゴリをどのように採点したか、また尺度の評価どのように行ったかを確認する。とくに尺度作成の経
9521 験がない場合は、関連書籍を参考に方法論を予習しておくことが望ましい。

9522 ■ **復習** IT 相関やアルファ係数は統計環境 R で簡単に計算できる。とくに psych パッケージにはこれらの
9523 関数がすでに準備されている。サンプルデータを使ってこれらを計算してみよう。

9524 F.3 テスト理論と因子分析

9525 F.3.1 授業内容

9526 科目の中でのこのコマの位置づけ

9527 目に見えないものを測定するという意味で、テスト理論は心理学の測定と関係が深い。前回は社会心理学
9528 における態度の測定を前提に議論したが、測定に関してはテスト理論で一般的に議論できる。

9529 真のスコアと誤差とに分解すること、誤差の基本的な仮定を確認した上で、古典的テスト理論を項目と被験
9530 者の特性に分割することで因子分析モデルに展開されるところを見る。また、多因子モデルに拡張した上で、
9531 その数理的展開から、信頼性と妥当性に言及できることを確認する。数式の展開は代数の基本的な特徴を確認
9532 すれば問題なくフォローできるはずである。

9533 コマ主題細目

9534 **古典的テスト理論** 古典的テスト理論についての復習である。その基本モデルについて触れ、その平均値と
9535 分散が意味するところから測定モデルの意味するところ(誤差が相殺しあうこと)と信頼性の定義が導
9536 出できることを改めて確認しておく。

9537 **因子分析モデル** 因子分析モデルは、古典的テスト理論のモデルを拡張したものである。まずは単因子モデル
9538 を例に、項目特性と被験者特性が分離されたことを確認する。その上で、性格検査や知能検査など
9539 の歴史に触れながら、多因子モデルについて解説する。多因子モデルを例に記号や添字を確認して
9540 おく。

9541 → 小杉 (2018) の Pp.173–177

9542 **因子分析の第 2 定理** 因子分析モデルを展開することで、相関係数が因子負荷量の積和で表現できるこ
9543 と、因子得点とその仮定から計算上消えることを確認する。得られた指揮は因子分析の第二定理と呼
9544 ばれ、妥当性に関する議論がここから導かれることをみる。

9545 **因子分析の第 1 定理** ある項目自身の相関係数を考えることで、因子分析の第一定理にたどり着く。ここで
9546 共通因子の二乗和を共通性と呼ぶことにすると、信頼性の考え方が項目レベルで行われるように発展
9547 したことが確認できる。

9548 → 小杉 (2018) の Pp.173–177

9549 キーワード

- 9550 • 古典的テスト理論
- 9551 • 因子分析法
- 9552 • 因子分析の定理

9553 F.3.2 授業情報

9554 ■コマの展開方法 講義

9555 予習・復習課題

9556 ■**予習** 一年時に信頼性・妥当性について、あるいは古典的テスト理論について学んだことを復習し、どのよ
9557 うな概念であったかを再確認しておくことが望ましい。

9558 ■**復習** テスト理論と因子分析モデルの関係について、因子分析モデルはどこが新しく何を改定しようとした
9559 のかについて、自分なりの言葉で説明できるようになろう。

9560 F.4 現代テスト理論

9561 F.4.1 授業内容

9562 科目の中でのこのコマの位置づけ

9563 テストの理論も目に見えないものを測定するという意味では、心理学と同じモデルを実践する領域である。
9564 心理学的尺度作成法の発展には、テスト業界における理論的展開の位置付けを知ることが役に立つ。

9565 因子分析によって項目と被験者の特徴を分離して考えることができるようになった。ここで学力テストに目
9566 を向けると、単因子でよいことと従属変数がバイナリになっていることがわかる。

9567 この特殊な測定法についてのモデルを考えるために、まずは通過率の概念を導入したうえで、累積正規
9568 分布とその近似としてのロジスティック曲線、および 1,2,3PL モデルを紹介する。これらのテストは新しいテ
9569 スト理論とよばれるが、それはこれまでのテスト理論に含まれていた集団に依存した測定であったこと、完全
9570 データに限定されていたことなどを乗り越えられるからである。もちろんテストの等価がしやすいという側面
9571 もある。

9572 テスト理論の展開としての項目反応理論と、因子分析モデルとの相同性を強調することで、見えないものを
9573 測定しようとするアプローチという意味では同じであったことを確認する。

9574 コマ主題細目

9575 **因子分析とテスト理論** 単因子モデルの特殊事例として、学力テストの例を考える。学力テストの性質か
9576 ら、因子構造よりも因子得点に注目するという強調点の違いはあるが、因子分析モデルの一環として
9577 捉えることを強調する。

9578 → [高橋 \(2002\)](#) は最も平易なテスト及び現代テスト理論への入門書である。最初の数ページだけで
9579 も参考になる。

9580 **通過率と累積正規分布** 学力テストの分析例として、通過率の計算から累積正規分布へとつなげる。累積
9581 正規分布をそのまま確率モデルに繋げてもよいが、ロジスティック曲線を使う方が関数の形が簡単で
9582 あり、こちらの方が実際には使い勝手が良い。ベルヌーイ分布を用いた線形回帰モデルの文脈で考え
9583 れば、ロジスティック回帰分析をしていることでもあることに言及する。

9584 → 通過率については[豊田 \(2012\)](#) の Pp.1-8 を、ロジスティック関数と累積正規分布の関係につい
9585 ては[加藤他 \(2014\)](#) の Pp.81-83 を参照

9586 **項目母数の特徴** ロジスティック曲線を導入することで、関数の変形がたやすくなった。ここでは 1PL,2PL
9587 ロジスティックモデルを導入し、どの項目母数が関数の位置や形をどのように変えるか、そしてそれが
9588 意味するところを理解する。モデル的には 5 母数モデルまで考えられるが、実際にはせいぜい 3PL
9589 モデルである。この講義では後の因子分析との対応関係も考えるため、2PL モデルまでの紹介に留
9590 める。

9591 → 項目母数については豊田 (2012) の Pp.31-34, 加藤他 (2014) の Pp.71-80 が参考になる。

9592 **被験者母数の特徴** 項目母数が明らかになった状況に置いて、どのように被験者母数を推定するかを考え
9593 る。ここで ICC から逆算的に被験者母数の位置がどこにあるか、該当領域を絞り込んでいく尤度関
9594 数を視覚的に確認する。この方法を使うと、すべての項目についての回答が得られていないと推定で
9595 きないといった不便がなく、また被験者母数の位置によっては ICC がそれほど有用な情報を与えてく
9596 れないこともある。これらの点は、完全情報最尤推定や情報関数にもつながるため、しっかりと理解し
9597 ておくことが必要である。

9598 → 被験者母数の絞り込みについては、小杉・清水 (2014) の Pp.171-172. が参考になる。

9599 キーワード

- 9600 • 通過率
- 9601 • ロジスティックモデル
- 9602 • 被験者母数の推定について

9603 F.4.2 授業情報

9604 ■コマの展開方法 講義

9605 予習・復習課題

9606 ■予習 テストの前提となる標準正規分布について復習しておく。とくに R を使って出力できる確率密度、
9607 確率点、累積確率など手を動かして予習しておくが良い。

9608 ■復習 適当なグラフ描画ツール (R でよい) をつかって、ロジスティックモデルを描写し、項目母数をどのよ
9609 うに変えるとどのように曲線の形が変わるかを確認してみよう。

9610 F.5 現代テスト理論その 2

9611 F.5.1 授業内容

9612 科目の中でのこのコマの位置づけ

9613 項目反応理論の数学的特徴を踏まえることと、現代的尺度構成法の理論的基礎を学ぶ。

9614 項目反応理論の導入によって、被験者母数と項目母数が完全に分離され、項目の特徴を細かく記述でき
9615 るようになった。また、項目の特徴がわかればテストの実践方法も変わってくる。ひとつは CAT に代表
9616 されるように、ダイナミックに出題を変化させることができるようになること、そうしたうえでもテストの平均点
9617 が事前にコントロールしうることなどが示される。項目から得られる情報という観点から項目情報曲線が、項
9618 目情報曲線の累積からテスト情報曲線が導出される。

9619 つづいてこのテスト理論の発展形として、多段階モデルに拡張可能なことをみる。とくに段階反応モデル
9620 は、リッカートのシグマ法のように段階反応をモデルかできるという意味で、現代的リッカート法であるともい
9621 える。段階反応モデルを用いることで、適切な反応段階のチェックをできるなど、応用的側面が高いことを確
9622 認する。

9623 またテスト理論は因子分析の特殊系であるという扱いだだったが、多段階、多因子へと展開することで再び

9624 因子分析モデルに統合されていくことを確認する。

9625 コマ主題細目

9626 **現代テスト理論の特徴** 現代テスト理論の特徴は、項目母数と被験者母数の分離、完全情報最尤推定、項目情報曲線による信頼性の表現、項目プールがあれば事前にテストの平均点を設計できることがあげられる。また Computer Adopted Test の形式を用いることでテストのあり方そのものも変わってしまう。ただし実際には、膨大な項目プールが必要であること、事前に項目母数を準備しておく必要があること、その他「公平性のために新しいテストでなければならない」という信念などが弊害となって実践的には敷居が高いことなどを解説する。

9632 → 古典的テスト理論との比較については、加藤他 (2014) の Pp.67–69、あるいは豊田 (2012) の前書きが十分に詳しい。

9634 **段階反応モデル** テスト理論はバイナリデータに対する分析だが、多段階の反応に拡張する方法がいくつか考えられている。1 つは段階反応モデルとよばれるもので、これを使うと適当な反応段階数がデザインできるなど利点は大きい。またその考え方はリッカートのシグマ法を洗練したものであるとも言え、せめてこうした方法を使わないと多段階反応を適当に分析できていない。統計パッケージなどの実装も進んでいるので、計算コストはほとんど障壁にならない。また、ポリリック相関係数を用いた因子分析を実行すると、段階反応モデルのパラメータに変換できることから、因子分析とテスト理論が同じものであったことを再確認できる。

9641 → 豊田 (2012) の Pp.155–172 が詳しい。

9642 **因子分析の歴史と展開** 因子分析モデルもテスト理論も潜在変数モデルとしては同じであり、一方が単因子・二段階、他方が多因子・多段階であることが道を分つ。またその性質から、一方が因子得点に、他方が因子構造に着目するため、テストの構成についての考え方が異なることにも注意する。繰り返しになるが、統計パッケージ上の実装は進んでいるので、どちらを使うにしてもとくに苦勞することなく、積極的にカテゴリ軽モデルを推進していくべきである。

9647 キーワード

- 9648 • 項目情報曲線、テスト情報曲線
- 9649 • 段階反応モデル
- 9650 • 因子分析モデルとテスト理論

9651 F.5.2 授業情報

9652 ■コマの展開方法 講義

9653 予習・復習課題

9654 ■**予習** 項目反応理論、とくに 2PL モデルによる因子得点の算出方法を確認しておくと同時に、心理尺度ではどのように尺度値を定めていたかについて復習しておく。

9656 ■**復習** 信頼性についての考え方が、古典的テスト理論、因子分析論、現代テスト理論を通じてどのように変わってきたかを確認しておこう。

9658 F.6 行列計算の基礎

9659 F.6.1 授業内容

9660 科目の中でのこのコマの位置づけ

9661 テスト理論や因子分析モデルの展開を理解した上で、さらに次のステップに進むためには、より数学的な構
9662 造の理解が必要である。ここまで因子分析モデルでは、因子得点をどのように算出するかが論じられていな
9663 い。また相関行列を分解して因子負荷量を算出するにあたって、どのように計算するかについては言及さ
9664 れてこなかった。これらの点を理解するための道具となるのが線形代数である。具体的には、行列の固有値
9665 分解を通じた解釈をすることで、因子分析、回帰分析など多変量データの方程式モデルを統一的に表現・理
9666 解できるようになる。そのための道具立てとして、線形代数の基礎知識を習得する必要がある。本講はこのよ
9667 り進んだ理解に向かうための、新しい数学ツールの導入を行う。

9668 線形代数は方程式を簡便的に表現するための表現法であり、行列の観点から新たに四則演算を定義し直
9669 すことで一般的な表現が可能になることを示す。

9670 コマ主題細目

9671 **行列とベクトル** 多変量データを行列とベクトルで表現することをみる。学ぶべき用語として、スカラー、縦
9672 ベクトル、横ベクトル、行列、正方行列、対称行列、対角行列、単位行列をあげる。

9673 **行列の四則演算** ベクトルとベクトルの和、行列と行列の和、スカラーとベクトルの積、スカラーと行列の積、
9674 縦ベクトルと横ベクトルの積、横ベクトルと縦ベクトルの積、行列と行列の積をみる。とくにサイズが変
9675 わることに注意が必要である。

9676 **行列による便利な表現** 連立方程式が行列で表現できることを見る。

9677 **逆行列と連立方程式** 行列の割り算に当たるのが逆行列である。逆行列は存在しないこともあるが、もし適
9678 当なものが見つかればそれは連立方程式の解を一気に計算ができることになる。

9679 キーワード

- 9680 • ベクトル, スカラー, 行列
- 9681 • 行列の四則演算
- 9682 • 連立方程式

9683 F.6.2 授業情報

9684 ■コマの展開方法 講義

9685 予習・復習課題

9686 ■予習 とくに予習の必要は感じないが、授業に参加するにあたってはノートの準備が必要である。

9687 ■復習 計算方法に慣れておく必要があるので、練習問題を繰り返して行うことで、とくに行列の積の計算
9688 ができるようになっておく。線形代数の入門書としては、数学のテキストとして読みづらさを感じるかもしれな
9689 いが、村上他 (2016) がよく、一冊手元に置いて演習をしながら進めると良い。

9690 F.7 行列による関係の表現

9691 F.7.1 授業内容

9692 科目の中でのこのコマの位置づけ

9693 線形代数についての基礎的なルールを習得する段階である。今回はより実践的・具体的に、データ行列を
 9694 どのように線形代数で表現できるかを考える。データ行列から分散共分散行列、相関行列へと形を変えるこ
 9695 とを学ぶ。つづいて線形モデル、とくに従属変数が明確な回帰分析モデルを行列で表現することを見、線形モ
 9696 デルとデザイン行列について考える。さらに因子分析モデルを行列で表現することを考える。行列で表現する
 9697 ことで、1つの式の中に第一、第二定理の両方を含んだ形で表現できることを理解する。

9698 コマ主題細目

9699 **データの行列表現** 実際に手にするデータセットは、表計算ソフトウェアの画面で見える行列形式の数値であ
 9700 るが、これを記号で表現することで一般的に扱うことができるようになる。添字に気をつけながら要素
 9701 ごとの表示をすることに加え、行列の計算をこのデータ行列に与えることによって、変数の平均や変数
 9702 ごとの平均偏差を持った行列が表現できる。平均偏差行列を用いると、行列の積の特徴から分散と共
 9703 分散を含んだ正方行列が作られることがわかる。また、データを標準化することで、標準化された行列
 9704 の積が相関行列を表すことになる。このように一般的に表現するために、これまでの行列計算の方法
 9705 が作られたのだと逆算的に理解すること、加えて行列のサイズに注目しながら、扱うデータの大きさが
 9706 イメージできるようになることが肝要である。

9707 → 岡太 (2008) の Pp.77-110

9708 **線形モデル** 行列表現の利便性は、データの変換だけにあるのではなく、統計モデルを表現する際にも生き
 9709 てくる。基礎で学ぶ線形モデルは、基本的にエレメントワイズな表記法であったが、行列を使うことで
 9710 単回帰も重回帰も同じ式で表現できることがわかる。このように表記の統一性があることが、線形代
 9711 数の利点である。また統一的な表記にするために、切片項にかかる列を追加するなどの工夫をするこ
 9712 ともにも注意する。これらの点は、R など統計ソフトウェアを扱う上でもヒントになることが多い。

9713 **デザイン行列** 基礎の段階で行った帰無仮説検定は、説明変数が離散変数であったことから、線形モデル
 9714 の特殊形に過ぎなかったことを再確認する。その上で、先の回帰分析を行列表記にしたように、離散
 9715 変数で説明する時の係数にかかる行列の形を確認する。この行列はとくにデザイン行列と呼ばれるこ
 9716 と、また自由度の関係から制約を加えた表現になるが、それがデザイン行列の中でどのように書き表
 9717 されるかを確認する。

9718 → J.Dobson (2008) の Pp.41-45 にごく簡単な紹介が、豊田 (2000) の Pp.47-62 には計画行列
 9719 として構造方程式の枠組みで説明されている。

9720 **因子分析モデルの行列表現** 因子分析モデルはここまでエレメントワイズで表現されていたが、同様に行列
 9721 表現にするとどのようになるかを確認する。とくに行列のサイズに注目することが重要である。というの
 9722 も、統計ソフトウェアを使っていると因子得点が表示されないことが少なくないが、行列の形で見ると
 9723 因子負荷量は項目数 × 因子数、因子得点は回答者数 × 因子数になることがより意識されやすいか
 9724 らである。他にも因子分析に関する特徴量が行列のどの要素にはいつているか、また因子分析の定理
 9725 が行列のどこで表現されているかを確認することが重要である。

9726 因子分析の行列的表現については → 芝 (1979) が良書だが、現在は絶版。同様の内容は小杉
9727 (2018) にもある。

9728 キーワード

- 9729 • データの行列表現
- 9730 • 分散共分散行列, 相関行列
- 9731 • デザイン行列

9732 F.7.2 授業情報

9733 ■コマの展開方法 講義

9734 予習・復習課題

9735 ■予習 行列の掛け算がメインになってくるので、計算方法並びに計算結果のサイズを確認する方法を見て
9736 おこう。

9737 ■復習 行列表現によって重回帰方程式が 1 つの形になることを確認する。平均値の差を見るために線形
9738 モデルが用いられることを確認する。また因子分析モデルを行列表現すると、一気に 2 つの定理が 1 つの式
9739 で表現できることを確認する。

9740 F.8 固有値と固有ベクトルと因子分析モデルの関係

9741 F.8.1 授業内容

9742 科目の中でのこのコマの位置づけ

9743 因子分析モデルを行列表現することで、いよいよ因子をどのように算出しているのかについての答えが明
9744 らかになる。

9745 因子負荷量を算出するためには、線形代数でいうところの固有値についての理解が必要である。まずは固
9746 有値と固有ベクトルを導入し、どのように計算するかを見る。とくに固有ベクトルはノルムが定まらないことを
9747 確認する。そこから、固有値と固有ベクトルがどのような性質を持っているかを幾何学的観点から確認する。
9748 正方行列が座標変換を行うためのものであると考えれば、固有ベクトルは変換行列の基底となることがわか
9749 るだろう。データ解析にあたって、相関行列の基底を求めるとはどういうことかをイメージするだけでも、因子
9750 分析の理解がまた一歩深まるだろう。

9751 コマ主題細目

9752 **固有値と固有ベクトル** 行列の固有値と固有ベクトルの性質を理解する。直感的には、正方行列がスカ
9753 ラーに変わることが、情報圧縮になっていると言えるだろう。また、行列のサイズと同じ数だけ固有値
9754 が見つかること、固有値の総和が元の行列のトレース trace になることを確認する。とくにデータ解析
9755 の領域では、分散共分散行列か相関行列が分析対象になることが基本であり、こうした対称行列の固
9756 有値は実数になること、相関行列のトレースは項目数と合致することを改めて確認することで、データ
9757 の情報圧縮になることについての直感的理解をめざす。

9758 **固有ベクトルを求める** 2×2 行列を例に、固有値と固有ベクトルを求める計算を行う。固有方程式を導入

9759 し固有値の計算を行うことは比較的簡単であるが、固有ベクトルの求め方が直感的にはわかりにく
 9760 い。というのも、固有ベクトルはその大きさが定まっておらず、要素同士の相対的な大きさを示すだけ
 9761 だからである。ここでベクトルのノルムを導入して標準化解を算出することを確認する。また行列のサ
 9762 イズが大きくなると方程式が高次になるため、一般解が得られないこと、結果的に近似解を求める計
 9763 算方法が開発されていることをみる。

9764 **固有値と固有ベクトルの幾何学的意味** 正方行列は一次変換行列であり、固有ベクトルはその基底であ
 9765 ることを単純な行列から理解する。固有ベクトルはノルムが定まっていないこと、すなわち方向性だけ
 9766 を持ったものであることを理解する。また固有値はその総和が元の行列のトレースと一致することか
 9767 ら、分散あるいは項目数 (相関行列の対角) を組み替えたものであり、固有値の大きさの順に考える
 9768 ことはすなわち、より明確な次元を抽出したことになることを確認する。

9769 → これについては平岡・堀 (2004) にアニメーション付きで説明されているのがわかりやすい。また長
 9770 沼 (2011) は固有値の章だけでなく、付録を読むとまた固有値と固有ベクトルの多角的な理解が
 9771 進む。

9772 **因子分析モデルの意味** 因子分析モデルは相関行列を固有値分解することであり、それはすなわち相関行
 9773 列の中にある基本的な次元・座標を求めることにある。すなわち複数人の反応パターンの共通要素を
 9774 取り出すということであり、これは心理学的アプローチをほぼ直接的に数学表現したものであることを
 9775 理解する。座標の回転についても触れ、仮定を緩めた場合の表現も理解する。

9776 キーワード

- 9777 • 因子分析モデルの行列表現
- 9778 • 固有値
- 9779 • 固有ベクトル
- 9780 • 固有ベクトルの幾何学的理解

9781 F.8.2 授業情報

9782 ■コマの展開方法 講義

9783 予習・復習課題

9784 ■予習 因子分析の基本モデル, 第一・第二定理の導出を復習しておこう。

9785 ■復習 因子分析モデルが何をやっているかを考えた上で、心理学における尺度の利用やその解釈におい
 9786 てどのような注意をしなければならないかを言語化してみよう。

9787 F.9 R をつかっての行列計算

9788 F.9.1 授業内容

9789 科目の中でのこのコマの位置づけ

9790 行列の計算は単純な計算ではあるが、要素の数が多くなるので反復回数が増え、また計算の法則も慣れ
 9791 るまでは難しい。人間にとってはミスが多くなりがちなこの計算が、計算機 (コンピュータ) は最も得意とする

9792 ところである。計算機は疲れることなく、単純な反復計算を瞬時にこなす。多変量データ解析は計算機の発展
9793 の歴史ともあり、昨今の計算機パワーは非常に複雑な統計解析も瞬時に答えを出すようになった。

9794 この行列計算は表計算ソフトにはできないことであり、統計環境 R のような、統計パッケージを利用するこ
9795 とになる。本項では、統計環境 R を用いて行列の基本的な計算を演習によって習得することを目的としてい
9796 る。また R で行列の計算ができることは重要ではあるが、実際に統計分析をする時にはより便利なパッケー
9797 ジを利用することになる。心理学関係の数値計算については、psych パッケージが便利である。これを導入
9798 し、記述統計量や信頼性係数など基本的な分析が便利になることを確認する。

9799 コマ主題細目

9800 R による行列計算 R についての基本的な使い方 (環境の準備, RStudio によるプロジェクト管理, パッ
9801 ケージの導入, 基本的な四則演算等) については習得済みであることを前提とする。行列計算にあ
9802 たっては、データをマトリックス型で保持している必要があり、また行列の計算は四則演算と異なるこ
9803 と、ベクトルの長さが時には再利用されることなど注意が必要な点がある。それらを踏まえて、データ
9804 の方を考えながら行列の四則演算を確認する。

9805 R によるデータの変換 R の行列計算を使って、前時までに行った raw data の変換計算, すなわち平均,
9806 平均偏差行列, 分散共分散行列, 相関行列などの計算プロセスを確認する。また, cov や cor 関数
9807 を使うとこれらが一気に計算されるが, 分散の関数には不偏分散が用いられていることに注意する必
9808 要がある。

9809 R による固有値計算 R の eigen 関数を使って, 固有値と固有ベクトルが計算される場所を確認する。固
9810 有ベクトルは標準化されていることに注意する。

9811 キーワード

- 9812 • 行列型
- 9813 • 行列関数

9814 F.9.2 授業情報

9815 ■コマの展開方法 R を使った演習

9816 予習・復習課題

9817 ■予習 R/RStudio を使った分析環境を再確認しておこう。またデータの読み込みや記述統計量などの算
9818 出関数を確認しておこう。

9819 ■復習 授業時間内に収まらなかったところがあれば、必ずキャッチアップしておくこと。いくつかの練習問題
9820 を実践し、エラーや警告がでてでも対応できるようになろう。

F.10 Rをつかった因子分析と尺度作成法

F.10.1 授業内容

科目の中でのこのコマの位置づけ

ここでは心理尺度を開発するような心理学研究を想定し、より実践的な順序に則って演習を進めていく。この講義の目標は、自らが質問紙調査を使った研究をした場合にどのような手順で行うかを理解し、実践できるようにすることである。具体的には前回導入した `psych` パッケージを用いて、さまざまな推定オプションを追加していくことで出力が変わっていくことを確認しながら進める。

コマ主題細目

psych パッケージ概説 心理学研究に用いられる便利な関数群である `psych` パッケージのマニュアルを見ながら、`describe.by` などの記述統計量関数、`alpha` や `omega` といった信頼性係数の関数を使ってロウデータの分析を行う。

調査研究の手順 心理尺度の作成研究の手続きを外観する。まず構成概念の設定、定義、妥当性を考えた上で、具体的な項目を選出し、テストデータを取る。探索的な因子分析によってその因子的妥当性を確認し、標準化のための本調査を行う。あるいは1次元性を確認した上で、IRTによって反応段階の確認、項目母数の確認、テスト情報関数の確認などが必要である。尺度の翻訳や検証的妥当性のチェックなどについては、構造方程式モデリングによる分析を行うのでここでは扱わず、参照するにとどめる。

共通性推定の問題 分析にあたって、改めて因子分析モデルの行列表現を提示し、行列の固有値分解によって因子負荷量が求められることを確認する。しかしその際、共通性をどのように推定するかの問題が残されていたことを確認し、そのためにいくつかの方法が提案されていることを理解する。これらは因子分析を行う上で、推定方法のオプション指定に関わってくる点であり、ソフトウェアが変わっても同様の指標が必要であることをみる。

→ [小杉 \(2018\)](#) の Pp.91–94.

fa 関数と探索的因子分析 探索的因子分析の手続きを `fa` 関数を使いながら考える。探索的因子分析の場合は因子構造、因子負荷量について何ら前提を置かないため、因子数の推定から始めなければならない。まずは `fa.parallel` 関数でスクリープロットを描画する。スクリープロットを読むときの形状について確認する。続いて因子数と共通性推定方法を定めた上で `fa` 関数を実行し、因子負荷量や共通性などアウトプットを確認する。続いて解釈を簡単にするために因子軸の回転を行うことを解説し、実行のために `rotate` オプションを追加することをみる。回転前の結果との比較、また直交回転と斜交回転の違いを確認する。

→ [小杉 \(2018\)](#) の Pp.81–91.

因子得点の算出 因子数と因子負荷量が明らかになると、そこから逆算的に因子得点を計算できる。`fa` 関数には `scores` オプションをつけることで、出力されたオブジェクトから因子得点を取り出すことができるのを見る。こうした方法とは別に、項目同士の素点の平均から因子得点を計算することもある。これは推定値を実体とすることの懸念が出発点であり、その長所と短所を把握しておくことが必要であ

9857 る。この簡便法は平均値情報を含んでいるため、尺度カテゴリに依拠した解釈が可能である。また取
9858 り出された因子得点と簡便的因子得点の相関を見ることを確認する。

9859 **因子分析の注意点** 因子分析を行う上で注意しなければならないのは、因子が実体としてあるのではなく、
9860 あくまでも準備された項目群の相関関係から得られる基底に過ぎないことを理解する点である。因子
9861 分析の流れの中では因子に命名することが 1 つの手順としてあるが、言葉として確定するとあたかも
9862 それがあるかのように考えられてしまうこと、それしかないように考えられてしまうことの危険性を理解
9863 する。心理尺度の呪いやてっちゃんの手品になってしまわないように注意し、常に元の項目群に戻って
9864 考える必要があることをしっかりと理解する。

9865 キーワード

- 9866 • 信頼性係数
- 9867 • 共通性
- 9868 • 因子負荷量
- 9869 • 因子得点
- 9870 • psych パッケージ
- 9871 • fa, fa.poly, fa.parallel 関数

9872 F.10.2 授業情報

9873 ■コマの展開方法 R を使った演習

9874 予習・復習課題

9875 ■予習 パッケージの読み込みや関数の結果を見る方法を確認しておこう。一年時のことを思い出して、1m
9876 関数を例に R の操作方を思い出しておく。

9877 ■復習 授業時間内に収まらなかったところがあれば、必ずキャッチアップしておくこと。いくつかの練習問題
9878 を実践し、エラーや警告に対応できるようになろう。心理学研究など心理学の専門雑誌を参考に、どのような
9879 分析結果がどのように報告されているかを確認しておくことも、理解を進める。

9880 F.11 R をつかった項目反応理論

9881 F.11.1 授業内容

9882 科目の中でのこのコマの位置づけ

9883 項目反応理論を実践的に理解する演習パートである。

9884 カテゴリカルな因子分析と数学的に同等ではあるが、より項目の特徴を広く表現できる項目反応理論の利
9885 用が、今後より重要なものになってくるだろう。

9886 ここではまずテスト理論の根本に立ち返り、二値単因子のデータを使って 1PL, 2PL モデルの分析を行う。
9887 分析結果は数値で見るとも重要であるが、ICC や IIC, TIC などを使って可視化するとより理解が深ま
9888 るだろう。多段階の反応についても、同様に GRM を実行し、閾値や識別力、IRCCC や IIC, TIC が描画
9889 できることを確認する。とくに IRCCC による反応段階の読み取り方には注意する。最後に多段階で多因子
9890 の場合、項目反応理論の文脈から言えば多次元 IRT になり専用のパッケージが必要になることを紹介しつ

9891 つ, カテゴリカル因子分析でも同様のことができることを確認する。

9892 コマ主題細目

9893 **項目反応理論の実際** 項目反応理論はテスト理論がその出自に当たるので, まずは二値データで単因子が
9894 想定できるような例を元に分析を行う。分析には `irtosys` パッケージや `ltm` パッケージを用いて,
9895 1PL モデル, 2PL モデルの演習を行う。項目母数の値と意味が, 具体的な設問に照らし合わせて考
9896 えることで, より実感をもって理解できるようになると思われる。とくに, ICC や IIC, TIC など可視化
9897 することでその意味が理解しやすくなるだろう。3PL モデルなどさらに拡張したモデルも利用可能で
9898 ある。

9899 **段階反応モデルの実際** 続いて単因子, 段階反応モデルの実践を行う。因子構造として, 前回の授業で扱っ
9900 た多因子の内, ある因子に限定して分析を行うこととする。段階反応モデル (GRM) の出力結果を数
9901 値だけでなく可視化することで, 項目の特徴がどのように表現されているかを考える。とくに反応段階
9902 の山が潰れているようなケースは, 適切な反応段階でなかったことを意味するので, 数値の置き換え
9903 など元データを修正しつつ分析し直すことを考える。これらを通じて, 適切な反応段階による調査法が
9904 必要であることを理解する。

9905 **カテゴリカル因子分析との対応** 多段階, 多因子の場合は `psych` パッケージの `fa` 関数にオプションを追
9906 加することでできる, カテゴリカル因子分析と同じである。出力結果について, これまでの相関係数を
9907 用いているものとの違いを確認する。また数値をどのように変換すれば対応するのかを見ることで, 数
9908 学的に等価であることを確認しておく。IRT の側面から多因子に拡張した, 多次元 IRT についても,
9909 `mirt` パッケージを利用すれば実行できる。この解析には計算時間がかかるが, 完全情報最尤推定の
9910 結果が得られることは利点である。

9911 キーワード

- 9912 • 1PL モデル, 2PL モデル, 3PL モデル
- 9913 • 段階反応モデル
- 9914 • カテゴリカル因子分析
- 9915 • `irtosys`, `ltm`, `psych`, `mirt` パッケージ

9916 F.11.2 授業情報

9917 ■コマの展開方法 Rを使った演習

9918 予習・復習課題

9919 **■予習** `irtosys`, `ltm` パッケージを事前にインストールして環境を整え, データファイルの読み込みなど
9920 R/RStudio の基本的な使い方を確認しておこう。とくに R の `data.frame` 型に含まれる変数が, `numeric`
9921 なのか `factor` なのかによって挙動がかわることがある。変数の型についても再確認しておこう。

9922 **■復習** 本講で習ったパッケージを使って, 具体的なデータを因子分析, IRT, カテゴリカル IRT などいく
9923 つかのモデルで分析し, それぞれの違いを確認しておこう。

9924 F.12 構造方程式モデリング

9925 F.12.1 授業内容

9926 科目の中でのこのコマの位置づけ

9927 構造方程式モデリングは、因子分析と回帰分析を統合して扱う、総合的分析モデルである。言い換えれば、
9928 これまでの多くの多変量解析モデルのほとんどは、構造方程式モデルの下位モデルとして表現できる。ここで
9929 はこれまでのモデルを統合した、より現代的でより上位のモデルである構造方程式モデリングを理解すること
9930 で、すべての多変量解析を網羅的かつ俯瞰的に捉えることが狙いである。

9931 構造方程式モデリングを理解するには、変数の種類と関係性の区分に注意したパスダイアグラムの描き方
9932 を知ることが早い。パスダイアグラムを用いると、回帰分析と因子分析は説明変数が観測変数なのか潜在変
9933 数なのかといった違いであることが明らかである。また因子分析と似た主成分分析がどのように表されるか
9934 も、パスダイアグラムを見れば一目瞭然である。

9935 パスダイアグラムには変数の尺度水準までかきこまれることはないが、ここに注意していろいろなモデルを
9936 描画すると、それがかつて多変量解析においてさまざまな名称で呼ばれた分析方法であったことがわかる。
9937 あるいは、今後どのようなモデルが開発される可能性があるか、どのようなモデルをどのように希釈すれば良
9938 いかもイメージできる。

9939 加えてこの統合的なモデルがなぜそうした複雑なモデルを表現できるのかについても、モデルを方程式で
9940 描画し、行列で考えることで、モデル行列とデータ行列を近づけることと理解できる。この観点から、データに
9941 モデルを当てはめる適合度の考え方が改めて理解されるだろう。

9942 コマ主題細目

9943 **パスダイアグラムの書き方** これまで学んできたモデルを図で表現することを学ぶ。そのためには、変数を
9944 観測変数と潜在変数に区別することと、変数間関係を因果関係と相関関係に区別する必要がある。
9945 観測変数を矩形、潜在変数を楕円形、因果関係を一方向矢印、相関関係を双方向矢印で表現すること
9946 で、因子分析や回帰分析が図で表現できることを学ぶ。

9947 **パスダイアグラムによるさまざまなモデル** 因子分析と回帰分析をパスダイアグラムで表現したことで、
9948 この両者を統合するような表現ができること、また潜在変数同士の関係を記述する、構造方程式を描
9949 画できるようになったと言える。因子分析と似た手法とされる主成分分析や、尺度水準の違いによるさ
9950 まざまな統計モデルを表現する方法を手に入れたことになる。この手法を総称して、構造方程式モデリ
9951 ングと呼ぶ。

9952 → 小杉・清水 (2014) の Pp.7-10

9953 **構造方程式モデルによる未知数の推定** 構造方程式モデリングでは、パスダイアグラムでも表現されるが、
9954 変数間関係を方程式で書くこともできる。方程式で描画することで、構造方程式も潜在変数の方程式
9955 と観測変数の方程式、それらが入れ子になった方程式で描画できることがわかる。またこれらのモデル
9956 を行列のイメージで捉えると、最終的には分散共分散行列という実態を持った数字に対して、未知
9957 数で描画された方程式を接続したことが直感的にわかるだろう。未知数の増え方と分散共分散行列
9958 の要素の増え方を比較すると後者が圧倒的に早く、未知数よりも既知数が多い方程式は解くことがで
9959 きるという原理から、未知数が推定しうることを理解する。

9960 → 小杉 (2018) の Pp.191-193.

9961 **適合度によるモデルの評価** データ行列とモデル行列をイコールで結んだ方程式を解くことが、未知数を求
 9962 める根本的な原理であるが、このことからモデルがデータとどの程度合致しているかという適合度が、
 9963 モデル評価の統合的観点として浮かび上がってくる。回帰分析では R^2 であったが、因子分析をはじ
 9964 めとしたさまざまな多変量解析モデルも、この評価次元で考えることができる。ただしその指標にはい
 9965 くつかの特徴があり、これらを総合的にみて評価するという実践的ノウハウも確認する。

9966 → 小杉・清水 (2014) の Pp.10–12

9967 キーワード

- 9968 • パスダイアグラム
- 9969 • 観測変数と潜在変数
- 9970 • 因果関係と相関関係
- 9971 • 構造方程式モデリング
- 9972 • 適合度

9973 F.12.2 授業情報

9974 ■コマの展開方法 講義

9975 予習・復習課題

9976 ■予習 回帰分析と因子分析という2つの分析方法についてはすでに学んでいるが、この両者の共通点と
 9977 相違点どこにあるかを事前に考えてみよう。外的な基準の有無、説明変数の種類の観点から、自分の言葉
 9978 で表現できるようになっていると良い。

9979 ■復習 これまで学んださまざまな統計モデルを、構造方程式モデリングの表記法に則ってパス図を書いて
 9980 みよう。またさまざまな尺度水準の組み合わせからなるモデルを考え、それらがどのような意味を持つのかと
 9981 推論するのも理解の助けになる。

9982 F.13 Rによる構造方程式モデリング

9983 F.13.1 授業内容

9984 科目の中でのこのコマの位置づけ

9985 これまでの流れと同じで、統計技術の理論を知っただけではなく、自分で実際に計算できる演習を経てこ
 9986 そ理解が深まるということから、本講ではRをつかって実際に構造方程式モデリングを解くことを演習的に
 9987 学ぶ。構造方程式モデリングを実装するパッケージは複数あるが、最も応用範囲がひろいlavaanパッケー
 9988 ジを用いることにする。

9989 まずは観測変数だけからなる簡単なパス解析を行う。データの入力の仕方、方程式の設定、関数の使い方
 9990 などを一通り習得する。続いて潜在変数を含んだモデルによる解析を行う。モデルの適合度や修正指数を参
 9991 考に、徐々にモデルを書き換えていく手順を学ぶ。注意すべきは、適合度を上げることが目的になって、不
 9992 自然な仮定やパスをおいてしまうことである。あくまでも具体的かつ妥当なモデリングを心がけるべきである。

9993 オプションな設定になるが、観測変数がカテゴリカルである場合や、推定方法の選択なども確認する。最
 9994 後に、R以外の統計パッケージによる構造方程式モデリングの実践例がいくつか紹介される。

9995 コマ主題細目

9996 **方程式の入力** まずは観測変数同士の関係をパスでつなぐモデルで練習する。パス解析は回帰分析の繰り返しで実行することもできるが、構造方程式モデリングによってパスの繋がりを 1 つのモデルで表現し、適合度も統一できるなどの利点がある。観測変数だけからなるモデルの結果と、実際に 1m 関数で実行した結果と比較すると良い。またパッケージにもよるが、自動的にパスダイアグラムを描画してくれるものもある。方程式とそのパスダイアグラムによる表現の対応を確認する。

10001 → 小杉・清水 (2014) の Pp.55-60

10002 **測定モデルの実践** 因子分析モデルを SEM 上で実行してみる。探索的因子分析と違い、どの項目にどの因子が影響しているかを固定したモデリングが可能であり、この検証的因子分析による結果と、いわゆる因子分析関数との結果を比較することで、パスが引かれていないところはその係数が 0 であるという強い仮定をおいていることを確認する。また尺度作成の観点からは、検証的因子分析をすることで因子的妥当性や弁別的妥当性、収束的妥当性を検討することもできる。さらに同じモデルを別のデータに適用することで多母集団同時分析を行うことになる。このように、モデルの暗黙の過程や、モデルとデータの適合という側面とくに注意する。さらに測定モデルと測定モデルをつなぐ、構造方程式を扱ったモデルへと拡張する。

10010 → 小杉・清水 (2014) の Pp.87-90

10011 **実践上の注意点** ここまでを通じて、一通りモデルを作成できるようになった。とくに構造方程式を踏まえると、心理的実体同士の関係を描画したと解釈できるため、心理学的概念間の関係を記述できることは魅力的に映るかもしれない。しかしデータを越えての解釈はご法度であり、潜在変数が心理的実在であるかどうかの議論は、理論的背景や測定の適切さ、標準化されないスコアが実際にどのように変化すれば何が言えるのか、といったところに一足飛びに行かぬよう注意する必要がある。またモデル改良のステップにおいて、適合度や修正指数を過度に参照していないか、注意する必要がある。

10017 **そのほかの統計パッケージ** 構造方程式モデリングの利点は、モデルを可視化したことにもある。たとえば AMOS は GUI でモデルを作成できる。他に Mplus はカテゴリカルな変数にも対応しているし、高度に複雑なモデルであっても表現が可能である。

10020 キーワード

- 10021 • 測定方程式
- 10022 • 構造方程式
- 10023 • 潜在変数を含んだモデル
- 10024 • 多母集団同時分析
- 10025 • 適合度

10026 F.13.2 授業情報

10027 ■コマの展開方法 講義

10028 予習・復習課題

10029 ■**予習** 構造方程式モデルは、数式レベルでの理解は難しいが実際は統合的なものであり、回帰分析や因
10030 子分析をその下位モデルとして含んでいる。改めて、回帰分析や因子分析を単体で行った場合にどのような
10031 出力がなされるのか確認しておく、同じものを構造方程式で実践したときの違いが明確に意識できるように
10032 なる。

10033 ■**復習** さまざまなモデルを試すなかでは、エラーや警告が出ることもある。そうしたエラーや警告の意味を
10034 理解し、またそれに対応するためにはどのような方法が取れるかを考える必要がある。まずは手元のデータを
10035 用いて、これらの練習を行うと良い。

10036 F.14 双対尺度法

10037 F.14.1 授業内容

10038 科目の中でのこのコマの位置づけ

10039 ここまでは尺度化によって与えられた数値を元に、因子構造を検証したり（潜在）変数間関係を記述したり
10040 することをみてきた。ここでは尺度化によって与えられる数値に改めて注目し、得られた情報を最も有効に活
10041 用するように数字を割り振る、数量化の考え方へと考え方を進める。

10042 因子分析や構造方程式モデリングで用いられる相関関係や共分散関係は、いずれも与えられた数字がど
10043 の程度の直線的関係にあるかという観点からモデルを組み立てていた。そのモデルの表現力は非常に豊富
10044 なので、スタートとなる変数同士の直線性（相関係数を使うこと）を改めて問うことがなかった、あるいは多少
10045 の違和感を抑えて進んでしまうことがある。改めてその線形性に注意を払い、逆に線形性を最大にするよう
10046 に数値をデータに付与するという発想を転換させるのが、数量化の手法である。

10047 数量化はいくつかの種類があるが、ここでは III 類を取り上げる。III 類は双対尺度法や対応分析とも呼
10048 ばれ、行と列の両方に数字を付与する。この方法によってさまざまな表現ができることを確認し、これが応用
10049 的側面ではテキストマイニングなどにも利用されているところを見る。こうした応用方法は臨床場面でも活き
10050 るものであり、どのようなデータにどのように応用できるかを考えることで、データに対して積極的に関わる姿
10051 勢を身につけることが期待される。

10052 コマ主題細目

10053 **直線的でない関係** 心理学的には中庸が良いことも少なくないが、であればカテゴリに付与される尺度値か
10054 らは U 字あるいは逆 U 字の関係がえられる。この関係は相関係数としては 0 に近くとも、無関係、無
10055 意味を表すものではない。そこから意味が取り出せないのであれば、数値化のルールを修正するべき
10056 であって、どのように関係性の高い数値を与えるかについては、行または列の平均からカテゴリの値
10057 を付け直すことである。ここでの目的は、直線性を取り出すためにカテゴリに数値を与える数量化とい
10058 う別の観点であることに注意する。

10059 → [西里 \(2010\)](#) の Pp.6-25.

10060 **林の数量化理論** 分析者にとって最も有用な情報が得られるようにカテゴリに数値を与える、という観点を
10061 数量化と呼ぶ。数量化の対象は離散変数や名義尺度水準の変数であってもよく、これらを用いた分析
10062 方法については、行動計量学の祖である林知己夫の多彩な手法をまとめた呼称である数量化 I 類、
10063 II 類、III 類などがある。I 類、II 類は重回帰分析や判別分析とも関係するが、数量化という観点か
10064 ら議論されていることに注意が必要である。数量化 I 類、II 類の応用例を外観することで分析方のイ
10065 メージを掴む。

10066

→ 小杉 (2019b) の Pp.189-195.

10067

双対尺度法による分析 数量化 III 類は開発者によって双対尺度法や対応分析など異なる名称を持つが、数学的には等価でいずれも目指すところは名義尺度水準の主成分分析である。リッカート法での尺度に値をつけた時のように、行及び列に含まれるカテゴリカルな区分に対してもっとも線型性が高くなるように数値を割り当てる。ここで行列計算にその目を向ければ、矩形行列に対する特異値分解によって、行の空間と列の空間を用意し、カテゴリに座標を与えることを意味していることになる。数量化 III 類、双対尺度法、対応分析の細かな違いにも注意しつつ、行カテゴリと列カテゴリを共通空間に表した図から何が読み取れるかを考える。

10074

→ 小杉 (2019b) の Pp.195-199.

10075

テキストマイニングへの応用 数量化理論の対象が名義尺度水準であることを考えれば、およそ言語化できたものはすべて多変量解析として分析できることになる。逐語録や自然言語の解析にはテキストマイニングと呼ばれる手法が用いられるが、この技術は形態素解析と多変量解析の組み合わせであり、多変量解析の元になる共変動が何で表されているかに着目すれば、統合的に解釈することが可能である。テキストマイニングには専門的なソフトウェアがあるが、軽く言及するにとどめる。

10080

テキストマイニングについては → 樋口 (2020) を参照せよ。

10081

キーワード

10082

- 数量化

10083

- 双対尺度法

10084

- 特異値分解

10085

- テキストマイニング

10086

F.14.2 授業情報

10087

■コマの展開方法 講義

10088

予習・復習課題

10089

■予習 構造方程式モデリングでさまざまなモデルを考えた際、変数が観測・潜在の別だけでなく尺度水準の組み合わせも変えて考えた場合、どのようなモデルがありえるかが想定できると思います。今回はとくに名義尺度水準のモデルになるので、名義尺度水準のモデルはどのようなものがあるかを想定してみてください。

10091

10092

■復習 名義尺度水準のモデルを手に入れたことによって、さまざまな分析の可能性が広がったのではないのでしょうか。今までおよそ統計的・数値的アプローチの対象にないと思われていたものに対しても、どのように、どこまでであればアプローチ可能で、どこからが限界になるのかを考えることが実践的には役に立つ視点です。統計モデルを使いこなすために、ぜひさまざまな応用例を考えてみてください。

10093

10094

10095

10096 F.15 多次元尺度構成法

10097 F.15.1 授業内容

10098 科目の中でのこのコマの位置づけ

10099 数量化まで学ぶことによって、カテゴリに適切な数値を割り振るという尺度化の原点に立ち返ることができ
10100 た。ここではさらに、心理尺度のような回答法ではない変数間同士の関係から、次元を取り出して分析する多
10101 次元尺度構成法について考える。この方法は実験刺激や知覚的反応、直感的判断などを対象にできるため、
10102 応用可能性が非常に高いだろう。

10103 多次元尺度構成法を理解するためには、まず実際の距離行列を分析して地図を再構成できるかどうかを
10104 見るところから始めるのが良いだろう。データとして与えられるのが距離行列であり、行動計量学ではこれを
10105 心理的な距離や意味的な距離が数値化されたものだと捉えることで、心理学的な地図を作っていると解釈し
10106 てきた。この仮定には最大限の注意を払いつつ、必ずしも計量的でない場合の数値化をする非計量的多次
10107 元尺度法に拡張することで、更なる心理学的用途が広がることを見る。なお、多次元尺度方は数量化 IV 類
10108 と同じである。

10109 多次元尺度法は分析の元が距離行列であり、その基底を固有値分解によって得るといいうみで、多変量
10110 解析としてはお馴染みの考え方であるともいえる。しかしデータが距離（を意味するもの）であれば良く、数
10111 学的にも簡単な拡張をすることで、個人差を表現するモデルに拡張することもできる。また心理尺度に対する
10112 考え方として、個人の内的な次元からの近さに応じて反応すると考える、展開型のモデルを使うことは、心理
10113 尺度の利用に新たな視点をもたらす。

10114 コマ主題細目

10115 **多次元尺度構成法** 多次元尺度構成法 (MDS) は、距離行列を元にした多変量解析の一種であり、距離関
10116 係から次元 (基底) を選び出し、対象に座標を与える方法である。まずは座標の復元例から考え、抽
10117 出する次元数をどのようにして求めるかといった基礎的な知識を得る。また、行列の考え方からみると
10118 正方対称行列の固有値分解であるから、これを確認するだけでも他の多変量解析と合わせた統合的
10119 な理解ができるだろう。因子分析モデルも多次元尺度法の一種であるということもできる。

10120 **距離と心理学のデータ** 距離行列があれば次元が抽出できることが分かったが、さて何を距離とみなすか
10121 を考えれば、非常に多くの可能性が広がるのがわかる。距離の定義は非負で対称性と三角不等式
10122 が成り立つことであり、さまざまな距離の定め方があるし、共分散や相関もその一種と考えることがで
10123 ける。元になるデータも尺度評定を用いる方法、刺激の混同率、代替価/連想価、刺激の汎化勾配、
10124 反応潜時、ソシオメトリックなデータなど、心理学のさまざまな領域で得られるデータが、距離とみなす
10125 ことができる。応用可能な領域が広いことを知ることで、統計モデルをハンドリングできることになる。

10126 → データの例に関しては高根 (1980) の Pp.14-27.

10127 **非計量多次元尺度法** 計量 MDS によって算出される座標は元のデータをうまく復元するが、心理学的な
10128 データの場合はデータの大小関係の表現 (順序尺度水準) がせいぜいであり、これに対応した非計量
10129 多次元尺度法が考案されている。この手法を用いることで、一対比較や順序比較などのより制限の少
10130 ないデータからであっても数量関係を導き出すことができる。

10131 → 計量・非計量多次元尺度構成法については、すでに絶版になったが岡太・今泉 (1994) がもっとも
10132 簡潔でわかりやすく説明している。手に入るところでは足立 (2006) の P.135-143, あるいは小杉

10133

(2019b) の P.199-203.

10134 **多次元尺度構成法の展開** 多次元尺度構成法で作られたものは地図である。地図には点を書き込んだり、
 10135 複数の地図を重ねたりできるように、多次元尺度法で得られた地図にも情報を追加したり、モデルを
 10136 展開するなどしてさまざまな応用的モデルを作ることができる。ここでは Prefmap や楕円モデル (非
 10137 対称 MDS), INDSCAL など応用例をいくつか示し、この技術の応用可能性を考える。

10138 **展開法** Coombs が考えた心理尺度の展開法は、サー斯顿法やリッカート法とはまた別の尺度化の考え方
 10139 を表している。この方法は被験者とカテゴリーの両方を地図上にプロットできる。具体的な分析例をみ
 10140 ながら、尺度やそれに数字を与える方法についての考え方を見る。

10141 → 展開型モデルについては、多少複雑な工夫が組み込まれているが、清水 (2018) が参考になる。

10142 キーワード

- 10143 • 多次元尺度構成法
- 10144 • 非計量多次元尺度構成法
- 10145 • 個人差多次元尺度構成法
- 10146 • 展開型多次元尺度法

10147 F.15.2 授業情報

10148 ■コマの展開方法 講義

10149 予習・復習課題

10150 **■予習** 数量化の関係を再確認しよう。すなわち数値に値を与えるというものであり、尺度のカテゴリだけで
 10151 なくより一般的な心理的刺激を考え、どういったモデルで表現できるかを考えると本講だけでなく後期の授業
 10152 にもつながる気づきを得るだろう。

10153 **■復習** 多次元尺度法によって、どのような分析ができるかを考えてみよう。とくに尺度法にかかわらず、実
 10154 験刺激からの反応を距離と見做せる関係にすれば分析でき、その結果をどのように解釈するかについても自
 10155 由度はかなり多い。紹介された発展的なモデルなどについても自分の研究関心にどのように応用できるか考
 10156 えてみよう。

10157 F.16 プログラミングの基礎

10158 F.16.1 授業内容

10159 科目の中でのこのコマの位置づけ

10160 後期の授業の主眼は「データから意味のある情報を取り出す」ことにある。すなわち、これまでのデータ駆
 10161 動型 (Data Driven) な発想を逆転し、データがどのように生成されたと考えるかという観点から、データ生
 10162 成モデル (Data generating Model) による理解を試みる。

10163 モデルは数学的表現がなされ、モデルの形成やデータとの照合 (fitting) には計算機の利用が必須であ
 10164 る。後期の初回となる今回の目的は、プログラミングの基本的な考え方を身につけ、次回以降の本格的な運
 10165 用に備えることである。プログラミングの基礎はコマンドによる命令であり、コマンドの書き間違いはエラーと

10166 なって帰ってくる。一見不親切に思えるが、即時反応により学習の効率は良く、コード補完機能などを用いる
 10167 ことで簡単なミススペルは回避することができる。小さなものから大きくしていくこと、1行ずつ実行すること
 10168 など、基本的な姿勢についても理解する。

10169 またプログラミング言語（高級言語）はいくつかあるが、文法的な違いを除けばその本質は代入、反復、条
 10170 件分岐である。この三点について理解し、基本的な書き方を学ぶ。実際に R でコードを書きながら、その挙動
 10171 について確認する。最終的な到達段階として、Fizz-Buzz 問題や行列計算ができるコードを書くこととする。

10172 コマ主題細目

10173 **プログラミングの基礎** プログラミングにあたって重要なことは「思った通りに」動くのではなく、「書いた通
 10174 りに」動くことである。ミススペルや大文字小文字の違いにも注意が必要である。コードのスペルチェッ
 10175 クや補完機能を活用し、また 1 行ずつ確認しながら進めるといったプログラミングに関わる基本的な
 10176 心構えについて理解する。

10177 **いくつかのプログラミング言語** R はプログラミング言語の一種であると言っても良い。プログラミング言
 10178 語には他にも Python や Basic, C 言語などがある。これらの基本的な関係について理解するととも
 10179 に、コンパイラとインタプリタという実行形式の違いについても理解する。この違いは今後確率的プロ
 10180 グラム言語を利用する際の知識として生きてくる。

10181 **高級言語の基本的な働き** 高級言語と呼ばれるプログラミング言語の基本的な働きは、代入、反復、条件分
 10182 岐である。R では<-や=で代入を、for や while で反復を、if や if_else で条件分岐を行う。こ
 10183 れらの表現は言語間を通じて共通なものが多いため、その基本的な振る舞いを確認することは技術
 10184 の一般化に役立つ。

10185 キーワード

- 10186 • プログラミング言語
- 10187 • コンパイル
- 10188 • 代入
- 10189 • 反復
- 10190 • 条件分岐

10191 F.16.2 授業情報

10192 ■コマの展開方法 講義/遠隔/演習

10193 予習・復習課題

10194 **■予習** 事前に環境の準備をしておく必要がある。環境の準備についてはいくつかの方策があり、これにつ
 10195 いては導入資料を参照しながら準備しておくこと。なお、環境準備中に問題が生じた場合はいち早く教員か
 10196 TA に相談し、実行できるようにしておくこと。

10197 **■復習** 反復計算の練習課題、条件分岐の練習課題など、複数の課題にしっかり取り組むこと。

10198 F.17 データ生成メカニズムとモデリング

10199 F.17.1 授業内容

10200 科目の中でのこのコマの位置づけ

10201 モデリングアプローチには実質的にベイズ推定が必須であり、そのためにはまず推定法としてのベイズ法の
 10202 位置付け、ベイズの定理の基礎、実践的方法としての MCMC 法について理解する。とくに、従来の心理統
 10203 計ではモーメント法と最尤法に言及されるにとどまるものが多い。確率モデルとしての表現は最尤法から導入
 10204 されており、尤度による推定を事後分布という確率分布で表現するようになったものとして、ベイズの定理を
 10205 位置付けると良い。MCMC 法は新しい方法であるが、その前に確率分布の特徴を記述するために乱数を
 10206 利用することができる、という事実を確認すれば、事後分布の記述との相性の良さも明らかである。このよう
 10207 に、ベイズ法を過度に新規でまったく異なるものであるという印象を与えることなく、従来の方法の延長線上
 10208 にあるものとして考えられるようにする。

10209 コマ主題細目

10210 **データ生成モデリング** 推測統計学が母集団分布における仮定から母数を推定することを目的とし、モー
 10211 メント法、最尤法、ベイズ法といったアプローチで推定を行うものであったことを再確認する。その上
 10212 で、帰無仮説検定など心理学一般で使われているモデルは得られた結果のみに基づいてモデル比較
 10213 をする形である。これは客観性を重視しデータにいかなる前提も置かないことを考えてのアプローチ
 10214 であり、いわばデータ駆動型の分析方法である。これに対して、データがどのように生まれてきたかを
 10215 考え、その仕組みに基づいて分析する方法がデータ生成モデリングである。データ駆動型分析法は実
 10216 験方法にそのメカニズムを埋め込んでおり、データ生成モデル駆動型分析法はメカニズムそのものを
 10217 検証する形になっている。

10218 →松浦 (2016) の Pp.18–25.

10219 **ベイズ推定の基礎** データ生成モデル駆動型分析にあたっては、未知なるパラメータが多くなるため、実質
 10220 的にベイズ推定を使うことになる。ベイズ推定の原理となる確率やベイズの公式を再確認することが
 10221 必要である。

10222 →Kruschke (2014) の Pp104–123 ほか枚挙に遑がない

10223 **MCMC** ベイズ推定は事後分布を用いて分析結果を考えることになり、これは結果がパラメータの確率分布
 10224 として得られることを意味する。そこで確率分布を分析する方法として、乱数を用いてアプローチする
 10225 ことを考える。乱数を用いるアプローチの利点は、積分計算が記述統計量で済むこと、周辺化分布に
 10226 ついても当該変数について考えるだけで済むこと、精度を上げる時にサンプル数を増やすだけで良い
 10227 ことなどが挙げられる。また乱数を用いるアプローチに対応したソフトウェアとして JAGS や Stan が
 10228 挙げられる。これらが確率的プログラミング言語が乱数発生機であることを踏まえ、簡単な実践を行っ
 10229 てみる。

10230 →Kruschke (2014) の Pp.147–194.

10231 **乱数によるアプローチの例** 簡単な確率分布を使って、乱数によるアプローチを実践する。サンプルサイズ
 10232 を増やすことで精度が上がること、記述統計量が確率分布の特徴を記述することを再確認する。とく

10233 に平均値, 中央値, パーセンタイル, 分散や標準偏差など, 分布の要約統計量の計算について, R ス
10234 クリプトで算出できるよう練習する。

10235 キーワード

- 10236 • データ生成モデリング
- 10237 • ベイズの定理
- 10238 • マルコフ連鎖モンテカルロ法
- 10239 • 乱数による近似

10240 F.17.2 授業情報

10241 ■コマの展開方法 講義/遠隔/演習

10242 予習・復習課題

10243 ■予習 環境の準備が整っていないものは, 急ぎ準備を行うこと。また 1 年次の心理学データ解析基礎にお
10244 いて, 確率分布や乱数を用いるアプローチについても言及はしているので, 振り返って「確率という数字の公
10245 理」を確認しておくことが望ましい。

10246 ■復習 R を用いて乱数を発生させ, 理論上の値の近似値になることを確認する。正規分布, ベルヌーイ分
10247 布, 二項分布, ポアソン分布など複数の既存関数を確かめることで, 一般化されない事後分布についての乱
10248 数発生機が手に入ったことの重要性に気づくことができるだろう。

10249 F.18 ベイジアンアプローチと確率的プログラミング 1

10250 F.18.1 授業内容

10251 科目中でのこのコマの位置づけ

10252 ベイズ推定によるモデリングアプローチを始めるにあたって, 比較的親近性の高い正規分布を用いた例を
10253 導入する。とくに分散/標準偏差に個人差を入れるモデルを用いることで, 平均値以外にも推定モデルが考え
10254 うることから, モデリングの対象とする領域の広さを意識させる。また確率的プログラミング言語を用いた初の
10255 演習でもあるから, コードの書き方, コンパイルなどの手順, 出力結果の診断と解釈など, 今後の分析に必要な
10256 技術的要素についても確認する。

10257 コマ主題細目

10258 7人の科学者 [Lee and Wagenmakers \(2013\)](#) より「7人の科学者」の例を紹介する。この例の利点として
10259 次の3点がある。1. 正規分布を用いていること, 2. 平均値以外のパラメータを用いていること, 3.
10260 小サンプルからの推論であり, ごく簡単なコードで実演できること。カバーストーリーからモデルを想像
10261 し, コードに落とししていくプロセスをたどりながらモデリングの実際を学ぶ。

10262 →[Lee and Wagenmakers \(2013\)](#) の Pp.48–50.

10263 Stan コードの書き方 カバーストーリーに沿ったモデル図(設計図)ができれば, 後は Stan の言語仕様に
10264 そって記述していただくだけである。ブロックによる分割, セミicolonによる一行の終わり, 変数の宣言と利

10265 用など言語仕様を外観したあとで、モデルブロックからパラメータ、データと逆順に書いていく書き方
10266 を試す。尤度と事前分布の違いにも注意し、コメントをつけながらコードを書いていく。

10267 **Stan を使った MCMC の実践** Stan コードは Stan ファイルに記載し、コンパイルは Stan そのものを用い
10268 るものである。これらファイル、インターフェイス (パッケージ)、実行エンジンなどの関係を明らかにし
10269 つつ、MCMC について学ぶ。

10270 →[Kruschke \(2014\)](#) の Pp.407–425.

10271 **MCMC の結果の診断** MCMC の結果は Stanfit オブジェクトとして得られるが、それがどう言った情報
10272 を持っているか、また MCMC の代表性、正確性、代表性など確認すべき点について理解する。

10273 →[Kruschke \(2014\)](#) Pp.180-194

10274 **MCMC の結果の解釈** 事後分布からの代表値として適切な MCMC 結果が得られたら、それを用いて結
10275 果の解釈を行う。結果はすべて分布として得られているので、確率分布をどのように代表するかにつ
10276 いて、前時の復習をしながら進める。また複数のパラメータの同時分布で結果が得られていること、す
10277 なわちあるパラメータの点推定値の組み合わせが、同時分布の適切な代表になっていない可能性に
10278 も注意が必要である。

10279 同時分布については →[Lee and Wagenmakers \(2013\)](#) の Pp.42–46.

10280 キーワード

- 10281 • 精度のモデリング
- 10282 • Stan
- 10283 • data ブロック, parameters ブロック, model ブロック
- 10284 • Rhat, 有効サンプルサイズ
- 10285 • EAP, MAP, MED
- 10286 • 同時分布

10287 F.18.2 授業情報

10288 ■コマの展開方法 講義/遠隔/演習

10289 予習・復習課題

10290 ■**予習** Stan を使う初めての演習になるので、プログラミングの心得 (思った通りに動くのではなく、書いた
10291 通りに動く) を再確認しておこう。

10292 ■**復習** データサイズや MCMC のサンプルサイズを変更するなどして、同じ乱数生成機の何をどう変えれ
10293 ばどう変化するのかを確認しておこう。

10294 F.19 モデリングの目から見た検定 1 ; 二群の平均値の差

10295 F.19.1 授業内容

10296 科目の中でのこのコマの位置づけ

10297 データ生成モデリングの観点を踏まえた上で、検定的アプローチとベイズ的アプローチの違いを学ぶ。

10298 二群の平均値の差の検定、いわゆる対応のない t 検定の場合は、同一の正規分布から得られたデータに
10299 対して平均値の差があると判断して良いかどうかを判断するという枠組みであった。これらの前提と判断基
10300 準を確認し、それがデータ生成モデルの観点ではどのように表されるかを検証する。ここで結果が分布として
10301 推定されること、差があるかないかというのは二群の推定された平均値の差であることから、生成量を使って
10302 平均値の差を出力することを考える。ここで効果量に改めて目を向けるとその理解が進む。

10303 コマ主題細目

10304 **t 検定の仮定** 二群の平均値の差を検定するときは t 検定が利用されるが、データが正規分布から得られ
10305 ているという仮定、分散が同じであるという仮定などを踏まえて設計図を書き、これを Stan で表現す
10306 ることを考える。あらためて t 検定のやり方や結果と比べてみることで、モデリングがデータ生成メカ
10307 ニズムに注目していること、パラメータの推定を行なっていることなどが確認できるだろう。また等分散
10308 性の仮定を外す方法についてもすぐに応用ができる。

10309 帰無仮説検定と二群の差の検定について、→ [山田・村井 \(2004\)](#) や一年次の資料をもとに確認して
10310 おく。

10311 **差の分布** 検定は推測に加えて判断を行っていた、ということを改めて確認するとともに、帰無仮説検定
10312 では母平均の差をターゲットにしていたことを確認する。MCMC は母集団からの代表値であるの
10313 で、推定された結果を使って差を表現することができる。これは R 側で得られたサンプルで行っても
10314 良いし、Stan の生成量を使っても良い。ここで generated quantities ブロックの考え方を導入
10315 し、平均値の差の分布を確認すること、帰無仮説検定が差の分布の一点についての仮説であったこと
10316 を確認する。また一方が他方よりも大きくなる確率はどれぐらいかとか、一方と他方が c 以上に違っ
10317 ている確率はどれぐらいか、と言ったことが生成量を使って計算することができるようにもいえる。

10318 **帰無仮説検定を省みる** ここまでくると、帰無仮説検定のロジックや考え方について別の視点から見るこ
10319 とができるようになるであろう。まずは帰無仮説と対立仮説という対立のさせ方の不平等さである。帰無
10320 仮説は一点についての仮説であり、対立仮説はそれ以外であればなんでも良い、という非対称な関係
10321 になっていた。それを省みると、差があるかないかといった二値判断に陥ることがいかに危険であるか
10322 がわかるだろう。また量的な判断ができないことから、効果量を合わせて報告することが望ましいとき
10323 されている。効果量とは、標準化された差の大きさのことであり、生成量を使って簡単に算出することが
10324 できる。また方向性を持った検定について、片側・両側検定などで考えられてきたが、生成料を使えば
10325 自然にそれが検証できることがわかる。ただしこれらの検証の仕方は、今回のデータと仮定されたモデ
10326 ルという前提の上で成立する程度であって、過度な一般化にはならないように注意する必要がある。

10327 キーワード

- 10328 • t 検定
- 10329 • 生成量

- 10330 • 効果量
- 10331 • 片側検定, 両側検定

10332 F.19.2 授業情報

10333 ■コマの展開方法 講義/遠隔/演習

10334 予習・復習課題

10335 ■予習 Stan の基本的なブロック構成, 設計図からコードに落とすやり方を確認しておこう。

10336 ■復習 データのサイズが変わるだけでなく, 平均値の差, 効果量が変化したときの t 検定の結果とベイズ
10337 推定の結果がどのように変わるのか, さまざまなケースを想定して「遊んで」みるとよい。加えて, そのほかの
10338 仮説検定がどのようなデータ生成メカニズムで表現できるかを考えることは, 次回以降の準備にもつながる。

10339 F.20 モデリングの目から見た検定 2 ; パタメータの世界とデータの 10340 世界

10341 F.20.1 授業内容

10342 科目の中でのこのコマの位置づけ

10343 データ生成メカニズムの観点から帰無仮説検定を省みた場合, 拙速な結論に飛びつかずに慎重な議論が
10344 できることを確認した。また事後予測分布を作ることで, 柔軟な仮説を考えられることなども示された。

10345 本時は, 同じく事後予測分布を使いながら, パラメータでなくデータのレベルでの比較ができることに言及
10346 し, 実質的に差があるとはどういうことであるかを考える。パラメータの世界, データの世界を分けて考えられ
10347 るように注意を促す。まずは事後予測分布をみることで, モデルが現在のデータを正しく再現しているかを見
10348 ることで, 視覚的にモデルの正しさが検証できることを確認する。その上で, 新しく作られた分布の特徴から,
10349 閾上率や優越率などを計算することができる。これらはデータに基づく予測であるから, より具体的で実感を
10350 得やすい予測として使えるだろう。翻って, 仮想データを生成して検証する, パラメータリカバリの手法を学
10351 ぶ。この方法では真値やサンプルサイズを自由に設定し検証できることから, 例数設計に応用することが可能
10352 である。ベイズ推定をしない場合であっても, シミュレーションによる例数設計が有用であることを理解する。

10353 コマ主題細目

10354 事後予測分布 推定値をつかって, 新たにデータを生成した場合どのようなことが言えるか。事後予測分布
10355 を描くことでモデルの正しさが確認できる。事前分布の特徴を反映した, 事前予測分布についても触
10356 れる。

10357 事前と事後の予測については →[Lee and Wagenmakers \(2013\)](#) の Pp.38–42 が詳しい。

10358 データレベルの仮説 これまで考えてこられた仮説は, パラメータについての仮説であった。一方, 事後予
10359 測分布が新しいデータを作っているのであれば, そこからデータレベルの仮説を考えることもできる。
10360 閾上率や優越率といった, データのレベルでの仮説を検証したり検討したりすることを考える。これら
10361 の視点は帰無仮説検定およびモーメント法による算出よりもわかりやすいかもしれないし, 統計的に
10362 差があるということがどの程度意味のあることなのかを実感するのも役立つだろう。

10363 優越率, 閾上率については, → 豊田 (2016), Pp.69–70.

10364 **パラメータ・リカバリ** 事後予測分布は乱数発生による新しいデータの生成である。であれば乱数発生
10365 のアプローチは, 理論に従う仮のデータを生成することができる, ということでもある。シミュレーショ
10366 ンとして仮にデータを作ってみて, サンプルサイズがどの程度であればどの程度正確な推定ができる
10367 のか, と言った理論的な検証をすることができる。これはパラメータ・リカバリという試みでもあり, モデル
10368 が複雑になって行ったときに正しく機能するかどうかをチェックする方法でもある。また, サンプルサ
10369 イズも自由に変えることができるのだから, どの程度のサンプルがあればどのような結果が得られるの
10370 か, と言ったシミュレーション, あるいは実験前のサンプルサイズ設計にもつながる。

10371 モデリングの基礎的手順について, → 松浦 (2016) の Pp.12–81.

10372 キーワード

- 10373 • 生成量
- 10374 • パラメータ・リカバリ
- 10375 • 事後予測分布
- 10376 • 例数設計

10377 F.20.2 授業情報

10378 ■コマの展開方法 講義/遠隔/演習

10379 予習・復習課題

10380 ■予習 generated quantities ブロックの書き方について, 数値を色々変えて確認しておくといよ。

10381 ■復習 さまざまなサンプルサイズ, 効果量の仮想データを生成し, 帰無仮説検定やベイズ法による推定を
10382 繰り返すことで, 各手法の長所や短所を考えることができる。遊び心を持って, さまざまな状況生成して, 実際
10383 に試してみるこ。

10384 F.21 モデリングの目から見た検定 3 ; 多群の平均値差を求めるモ 10385 デル

10386 F.21.1 授業内容

10387 科目中でのこのコマの位置づけ

10388 ここまで generated quantities ブロックの活用で, 事後分布, 事後予測分布を生成できること, パラ
10389 メータの世界, データの世界それぞれでの仮説が検証できることを見てきた。またパラメータリカバリの方法
10390 を見ることで, データ生成モデルを分析前に活用する方法についても見た。

10391 続いて多群の平均値差を求めるモデルを考える。まずは一要因 3 水準の Between モデルから, 三つの平
10392 均値をバラバラに求めること, 生成量から差分を計算することを考える。データやパラメータレベルでの仮説
10393 的な検証方法を再確認し, 帰無仮説検定の枠組みで考えなければならなかったアルファ水準のインフレ問題
10394 が生じないことを, 確率の考え方の違いに即して理解する。続いて差をモデルに組み込む方法を考える。ここ
10395 でパラメータの数に制約をかける方法として transformed parameters ブロックの使い方を導入する。ま

10396 たパラメータの数の制約は、自由度の概念と深く関係していることへの洞察を得る。またパラメータリカバリの
10397 コードがリバースエンジニアリングのコードと同じもので、鏡合わせの関係にあることを確認する。続いて交
10398 互作用が含まれるモデルを考える。ここで制約からどれだけのパラメータが必要か、どのように組み上げるか
10399 を学ぶことができる。

10400 コマ主題細目

10401 **要因計画モデル** 一要因 3 水準 Between モデルを考える。対応のない二群の時のように、これは素直に
10402 3 群のモデルとして表現できるし、群間の差を生成量として計算できることを再確認する。また検定と
10403 違って差の大きさを直接検証すること、確率的判断を含まないことから、アルファ水準のインフレ問題
10404 に悩む必要がないことがわかる。この考え方は、確率の捉え方の違いにもつながることに留意する。

10405 **パラメータの変形と制約** 三群のうち、ある群を基準にした差分を直接パラメータとして推定するモデルを
10406 考えると、二つの差分を計算することができる。またある群を基準におかなくとも、全体平均を基準に
10407 置くことができるが、その場合は差分のパラメータに制約をかける必要がある。これらの制約を含んだ
10408 モデルを、transformed parameters ブロックを使って作ることを確認する。ここでパラメータの自
10409 由度について理解する。

10410 **モデルの洗練** 技術的な問題であるが、多群モデルの場合は一般的に描画するためにも、整然データの形
10411 式に整えておくことが望ましい。書いたモデルの一般化という観点から、変数にできるところは変数に
10412 するなど、コードの洗練を試みる。ここで群を識別する変数を導入することは、今後の個人内反復測定
10413 モデルにも応用できる点であるので、しっかり理解する。

10414 **パラメータリカバリ** これらのモデルのパラメータリカバリから、仮想データはデータ生成モデルを逆転さ
10415 せるだけで出来上がることがわかり、要因計画を裏側から眺めるような、新しい観点からの理解が進む
10416 と考えられる。

10417 キーワード

- 10418 • 要因計画
- 10419 • 整然データ
- 10420 • 生成量
- 10421 • パラメータリカバリ

10422 F.21.2 授業情報

10423 ■コマの展開方法 講義/遠隔/演習

10424 予習・復習課題

10425 ■予習 パラメータリカバリの必要性など、データ生成モデルのアプローチにおける標準的な手順を再確認
10426 する。また対応のない二群の平均値差の検定、要因計画の検定についてこれまでの復習をしておくことで、今
10427 回の内容の理解が深まるだろう。

10428 ■復習 要因数が増えた場合どのようになるか、またその都度 Stan モデルを書き換えなくても良くなるよう
10429 な一般的な書き方について、自分なりに試行錯誤することが望ましい。

10430 F.22 モデリングの目から見た検定 4 ; 対応のある群の比較

10431 F.22.1 授業内容

10432 科目の中でのこのコマの位置づけ

10433 ここまで Between 計画についてのモデル化を勧めてきたが、ここからは Within モデルについて考えるこ
10434 とにする。Within モデルは相関係数の考え方と、階層モデルへの入り口として位置付けることができる。

10435 Within モデルは対応がある場合であり、変数間に相関関係を想定することになる。相関関係のあるデー
10436 タ生成分布は、多変量分布であり、正規分布を多次元正規分布に拡張する必要がある。

10437 多変量正規分布が必要とするパラメータはベクトルと行列であることから、Stan におけるベクトル型、行列
10438 型など特殊な型についての理解を進める。

10439 コマ主題細目

10440 **対応のある群** 対応のあるデータというのは、反復測定あるいは測定間に相関関係が想定されるデータで
10441 として考えることができる。データ間に相関関係がある場合、多次元正規分布から出てくるモデルに
10442 なることから、これを使ったデータの書き方を習得する。まずは多次元正規分布の考え方とその表現
10443 方法を理解する。

10444 **ID を持ったデータ構造** 既に識別変数をもった整然データのモデル化については習得している。コード化
10445 するときは群の識別 ID ではなく個人の識別 ID を持ったコードを作成することになる。代入された変
10446 数がさらに代入されるという入子構造のプログラミングに慣れる。

10447 **個人差と変化量を想定した書き方** モデルを 3 群以上の比較にすることを考えると、多次元モデルをさら
10448 に広げることになるが、別の表現の仕方としてベースラインとそこからの変化量として表現できること
10449 になる。この時、ベースラインの散らばりすなわち個人差は、別の分布から出てきていることになる。この
10450 表記方法は今後の階層モデルにつながる観点でありことに言及する。

10451 キーワード

- 10452 • Within 計画
- 10453 • 分散共分散行列
- 10454 • ベクトル型と行列型
- 10455 • 多次元正規分布

10456 F.22.2 授業情報

10457 ■コマの展開方法 講義/遠隔/演習

10458 予習・復習課題

10459 ■予習 対応のある t 検定や Within モデルなど、伝統的な方法による分析方法について復習しておく。

10460 ■復習 パラメータリカバリや身の回りのデータを使って、今回のモデルが具体的にどのように使うことが
10461 できるかを考えておく。

10462 F.23 モデリングの目から見た検定 5 ; カテゴリカル分布をつかって

10463 F.23.1 授業内容

10464 科目の中でのこのコマの位置づけ

10465 コマ主題細目

10466 これまでは連続変数についてのモデルばかり扱ってきたが、度数など心理学で扱うデータの中にはカテゴリカルなものも少なくない。これらを帰無仮説検定の文脈で扱う時は、 χ^2 検定が有用であるが、カテゴリカルな分布を使ったベイジアンアプローチももちろん、結果の解釈には有用かつ直感的である。

10469 ここでは連続分布と離散分布の違いを確認し、いくつかの代表的な離散分布を演習とともに学び、カテゴリカルな出力変数の分析にも確率モデルが有用であることを確認する。

10471 **離散的な分布** 変数には離散・連続の違いがあり、確率分布でも確率質量と確率密度の違いがある。ここでは代表的な離散確率変数であるベルヌーイ分布、二項分布、多項分布を導入する。

10473 → 松浦 (2016) の Pp.82, 83, 85–89.

10474 χ^2 検定 カテゴリカルな変数の検定についてはこれまで扱って来なかったため、ここで改めて χ^2 分布を使った検定の例を導入する。 χ^2 検定は比率 (割合)、独立性、関連の強さ、適合度の検定などに用いられることを確認する。帰無仮説のおき方に注意すること、この検定がモデル適合度などの文脈でも望まれることに言及する。

10478 → 山内 (2010) の Pp.189-197.

10479 **カテゴリカル分布のモデリング** 確率分布を用いたアプローチをすることで、母比率を直接検証したり、連言命題が成立する確率など生成量を使ってさまざまな検証ができることを確認する。

10481 → 豊田 (2016) の Pp.136–163.

10482 κ 係数の算出 変数間の関連の強さを見る指標として κ 係数がある。一致率の係数ともして知られており、記述統計的アプローチでも算出できるものではあるが、確率モデルで表現する場合は工夫が必要である。

10485 → Lee and Wagenmakers (2013) の Pp.56–59.

10486 キーワード

- 10487 • 離散分布と連続分布
- 10488 • ベルヌーイ分布
- 10489 • 二項分布
- 10490 • 多項分布
- 10491 • クロス集計表
- 10492 • κ 係数

10493 F.23.2 授業情報

10494 ■コマの展開方法 講義/遠隔/演習

10495 予習・復習課題

10496 ■予習 確率変数が連続的か、離散的かということがどういう違いであるのかについて、データ解析基礎の
10497 確率に関する資料などを参考に確認しておく。また尺度水準による数値データの分類についても見直してお
10498 くと良い。

10499 ■復習 カテゴリカルな分類、集計に関しては身の回りに多くのデータ例がある。たとえば官公庁の統計資料
10500 などをもとに、母比率の推定や連言命題が成立する確率など、さまざまな仮説を自ら立てて検証すると良い。

10501 F.24 一般化線形モデル

10502 F.24.1 授業内容

10503 科目の中でのこのコマの位置づけ

10504 コマ主題細目

10505 ここまで要因計画が線形モデルと同一であること、すなわち一般線形モデルについて議論されてきた。あら
10506 ためて回帰分析の確率モデルを考えると、平均に構造を入れたモデルという意味で同じであることが確認で
10507 ける。ここで確率分布を違う形に変えることでより一般的な線形モデル、一般化線形モデルに拡張すること
10508 ができる。確率モデルによっては結果変数の型の違いによって確率分布が変わり、確率分布によってはパラメー
10509 タの取りうる範囲が定まるのでリンク関数によって変換する必要があること、結果を解釈するときはリンク関
10510 数を經由して分析されていることなどに注意が必要である。まずはパラメータの数が少ないロジスティック回
10511 帰について学ぶ。

10512 一般線形モデル 正規分布の平均構造を導入するという意味で、回帰分析のベイズ推定はこれまでと同様
10513 に実施、解釈することができる。事後分布や事後予測分布などを使って、最尤推定のモデルと異なる
10514 点を確認しておく。

10515 データに合わせた確率分布 データの形によっては確率分布の形を変える必要がある。まずはベルヌーイ
10516 分布によるロジスティック回帰分析を通じて、確率分布関数を選択できること、そのためにパラメータ
10517 の形をリンク関数を經由し得て変換することを学ぶ。

10518 リンク関数とパラメータの解釈 リンク関数を介して線形モデルを考えるので、独立変数が一単位増える
10519 ことがそのまま従属変数が一単位増えることにはつながらない。これらの関係を知るために、リンク関
10520 数、逆リンク関数の関係を辿って考える。ロジスティック回帰の場合は、オッズ、オッズ比などの用語に
10521 も触れることになる。またリンク関数による一般化が理解できれば、同様の考え方で他のさまざまな離
10522 散確率分布に応用できることが用意に想像できるだろう。

10523 キーワード

- 10524 • 一般化線形モデル
- 10525 • ベルヌーイ分布
- 10526 • ロジスティック回帰分析
- 10527 • リンク関数
- 10528 • オッズ比

10529 F.24.2 授業情報

10530 ■コマの展開方法 講義/遠隔/演習

10531 予習・復習課題

10532 ■予習 確率変数が連続的か、離散的かということがどういう違いであるのかについて、データ解析基礎の
10533 確率に関する資料などを参考に確認しておく。また尺度水準による数値データの分類についても見直してお
10534 くと良い。

10535 ■復習 離散的なデータについて、身近な例を考えてみると良い。また今回導入した分布関数以外にも応用
10536 を考えることができるから、確率分布とリンク関数の一覧を参考にさまざまなモデルに思いを馳せてみると
10537 良い。

10538 F.25 階層線形モデル

10539 F.25.1 授業内容

10540 科目の中でのこのコマの位置づけ

10541 ここまで線形モデル、一般化線形モデルの例を見てきたが、久保 (2012) の例にならって一般化線形混合
10542 モデル、階層ベイズモデルへとモデルを展開させていく。

10543 一般化線形混合モデルについては、Within デザインですすでに対応しており、正規分布以外の確率分布を
10544 使うことで一般化可能である。ここに切片や傾きなど、係数の方に分布が混ぜ合わせられることで階層化され
10545 たモデルとなる。ネストされたデータの具体例として、反復測定と大規模調査の二種類を取り上げ、また階層
10546 モデルの設計図を書いてから分析コードを書く手順を確認する。

10547 コマ主題細目

10548 一般化線形混合モデル Within デザインの分析モデルを再確認するところから考える。このモデルは個体
10549 差のような個人ごとに変わる要因が含まれており、ここで変量効果と固定効果の違いを考える。また
10550 従属変数が正規分布でないモデルにすることで、一般化線形モデルと考えることができる。

10551 ネストされたデータ すでに反復測定データの場合が該当するが、データが階層性を持っているネストさ
10552 れたデータの例として、プロ野球データや大規模調査の例を考える。これらに対して、階層化しない分
10553 析とする分析とで解釈が異なる例を挙げ、階層的なデータに対して適切な分析が必要であることを理
10554 解する。

10555 階層線形モデル 階層化されたモデルを数式的に理解する。名称として、レベル 1/2 の効果、個人/集団レ
10556 ベルの効果と呼ばれることもあることを確認する。またモデルの設計図をかき、それをコードに起こす
10557 ことで分析ができることを確認する。個別の回帰直線を引く場合に比べて、縮小が起こっていることを
10558 モデル比較を通じて確認する。

10559 キーワード

- 10560 ● 混合モデル
- 10561 ● ネストされたデータ
- 10562 ● 固定効果

- 10563 • 変量効果
- 10564 • 階層モデル

10565 F.25.2 授業情報

10566 ■コマの展開方法 講義/遠隔可/演習

10567 予習・復習課題

10568 ■予習 反復測定デザインの分析例を確認しておく。そこで分布が混合されていることを確認しておく。

10569 ■復習 階層線形モデルが応用できるようなデータ例を身の回りから考えてみるとよい。その上で、データが
10570 どのような背景から生成されているかの設計図を書き、設計図からコードに起こすという手順を一步步確
10571 認しておこう。

10572 F.26 混合分布モデル

10573 F.26.1 授業内容

10574 科目の中でのこのコマの位置づけ

10575 ここまで GLMM, HLM とさまざまな分布を組み合わせて利用するモデルについてみてきた。

10576 今回は混合正規分布モデルと 0 過剰ポアソンモデルを考える。これまではデータが全体的に均質的である
10577 ことを想定していたが、そもそも異なる種類のデータが混じり合っていると考えられる場合、異なる分布をあ
10578 てがう方がよい。ここでどちらの分布に属するかがある種の確率で代わり、それに従って続くモデルが異なる
10579 という、離散確率分布による条件分岐をデータ生成メカニズムに導入することになる。またこれを Stan で実
10580 行する場合には、離散変数が直接扱えないことから、すべての可能な場合を数え上げて足し合わせるという
10581 周辺化して消去するというテクニックを利用することになる。技術的に高度な側面もあるが、条件分岐をデー
10582 タ生成メカニズムに組み合わせることができれば表現力は一気に広がることになる。

10583 コマ主題細目

10584 **混合分布モデル** 混合分布モデルは確率的なクラスター分析 (Model Based Clustering) でもある。階層
10585 線形モデルと違って、クラスター分析はデータの背後にあるクラスが明示的に示されておらず、データ
10586 の適合から考えることになる。まずデータの可視化によってまずその可能性に気づき、これをどのよう
10587 にモデル化できるか、そのアイデアを理解する。具体的には、ベルヌーイ分布など離散変数のパラメー
10588 タによって条件分岐が発生し、各条件のもとで確率モデルが描かれることになるが、この確率モデルを
10589 どのように統合するかということを設計図の段階で理解する。設計図での理解を踏まえて、Stan での
10590 実装レベルでの理解に進むことができるのだから。

10591 → [松浦 \(2016\)](#) の Pp.209–213.

10592 **周辺化消去** Stan は離散確率分布を直接モデルの中に組み込むことはできない。そこで `log_sum_exp` と
10593 いう特殊な関数を使って、すべての場合わけを行った離散モデルの統合を考える。まずターゲット記法
10594 について理解し、続いて周辺化消去の書き方を理解する。また混合正規分布モデルの場合、ラベルス
10595 イッチングの問題が生じることが考えられるから、`ordered vector` の型を利用することが多い。こ
10596 うした特殊な関数やベクトルについても理解を深める

10597

→ 松浦 (2016) の Pp.203–208.

10598

10599

10600

10601

0 過剰ポアソン データの性質を考えると、必ずしも正規分布モデルばかりではなく、離散変数など一般的なモデルまで拡張することが可能である。そこで野球データなど具体的な分布情報とともに、ゼロ過剰ポアソン分布を使ったモデルを考える。データ生成のメカニズムがより具体的、実践的に考えることができるので、モデリングによる分析の自由度が高まることを実感できるだろう。

10602

→ 松浦 (2016) の Pp.82, 83, 85–89.

10603

キーワード

10604

- クラスタ分析

10605

- 混合分布モデル

10606

- 周辺化消去

10607

- ゼロ過剰ポアソン分布

10608

F.26.2 授業情報

10609

■コマの展開方法 講義/遠隔可/演習

10610

予習・復習課題

10611

■予習 データの描画から気づかされることは無数にある。探索的にデータプロットができるように、ggplot やデータハンドリング技術について復習しておくことが望ましい。

10612

10613

■復習 混合分布モデルが応用できるようなデータ例を身の回りから考えてみるとよい。その上で、データがどのような背景から生成されているかの設計図を書き、設計図からコードに起こすという手順を一步ずつ確認しておこう。

10614

10615

10616

F.27 確率的プログラミングの応用 1; 項目反応理論

10617

F.27.1 授業内容

10618

科目の中でのこのコマの位置づけ

10619

10620

10621

10622

10623

10624

10625

ここまでで線一般形モデル、一般化線形モデル、階層線形モデル、混合分布モデルと定型的な分析モデルについて一通り学んできた。以後はアラカルト的に、確率的プログラミングの応用による柔軟なモデリング例のトピックスを取り上げる。最初に扱うのは、既に前期に学んだ項目反応理論のモデリングである。尺度作成の文脈で、理論的概要は一通り説明が終わっているところであるが、改めて確認するとともに確率的モデリングとして実装する。確率モデルとして考えると、0/1 の反応に対するロジスティック回帰の応用であり、実装自体は既有知識の応用で可能であろう。コーディングのポイントとして、long 型データ (tidy data) にしておくことで欠損値が含まれる場合も対応できるようになることが挙げられる。

10626

コマ主題細目

10627

ロジスティックモデルの復習 本講では前期のうちに、ロジスティックモデルについての理論的説明と、R の ltm パッケージによる実践例を解説済みである。とはいえ、以前学んでから時間が空いているの

10628

10629 で、あるいは後期のみ履修する学生もいることが考えられるので、授業の冒頭で 15 分程度の時間
10630 をかけて、大まかな理論・モデルの復習をしておく必要があるだろう。

10631 **ロジスティック回帰モデルでの実装** ロジスティック回帰分析を思い出しつつ、1PL,2PL,3PL モデルそ
10632 れぞれを transformed parameters ブロックで記述することを演習で学ぶ。

10633 **整然データでの分析** データを整然データの形にして分析することで、欠損値が含まれないデータセットを
10634 作って分析に応用することができる。ここでは個人と項目それぞれを識別する変数が必要になるが、こ
10635 れまで学んできた技術で十分対応可能であると考えられる。

10636 キーワード

- 10637 • 項目反応理論
- 10638 • 1PL ロジスティックモデル
- 10639 • 2PL ロジスティックモデル
- 10640 • 3PL ロジスティックモデル
- 10641 • 整然データ

10642 F.27.2 授業情報

10643 ■コマの展開方法 講義/遠隔可/演習

10644 予習・復習課題

10645 ■予習 これまでの知識や技術を組み合わせて問題に対応することになる。項目反応理論とロジスティック
10646 モデル、GLM におけるロジスティック回帰分析、データハンドリングにおける整然データの考え方など、これ
10647 までの資料に戻って復習しておくとうい。

10648 ■復習 自分で描いたモデルが R のパッケージが出す答えとどの程度一致するのかを確認しておこう。また
10649 欠損値がある被験者の被験者母数は、その確信区間が広くなると考えられる。なぜそうなるかを改めて考え、
10650 実際のデータ適用例で確認しておこう。

10651 F.28 確率的プログラミングの応用 2; 変化点と折線回帰

10652 F.28.1 授業内容

10653 科目の中でのこのコマの位置づけ

10654 コマ主題細目

10655 変化点検出は、時系列的なデータの中に異なる二つの平均値を持つ群があることをモデリングする手法で
10656 ある。とくにある時点から異なる群に属する、という系列的な意味があることと、変化点があるとすればどのあ
10657 たりになるかという「変化点の位置的不明確さ」を確率分布で表現し、データから検出するという観点では、確
10658 率モデルの表現の自由さとデータとの接合を許す確率的プログラミング言語の面白さを味わうには最良の材
10659 料である。

10660 まずは混合分布モデルのように、二つの群を分類するモデルを再確認し、その上で時系列的なデータとい
10661 う既有知識から「変化点」という考え方の導入、モデリングへと繋げる。またデータによっては、一定の点を期
10662 に線形モデルの傾きが変わるような表現が可能なのがある。この変化点と回帰分析を融合させた、折線回

10663 帰モデルを考えることで、固定的なモデルを超えた柔軟なモデリングが可能であることを理解する。
 10664 ただしここで使うデータは時系列的なものであるから、一般的な回帰分析の前提であるサンプルの独立性
 10665 がない。その意味で不適切なモデルであることに注意し、次の時系列分析へと繋げる。

10666 **混合分布モデル** データは可視化することが重要であり、見れば明らかに異なる状態の混合であることがわ
 10667 かる場合がある。具体例とともに可視化を行い、またこれまで学んだ混合分布モデルで表現できるこ
 10668 とを再確認する。ここで用いるデータは、小杉の体重記録データを用いる。

10669 **変化点検出** データの横軸が時系列的な意味を持つのであれば、時空を超えて二つの群が混合している
 10670 というのは不自然な前提である。そこで横軸に時系列的な意味を置くと、ある時点から状態が変化した
 10671 ものとして考えることができる。ここでその時点が「いつ」であるのかは不明であるが、わからないこと
 10672 を確率で表現するのが確率モデルのおもしろい点である。変化点を確率的パラメータとし、その前後
 10673 で群が異なるというモデルは、変化点検出のモデリングと言われる。このモデリングはポリグラフ検査
 10674 など、実践的な場面での利用価値も高い。

10675 →Lee and Wagenmakers (2013) の Pp.59–61, 松浦 (2016) の Pp.238-245

10676 **折線回帰** 平均点の位置が変わるだけでなく、変化の傾向が明らか場合は線形モデルを当てはめること
 10677 ができる。変化点の前後で傾きが変わるような線形モデルは、折線回帰とも呼ばれる。折線回帰モデ
 10678 ルの実装については、変化点と回帰モデルを組み合わせたとで、折れる点を繋げる数学的補正を加
 10679 えたモデルへと修正する。最後に、説明変数が時点であることから回帰分析の前提として標本の独立
 10680 性が担保されていない問題を指摘する。

10681 キーワード

- 10682 • 混合分布モデル
- 10683 • 変化点
- 10684 • 折線回帰
- 10685 • 時系列分析

10686 F.28.2 授業情報

10687 ■コマの展開方法 講義/遠隔可/演習

10688 予習・復習課題

10689 ■予習 混合分布モデルの応用になるので、混合分布モデルの基本的な書き方や解析方法について、第
 10690 F.26 講を復習しておくことが望ましい。

10691 ■復習 折線モデルが応用できそうなデータを見つけて、自分なりに実践してみると理解が深まるだろう。と
 10692 くに折れる点が多数あるモデルや、折れる点の数を検出するモデルへと拡張するなど、モデル展開の可能性
 10693 をかんがえることもできる。

F.29 確率的プログラミングの応用 3; 状態空間モデル

F.29.1 授業内容

科目の中でのこのコマの位置づけ

前時に時系列的なデータを導入したが、時系列的な性質を無視したモデリングになっていた。時系列的な分析方法は、心理学においてもウェアラブル端末の利用や SNS など公的なデータを分析することなどにも利用できるため、非常に有用なものになりうる。しかしデータの特徴として非独立性の問題、周期性やトレンドの存在などがあり、周波数解析をおこなったり多次元の行列分解などが必要である。中でも状態空間モデルは比較的シンプルであり、とくにベイズアンプローチで実装が容易になったと言えるだろう。

ここでは状態空間モデルの基本的な考え方を導入し、モデリングについて解説する。ここで状態と観測の分離を行い、とくに観測が行われていない点があっても分析できること、観測が行われていない点をパラメータとして保管することに言及する。観測が行われていない点が保管できるのであれば、未来の時点についても予測が可能になるということである。ホワイトノイズモデルでそれを行うと、確信区間が広がっていくことが観測される。そこでトレンドを入れたモデルにすることで、さらに予測の形を変えられることを学ぶ。

そのほかにも季節項など、時系列特有の情報を組み込むことや、二次元に展開することで空間データの分析にも応用できることに言及する。

コマ主題細目

時系列データの特徴 時系列的なデータがどのように得られ、どのようなシーンで利用可能であるかを概観する。ここで時系列データはサンプルの独立性が満たされていないという問題があるため単純な線形回帰は不適切であること、また周期性やトレンド、介入効果が出てくるまでの期間など独自に考えなければならぬことがいろいろ含まれている。これまで研究されてきた領域や研究方法について概観する。

状態空間モデル さまざまな分析方法がこれまで考えてこられているが、状態空間モデルはその中でも比較的簡単な数理的構造を持ち、またベイズアンモデリングを利用することでかつての分析モデルが必要としたスムージングなどを、特段意識することなく分析できる。状態と観測というモデルの基本構造を提示し、これらがどのように実装可能かをみる。

→ [松浦 \(2016\)](#) の Pp.229-235, [馬場 \(2019\)](#) の第 5 部

欠損値の補間 観測時点には欠損が含まれることもあり、これをパラメータとして推定・補間することができる。またこれが可能であるということは、未来の時点を欠損値として考えれば予測ができることにもなる。プログラミングの工夫により、欠損を補間するようなコードの書き方を学ぶ。また単純なホワイトノイズモデルであればあまり予測として意味がないが、トレンドを考えることで時系列的な影響について考えることができる。

状態空間モデルの展開 状態空間モデルは、説明変数を加えた回帰モデルに応用したり、周期性をモデリングすることなども可能である。さらに時系列は一次元的であるが、二次元にも広げると空間分析にも利用が可能であることに言及する。これからの心理学は、時系列や空間など状況変数をより積極的に取り組んだモデルも利用するようになるだろう。

10729 キーワード

- 10730 • 状態空間モデル
- 10731 • トレンド
- 10732 • 補間

10733 F.29.2 授業情報

10734 ■コマの展開方法 講義/遠隔可/演習

10735 予習・復習課題

10736 ■予習 時系列データを、時間を独立変数とした回帰分析にすることでどういった問題があるのかについて、回帰分析や確率モデルの前提などを考えて振り返っておくことが良い予習になるだろう。

10738 ■復習 身の回りの身近ところからでもデータを取ることができるのが、時系列データのおもしろいところでもあるので、応用可能なデータを探して分析してみると良い。可能であれば今日からでも、時系列的なデータを取り始めると、長期的に見て非常に興味深い分析ができるようになるだろう。

10741 F.30 モデル比較

10742 F.30.1 授業内容

10743 科目の中でのこのコマの位置づけ

10744 最終回となるこの回では、これまで後期の授業で扱ってきたさまざまなモデルについて総括し、ベイジアンモデリングの心理学的位置付けについて解説する。加えて最後の話題提供として、モデル比較について言及する。帰無仮説検定の代わりとして考えるのであれば、区間推定を使った比較が必要であるし、ベイジアンモデリングの観点からはモデル比較になる。ベイズファクターによるモデル比較とそれを実行するためのブリッジサンプリング法、また WAIC など予測的観点から評価する方法があることなどに言及する。

10749 コマ主題細目

10750 **ベイジアンモデリング** 一般的に「モデリング」という観点から、これまでの授業内容だけでなく心理学における研究法としてのその意味や意義を考える。帰無仮説検定の Alternative として利用するだけでなく、心理学的メカニズムをより具体的に、緻密に記載するために数学的方法を用いることで、心のメカニズムの理解を深めることができるかもしれない。そこに含まれる仮定や前提について、自覚的に記述する必要があることがモデリングの利点であり、MCMC をはじめとするベイズ統計の技術は、それを可能にしてくれる方法論的補助にすぎない。

10756 **帰無仮説検定の代案** 帰無仮説検定の代わりにモデリングアプローチを取ることの利点はさまざま挙げられるが、 p 値のように「ここだけ見ておけば良い」というような機械的判断ができる基準がない。むしろそうした機械的判断の弊害が指摘されてきているのであるが、代案としてはどのような基準があるのかはドメイン知識に基づく必要がある。むしろパラメータだけでなくデータの観点から考察できるようになったことや、実質的な値に基づいて考えられることを利点と捉えつつ、判断基準としての ROPE などについて一瞥する。

10762 パラメータ推定かモデル比較かについては、→[Kruschke \(2014\)](#) の Pp.341–361.

10763 **モデル比較** モデルとその有用性を考えるにあたって、パラメータ推定かモデル比較かという二つのアプ
10764 ローチがあり得る。後者については、階層モデルによるもの、ベイズファクター、WAIC が考えられる。
10765 ただし WAIC については、[渡辺ベイズ理論](#)ともいうべき、より包括的なベイズ理論の枠組みで捉え直
10766 す必要があり、心理学研究にこの枠組みがどれほど有用かについては、いまだに結論が出ていない
10767 ところでもある。ここでは概略的にその特徴に触れるにとどまり、受講生諸君の今後の活躍に期待し
10768 たい。

10769 WAIC については → [浜田他 \(2019\)](#) が丁寧である。

10770 キーワード

- 10771 • ROPE
- 10772 • ベイズファクター
- 10773 • サヴェージ・ディッキー法
- 10774 • ブリッジサンプリング

10775 F.30.2 授業情報

■コマの展開方法 講義/遠隔可/演習

参考文献

- Abelson, R. P. (1954). A technique and a model for multi-dimensional attitude scaling. *Public Opinion Quarterly*, 18(4), 405–418.
- 足立 浩平 (2006). 多変量データ解析法—心理・教育・社会系のための入門 ナカニシヤ出版, 単行本版, 171
- 岡太 彬訓 (2008). データ分析のための線形代数 共立出版, 単行本版
- Allport, G. W. (1967). Attitudes. In Fishbein, M.(Ed.) *Readings in Attitude Theory and Measurement*(pp. 3–13). New York: John Wiley & Sons Inc
- Amrhein, V., Greenland, S., & McShane, B. (2019). Scientists rise up against statistical significance.
- 馬場 真哉 (2019). 実践 Data Science シリーズ RとStan ではじめる ベイズ統計モデリングによるデータ分析入門 (KS 情報科学専門書) 講談社
- Chalmers, R. P. (2012). mirt: A Multidimensional Item Response Theory Package for the R Environment. *Journal of Statistical Software*, 48(6), 1–29, DOI: <http://dx.doi.org/10.18637/jss.v048.i06>.
- Epskamp, S. (2021). semPlot. <https://github.com/SachaEpskamp/semPlot>.
- 藤原 武弘 (2001). 社会的態度の理論・測定・応用 関西学院大学出版会
- Gelman, A. et al. (2006). Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper). *Bayesian analysis*, 1(3), 515–534.
- Grimm, L. G., & Yarnold, P. R. (1994). *Reading and Understanding Multivariate Statistics*: American Psychological Association. (グリム, L.G.・ヤーノルド, P.R. 小杉 考司・高田 菜美・山根 嵩史 (訳)(2016). 研究論文を読み解くための多変量解析入門 基礎篇: 重回帰分析からメタ分析まで 北大路書房), URL: <http://amazon.co.jp/o/ASIN/4762829404/>
- Grimm, L. G., & Yarnold, P. R. (2001). *Reading and Understanding More Multivariate Statistics*: American Psychological Association. (グリム, L.G.・ヤーノルド, P.R. 小杉 考司・高田 菜美・山根 嵩史 (訳)(2016). 研究論文を読み解くための多変量解析入門 応用篇: SEM から生存分析まで 北大路書房), URL: <http://amazon.co.jp/o/ASIN/4762829439/>
- Gronau, Q. F., Singmann, H., & Wagenmakers, E.-J. (2020). bridgesampling: An R Package for Estimating Normalizing Constants. *Journal of Statistical Software*, 92(10), 1–29, DOI: <http://dx.doi.org/10.18637/jss.v092.i10>.
- Guilford, J. (1954). *Psychometric Methods*. New York: McGraw-Hill Book Company. (ギルフォード, J.P 秋重 善治(訳)(1959). 精神測定法 培風館)
- 南風原 朝和・芝 祐順 (1987). 相関係数および平均値差の解釈のための確率的 な指標 教育心理学研究, 35(3), 259–265, DOI: http://dx.doi.org/10.5926/jjep1953.35.3_259.

- 浜田 宏・石田 淳・清水 裕士 (2019). 社会科学のためのベイズ統計モデリング 朝倉書店, URL: <http://amazon.co.jp/o/ASIN/4254128428/>
- 豊田 秀樹 (2012). 項目反応理論 [入門編] (第2版) (統計ライブラリー) 朝倉書店, 第単行本版
- 樋口 耕一 (2020). 社会調査のための計量テキスト分析—内容分析の継承と発展を目指して【第2版】 KH Coder オフィシャルブック ナカニシヤ出版, 第単行本版, 264, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4779514746>
- 平岡 和幸・堀 玄 (2004). プログラミングのための線形代数 オーム社, 第単行本版, 355, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4274065782>
- 清水 裕士 (2018). 阪神ファン—巨人ファンの2大精力構造は本当か 豊田秀樹 (編) たのしいベイズモデリング (pp. 21–32) 北大路書房
- Holzinger, K. J., & Swineford, F. (1939). A study in factor analysis: The stability of a bi-factor solution.. *Supplementary educational monographs*.
- 池田 功毅・平石 界 (2016). 心理学における再現可能性危機:問題の構造と解決策 心理学評論, 59(1), 3-14, DOI: http://dx.doi.org/10.24602/sjpr.59.1_3.
- JASP Team (2021). JASP (Version 0.16)[Computer software]. URL: <https://jasp-stats.org/>.
- J.Dobson, A. (2008). 一般化線形モデル入門 原著第2版, 田中 豊・森川 敏彦・山中 竹春・富田 誠 (訳) 共立出版, 第単行本版, 280, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4320018672>
- 川端 一光・荘島 宏二郎 (2014). 心理学のための統計学入門 [心理学のための統計学 1]: ココロのデータ分析 誠信書房, URL: <http://amazon.co.jp/o/ASIN/4414301874/>
- 加藤 健太郎・山田 剛史・川端 一光 (2014). Rによる項目反応理論 オーム社, 第 Kindle 版版
- 小杉 考司 (2018). 言葉と数式で理解する多変量解析入門 北大路書房, URL: <http://ci.nii.ac.jp/ncid/BB27527420>
- 小杉 考司 (2019a). Rでらくらく心理統計: RStudio 徹底活用 講談社, URL: <http://ci.nii.ac.jp/ncid/BB27718917>
- 小杉 考司 (2019b). その他の他変量解析 楠見 孝・日本心理学会 (編) 公認心理師の基礎と実践 5 心理学統計法 (pp. 189–206) 遠見書房
- 小杉 考司・清水 裕士 (編) (2014). M-plusとRによる構造方程式モデリング入門 北大路書房, 第単行本版, 332, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4762828254>
- 小杉 考司 (2014). 学校適応感尺度 FIT の開発 研究論叢. 第3部, 芸術・体育・教育・心理, 64, 69-82, URL: <https://ci.nii.ac.jp/naid/120005596041/>.
- Kruschke, J. (2014). *Doing Bayesian data analysis 2nd Ed.* NewYork: Elsevier2nd ed. (クルシュケ, J.K 前田 和寛・小杉 考司 (監訳)(2017). ベイズ統計モデリング: R, JAGS, Stan によるチュートリアル 原著第2版 共立出版)
- Kruskal, B., Joseph (1964a). Multidimensional scaling by optimizing goodness of fit to a non-metric hypothesis. *Psychometrika*, 29(1), 1–27.
- Kruskal, B., Joseph (1964b). Nonmetric multidimensional scaling: a numerical method. *Psychometrika*, 29(2), 115–129.
- 久保 拓弥 (2012). データ解析のための統計モデリング入門 – 一般化線形モデル・階層ベイズモデル・MCMC (確率と情報の科学) 岩波書店, 第単行本版, 272
- Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian Cognitive Modeling: A Practical Course*: Cambridge University Press. (マイケル・D. リー・エリック・ジャン ワーゲンメイカーズ井関

- 龍太 (訳)(2017). ベイズ統計で実践モデリング: 認知モデルのトレーニング 北大路書房), URL: <http://amazon.co.jp/o/ASIN/4762829978/>
- 松浦 健太郎 (2016). Stan と R でベイズ統計モデリング (Wonderful R) 共立出版, URL: <http://amazon.co.jp/o/ASIN/4320112423/>
- 三中 信宏 (2018). 統計思考の世界 ~曼荼羅で読み解くデータ解析の基礎 技術評論社
- 宮川 雅巳 (1997). グラフィカルモデリング (統計ライブラリー) 朝倉書店
- 宮谷 真人・坂田 省吾・林 光緒・坂田 桐子・入戸野 宏・森田 愛子 (編) (2009). 心理学基礎実習マニュアル 北大路書房, 単行本(ソフトカバー)版
- 村上 正康・佐藤 恒雄・野澤 宗平・稲葉 尚志 (2016). 教養の線形代数 培風館, 単行本版
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *ETS Research Report Series*, 1992(1), i-30.
- 長沼 伸一郎 (2011). 物理数学の直観的方法〈普及版〉 講談社, 第 Kindle 版版, 301, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=B00JQYYCPA>
- 西村 武 (1977). 主観評価の理論と実際 テレビジョン, 31(5), 369-377, URL: <https://cir.nii.ac.jp/crid/1390001205397311616>, DOI: http://dx.doi.org/10.3169/itej1954.31.5_369.
- 西里 静彦 (2010). 行動科学のためのデータ解析—情報把握に適した方法の利用 培風館
- Norretranders, T. (2002). ユーザーイリュージョン—意識という幻想, 柴田 裕之 (訳) 紀伊國屋書店, (トール ノーレットランダーシュ)
- 岡太 彬訓 (2008). データ分析のための線形代数 共立出版, URL: <http://amazon.co.jp/o/ASIN/4320018591/>
- 岡太 彬訓・今泉 忠 (1994). パソコン多次元尺度構成法 共立出版, 単行本版, 174, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4320014723>
- 小野島 昂洋 (2021). lav2tikz.R, <https://github.com/onoshima/myfunction>.
- 小塩 真司 (2020). 性格とは何か より良く生きるための心理学 (中公新書) 中央公論新社, 第 Kindle 版版
- Partchev, I., Partchev, M. I., & Suggests, M. (2017). Package ‘irtoys’. *A collection of functions related to item response theory (IRT)*.
- Revelle, W. (2021). *psych: Procedures for Psychological, Psychometric, and Personality Research*, Northwestern University. Evanston, Illinois, URL: <https://CRAN.R-project.org/package=psych>, R package version 2.1.3.
- Rizopoulos, D. (2006). ltm: An R package for Latent Variable Modelling and Item Response Theory Analyses. *Journal of Statistical Software*, 17(5), 1-25, URL: <http://www.jstatsoft.org/v17/i05/>.
- Samejima, F. (1997). Graded response model. In *Handbook of modern item response theory* (pp. 85-100): Springer
- 芝 祐順 (1979). 因子分析法 東京大学出版会, 単行本版
- 清水 裕士 (2016). フリーの統計分析ソフト HAD: 機能の紹介と統計学習・教育, 研究実践における利用方法の提案 メディア・情報・コミュニケーション研究 (1), 59-73, URL: <https://ci.nii.ac.jp/naid/120005744983/>.
- 清水 裕士・荘島 宏二郎 (2017). 社会心理学のための統計学 [心理学のための統計学 3]: 心理尺度の構成と分析 誠信書房, URL: <http://amazon.co.jp/o/ASIN/4414301890/>
- 清水 裕士 (2021). 心理学統計法 (放送大学教材 1638) 放送大学教育振興会

- Stevens, S. S. (1946). On the theory of scales of measurement. *Science*, 103(2684), 677-680.
- 末永 俊郎 (編) (1987). 社会心理学研究入門 東京大学出版会, 第ハードカバー版
- 高橋 正視 (2002). 項目反応理論入門—新しい絶対評価 アイデア アイデア出版局, 第単行本版, 255, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4900561002>
- 高根 芳雄 (1980). 多次元尺度法 東京大学出版会, 第一版, 332, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=B000J8ABS0>
- 田中 良久 (1977). 心理学的測定法 東京大学出版会, 第単行本版
- 豊田 秀樹 (2000). 共分散構造分析 応用編—構造方程式モデリング (統計ライブラリー) 朝倉書店, 第単行本版, 303, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4254126611>
- 豊田 秀樹 (2007). 共分散構造分析—構造方程式モデリング 理論編 (統計ライブラリー) 朝倉書店, 第単行本版, 287, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4254126964>
- 豊田 秀樹 (2008). データマイニング入門 東京図書
- 豊田 秀樹 (2016). はじめての 統計データ分析 - ベイズ的(ポスト p 値時代)の統計学 - 朝倉書店, URL: <http://amazon.co.jp/o/ASIN/4254122144/>
- 豊田 秀樹 (2017). 実践 ベイズモデリング -解析技法と認知モデル- 朝倉書店
- 豊田 秀樹 (2018). たのしいベイズモデリング: 事例で拓く研究のフロンティア 北大路書房
- 豊田 秀樹 (2019). たのしいベイズモデリング 2: 事例で拓く研究のフロンティア 北大路書房
- 豊田 秀樹 (2020). 瀕死の統計学を救え! 朝倉書店
- Van Lissa, C. J. (2019). tidySEM: A tidy workflow for running, reporting, and plotting structural equation models in lavaan or Mplus.. <https://github.com/cjvanlissa/tidySEM/>.
- Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using leave-one-out cross-validation and WAIC. *Statistics and Computing*, 27, 1413-1432, DOI: <http://dx.doi.org/10.1007/s11222-016-9696-4>.
- Walter, I. (2012). スティーブ・ジョブズ 1,2, 井口 耕二 (訳) 講談社, 第新書版
- 山田 剛史・村井 潤一郎 (2004). よくわかる心理統計, やわらかアカデミズム・「わかる」シリーズ ミネルヴァ書房, URL: <http://ci.nii.ac.jp/ncid/BA68747748>
- 山内 光哉 (2010). 心理・教育のための統計法 サイエンス社, 第第 3 版, URL: <http://amazon.co.jp/o/ASIN/4781912354/>
- 永田 靖 (2005). 統計学のための数学入門 30 講 朝倉書店, 第単行本版
- シ (2016). 計算機言語のまとめノート 暗黒通信団, 第単行本版, 32, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4873100518>
- 千野 直仁・岡田 謙介・佐部利 真吾 (2012). 非対称 MDS の理論と応用 現代数学社, 第単行本版, 331, URL: <https://lead.to/amazon/jp/?op=bt&la=ja&key=4768704050>

索引

記号／数字

1 パラメータ・ロジスティックモデル	46, 114, 319
2 階差分のトレンド	354
2 パラメータ・ロジスティックモデル	46, 114, 319
3 パラメータ・ロジスティックモデル	114

A

AGFI	131
AIC	131
Amos	141

B

BIC	131
BUGS	186

C

CFI	131
Cohen の d	223
CRAN	369

E

EAP	219, 227, 228, 346
EAP 推定値	278, 301, 323

G

GFI	131
-----	-----

I

int 型	206
IRT	43
IT 相関	31, 44

J

JAGS	186
------	-----

K

K-平均法	305
ニット	368

M

MAP	228, 346
MAP 推定値	199, 282
matrix 型	206
MCMC	178, 283
MED	346
MED 推定値	278
Mplus	141

P

Preference Mapping	158
p 値	216

R

real 型	206
Rhat	233

RMSEA	131
Rhat	211

S

SRMR	131
------	-----

T

TLI	131
t 検定	126, 215, 233, 241, 287, 359

V

vector 型	206
----------	-----

W

Welch の補正	219
-----------	-----

あ

当て推量母数	114, 118, 320
α 係数	31, 113
閾上率	231
閾値	55, 120
位置母数	274, 290
一様分布	178
一般化線形混合モデル	283, 289
一般化線形モデル	18, 283, 289
一般線形モデル	74, 273, 283, 289
因果関係	123
因子間相関	85, 106
因子軸の回転	85, 103
因子的妥当性	32, 137
因子得点	37, 38, 46, 107, 122
因子負荷量	37, 38, 58, 61, 103, 107, 124, 136, 146
因子分析	15, 17, 32, 36, 81, 124
インタプリタ	187
隠匿情報検査	333
ウォード法	304
ウォームアップ	213
エディタ	187
エビデンス	177, 362
オッズ	285, 289
オッズ比	286
オブジェクト	169
重み	17

か

回帰分析	17, 73, 123, 177, 273
外生変数	139
階層性	295
階層線形モデル	261, 297
階層的クラスタリング	304
回転行列	85
カウント変数	290
確信区間	214, 228
確率質量	308
確率質量関数	263
確率的プログラミング言語	180, 185, 186, 205, 233, 359

確率分布	18, 161, 178, 179, 183, 204, 206, 207, 219, 228, 263, 308	項目プール	55
確率密度	308	固定効果	261, 297
確率モデル	176	固有値	79, 146
過剰識別	128	固有値分解	146, 151
型	206	固有ベクトル	79, 146
片側検定	224	固有方程式	82
カッパ係数	269	混合分布モデル	304
カテゴリ確率曲線	56	困難度	46, 56, 319
カテゴリカル因子分析	58, 121	困難度母数	114, 115, 117
カテゴリカル分布	306	コンパイラ	187
可読性	167	コンピュータ適応型テスト	50, 52
カルバック・ライブラー情報量	363		
間隔尺度水準	19, 28, 43, 121, 263	さ	
完全情報最尤推定	52	最高密度区間	200
観測度数	265	最小二乗法	99, 130, 177, 283
観測変数	123	採択	216
簡便的因子得点	108	最尤推定法	308
機械学習	17, 186, 304	最尤法	99, 103, 130, 175, 177, 277, 283, 320
棄却	216	作業フォルダ	390
危険率	216	残差	124
基準関連妥当性	32	シード値	180
期待値	182	ジオミン回転	106
期待度数	265	識別可能	128
基底	151	識別不可能	128
帰無仮説	216, 243	識別力	47, 56, 319
帰無仮説検定	218, 239, 358	識別力母数	114, 115, 117
帰無モデル	243	シグマ法	28, 30
逆行列	67, 73, 81, 91	事後確率最大値	213
逆リンク関数	289, 320	自己相関	341
客観性	175	事後分布	177, 229, 241
級内相関	298	事後予測分布	228, 229, 230
教師なし学習	304	事前・事後デザイン	251
共通因子	38	事前分布	177, 204
共通性	41	実験計画	175, 239
共通性の推定問題	99	実質的に等価な範囲	221, 361
共分散	22	質的変数	19
共分散構造分析	17, 123	弱情報事前分布	204
行ベクトル	61	尺度	25, 162
行列	62	尺度水準	18
行列式	82	尺度母数	274
虚偽検出	332	斜交回転	85, 105
距離	23, 151, 153, 304	主因子法	99
距離行列	151	重回帰分析	17, 73, 124
寄与率	100	修正指数	131, 138, 139
偶然誤差	36	収束的妥当性	32, 137
区間推定	214	自由度	130
クラスターリング	304	周辺化	179
グラフィカルモデリング	17	周辺化消去	307, 362
クロス集計表	263	周辺尤度	177, 178, 362
クロンバックのアルファ	31	縮小	299
群間計画	239, 251	順序ロジットモデル	126
群内計画	251	主成分分析	17, 99, 125
経験サンプリング	341	主成分法	99
係数	17	出版バイアス	141
形態素解析	149	順序尺度水準	19, 25, 58, 119, 125, 126, 143, 145, 155, 263
系統誤差	36	順序プロビットモデル	126
計量的多次元尺度構成法	155	状態空間モデル	342
系列範疇法	379	情報量規準	131, 363
検証的因子分析	32, 127, 137	信頼区間	214
効果量	223, 234, 361	信頼性	30, 36, 115
高級言語	186	数値モデリング	18
交互作用	289	数量化の理論	144
構成概念妥当性	32, 97	スカラー	62
構造方程式	126	スクリープロット	99, 102, 157
構造方程式モデリング	17, 123, 126, 265	スクリプト	187
項目情報曲線	53, 115, 116	スムージング	356
項目特性曲線	46, 115, 116	正規化定数	177, 178
項目反応カテゴリ特性曲線	56, 120	正規混合モデル	304
項目反応理論	30, 43		

正規表現	326		
正規分布	17, 36, 180, 303		
生成量	222, 257, 361		
整然データ	247, 258, 323		
正方行列	151		
正方対称行列	157, 252		
制約	128		
積率	18, 22		
絶対パス	390		
ゼロ過剰ポアソン	314		
線形代数	61		
線形モデル	126, 358		
宣言	206		
潜在変数	123		
尖度	21		
相関関係	123		
相関行列	63		
相関係数	23, 304		
相対パス	390		
双対尺度法	145		
測定方程式	126		
ソフトクラスタリング	305		
た			
ターゲット記法	308		
対応のある t 検定	251		
対応分析	145		
対角	63, 252		
対角行列	63		
対称行列	62, 151		
対数尤度	117, 210, 308		
態度	25		
タイプ 1 エラー	240, 360		
タイプ 2 エラー	237		
対立仮説	216		
多次元項目反応理論	59, 122		
多次元尺度構成法	151, 152		
多次元正規分布	251		
多次元展開法	162		
多重比較	240		
多段採点モデル	55		
妥当性	30, 32		
多変量解析	16		
多変量正規分布	251		
多母集団同時分析	137		
単位行列	63		
段階反応モデル	19, 30, 55, 119, 161		
探索的因子分析	127, 136		
単純構造の原則	41		
単純構造の原理	59, 137		
チャンク	366		
丁度識別	128		
直交回転	85, 104, 105		
通過率	44, 319		
データ生成モデリング	176, 203, 215, 223		
停止規則	361		
適合度	130, 265		
テキストマイニング	148		
デザイン行列	76, 273		
テスト情報関数	54		
テスト情報曲線	115, 116		
てっちゃんの手品	109		
テトラコリック相関係数	58		
点推定	278		
点双列相関係数	112		
転置	66, 90		
等価	52		
等現間隔法	26		
同時確率空間	210		
等裾区間	200		
特異値分解	146		
特異ベクトル	146		
独自性	41, 107		
独立性の検定	265		
トレース	80		
トレースプロット	211, 233		
な			
内生変数	139		
内的整合性信頼性	31		
内容的妥当性	32		
二項分布	181, 264, 287		
二次元正規分布	251		
日誌法	341		
ネイピア数	290		
ノルム	83		
は			
ハードクラスタリング	305		
ハイパーパラメータ	299		
配列	207		
パス解析	17, 135		
パス係数	124		
パス関	123		
パスダイアグラム	123, 126		
パラメータ	228		
パラメータリカバリ	234		
バリマックス回転	106		
半コーシー分布	204		
反復測定	251, 296		
判別分析	126		
非階層的クラスタリング	304		
非計量的多次元尺度構成法	154		
被験者母数	48		
標準得点	37		
標準偏差	21		
標本統計量	216		
比率尺度水準	19		
頻度主義統計学	287		
複雑性指標	107		
布置	157		
ブリッジ・サンプリング	363		
プログラミング言語	186		
プロマックス回転	106		
分割表	263		
分散	20		
分散共分散行列	63		
分散分析	18, 126, 239, 266, 359		
平均	20		
並行分析	100		
ベイジアンモデリング	357		
ベイズ推定	239, 283		
ベイズ推定法	308		
ベイズファクター	363		
ベイズ法	18, 175		
ベイズ統計学	287		
ベルヌーイ分布	181, 263, 279, 289		
変化点検出	332		
変数	16		
弁別的妥当性	32, 97, 137		
変量効果	261, 297		
ポアソン分布	287, 290		
補間	347		
ポリグラフ検査	332		
ポリコリック相関係数	58, 121		
ポリシリアル相関係数	58		
ホワイトノイズモデル	342		

ま	
マルコフ連鎖モンテカルロ法	178
マンハッタン距離	154
ミンコフスキー距離	154
名義尺度水準	19, 28, 126, 144, 151, 263, 273, 279
モーメント	22
モーメント法	175, 214, 224, 231, 287
モデリング	299, 319, 342, 357, 358, 359, 364
モデル比較	140, 361
や	
ユークリッド距離	304
有意差	359
有意水準	216
優越率	231
有効サンプルサイズ	211, 233
尤度	49, 177, 308
要因計画	239
ら	
ライフログ	341
ラッシュモデル	114
ラベルスイッチング	311
ランダム化比較実験	243
離散確率分布	279
リッカート尺度	119
量的変数	19
リンク関数	283, 289, 320
類似度	304
例数設計	234, 237
列ベクトル	61
連言命題	267
ロジスティック回帰	282, 320
ロジスティック関数	45, 284, 289, 319
ロジット関数	285
論理演算	173
わ	
歪度	21
漢字	
中央値絶対偏差	196