

日本語版への序文

私が本書を執筆しようと思った理由は、次世代の学生そして若手研究者のみなさんに、政治・経済・教育・公衆衛生などの多岐にわたる重要な社会問題を解決するにあたって、データ分析がいかに強力な道具であるか、ということを実感してもらいたかったからである。この本をきっかけにして、1人でも多くの優秀な若者が、計量社会科学(Quantitative Social Science)をさらに勉強し、将来様々な分野でデータ分析を社会のために活用してくれれば、本望である。

本書はプリンストン大学で私が教えてきた授業の題材を基にして書かれたものである。この授業は、主に社会科学に興味のある学部生を対象にした、データ分析入門を目的としている。したがって、本書は既存の教科書とは異なり、最初から実際に出版された論文で使われたデータを統計ソフトウェアを使って直接分析することによって、読者に計量社会科学の魅力を理解してもらうことを最も重要な課題としている。また、確率・統計理論の紹介を本の後半に置き、実験や世論調査などを用いた具体的な研究に最初に触れることによって、統計検定や信頼区間といった抽象的な概念とその必要性も本書の読者にとっては理解しやすくなっているはずである。英語版と異なり、日本語版は上下巻に分かれているが、下巻はテキスト、ネットワーク、そして地理データの分析といった新たな分野を扱っているだけでなく、確率論や統計理論の基礎もカバーしているので、ぜひ上下巻を通して勉強してもらいたい。

日本では、大学入試時から理系と文系の区別がなされており、理系の手法を用いて文系の学問を勉強するということがなかなかやりにくい環境がある。しかしながら、いわゆる「ビッグデータ」時代の現代社会においては、学際的なアプローチが必要不可欠になっている。社会問題をデータ分析を用いて解決していくためには、統計や機械学習の知識と社会科学の観点を有効に組み合わせることが大切であり、本書が学部・大学院教育における文理融合を大胆に進めていく1つのきっかけになることを期待している。

大学までを日本で過ごした私にとって、英語で書いた自分の本が日本語に訳

されるというのは、非常に感慨深い。運よく学部時代の交換留学中に、アメリカで計量社会科学に出会った私にとって、日本語で多くの日本の学生に、社会科学のためのデータ分析の魅力を伝える機会を得ることは長年の夢であった。その夢を現実のものにして下さった、慶應義塾大学の粕谷祐子先生、福岡大学の原田勝孝先生、そして大阪大学の久保浩樹先生にはこの場を借りて、感謝の意を表したい。特に粕谷先生は、英語での執筆が終わる前からこの本の意義をいち早く理解されて翻訳に取り組んでくださり、非常に勇気づけられた。

最後になるが、人生の良きパートナーであるクリスティーナ、そして小さい頃から自由奔放に育ててくれた両親に本書を捧げたい。

2017年8月

アメリカ合衆国ニュージャージー州プリンストンにて
今井耕介

目 次

日本語版への序文

1	イントロダクション	1
1.1	本書の概観	4
1.2	本書の使い方	10
1.3	Rの基礎	13
1.3.1	算術演算	14
1.3.2	オブジェクト	16
1.3.3	ベクトル	20
1.3.4	関数	23
1.3.5	データファイル	27
1.3.6	オブジェクトを保存する	32
1.3.7	パッケージ	34
1.3.8	プログラミングと学習のコツ	35
1.4	まとめ	38
1.5	練習問題	38
1.5.1	自己申告に基づく投票率のバイアス	38
1.5.2	世界人口の動態を理解する	40
2	因果関係	45
2.1	労働市場における人種差別	45
2.2	Rでデータを部分集合化する	51
2.2.1	論理値と論理演算子	51
2.2.2	関係演算子	54
2.2.3	部分集合化	56

2.2.4	簡単な条件文	60
2.2.5	因子変数	61
2.3	因果効果と反事実	64
2.4	ランダム化比較試験	68
2.4.1	ランダム化の役割	68
2.4.2	社会的プレッシャーと投票率	71
2.5	観察研究	76
2.5.1	最低賃金と失業率	76
2.5.2	交絡バイアス	80
2.5.3	事前・事後の比較と差の差分法	84
2.6	1変数の記述統計量	88
2.6.1	分位数	89
2.6.2	標準偏差	93
2.7	まとめ	95
2.8	練習問題	96
2.8.1	初期教育における少人数クラスの有効性	96
2.8.2	同性婚に関する意見の変化	99
2.8.3	自然実験としての指導者暗殺の成功	101
3	測 定	103
3.1	戦時における民間人の被害を測定する	103
3.2	Rで欠損データを扱う	106
3.3	1変量の分布をビジュアル化する	109
3.3.1	棒グラフ	110
3.3.2	ヒストグラム	111
3.3.3	箱ひげ図	116
3.3.4	グラフの印刷と保存	119
3.4	標本調査	120
3.4.1	ランダム化の役割	121
3.4.2	無回答とその他のバイアス発生要因	126
3.5	政治的分極化を測定する	131

3.6	2 変量関係の要約	133
3.6.1	散 布 図	133
3.6.2	相 関	137
3.6.3	Q-Q プロット	142
3.7	クラスター化	145
3.7.1	Rにおける行列	145
3.7.2	Rにおけるリスト	148
3.7.3	k 平均法	150
3.8	ま と め	156
3.9	練習問題	157
3.9.1	同性婚に関する意見の変化再考	157
3.9.2	中国とメキシコにおける政治的有効性感覚	159
3.9.3	国連総会における投票	162
4	予 測	165
4.1	選挙結果の予測	165
4.1.1	Rにおけるループ(繰り返し)	167
4.1.2	Rにおける一般的な条件文	170
4.1.3	世論調査からの予測	175
4.2	線形回帰	186
4.2.1	顔の見た目と選挙結果	187
4.2.2	相関と散布図	189
4.2.3	最小2乗法	191
4.2.4	平均への回帰	198
4.2.5	Rにおけるデータの結合	201
4.2.6	モデルの当てはまり	210
4.3	回帰分析と因果関係	217
4.3.1	ランダム化実験	218
4.3.2	重回帰モデル	222
4.3.3	不均一トリートメント効果	229
4.3.4	回帰分断デザイン	238

4.4	ま と め	244
4.5	練習問題	246
4.5.1	賭博市場に基づく予測	246
4.5.2	メキシコにおける選挙と条件付き現金給付プログラム	248
4.5.3	ブラジルにおける政府間移転支出と貧困削減	251

事項索引

R 索引

下巻の目次

5 発 見

- 5.1 テキスト・データ
- 5.2 ネットワーク・データ
- 5.3 空間データ
- 5.4 ま と め
- 5.5 練習問題

6 確 率

- 6.1 確 率
- 6.2 条件付き確率
- 6.3 確率変数と確率分布
- 6.4 大標本理論
- 6.5 ま と め
- 6.6 練習問題

7 不確実性

- 7.1 推 定
- 7.2 仮説検定
- 7.3 不確実性を伴う線形回帰モデル

7.4 ま と め

7.5 練習問題

8 次の一歩

表 目 次

1.1	swirl 復習問題	12
1.2	世界人口推定	21
1.3	アメリカの投票率データ	39
1.4	出生数・死亡数の推定データ	41
2.1	履歴書実験データ	46
2.2	論理積と論理和	53
2.3	因果推論の潜在的結果の枠組	66
2.4	社会的プレッシャー実験データ	74
2.5	最低賃金研究データ	77
2.6	STAR プロジェクトのデータ	97
2.7	同性婚データ	99
2.8	指導者暗殺データ	101
3.1	アフガニスタンに関するサーベイ・データ	104
3.2	アフガニスタンの村に関するデータ	124
3.3	議員の理想点に関するデータ	133
3.4	アメリカのジニ係数データ	139
3.5	同性婚に関するデータを作り直したもの	158
3.6	CCAP サーベイ・データ	158
3.7	ヴェニエット形式のサーベイ・データ	161
3.8	国連の理想点データ	163
4.1	2008 年アメリカ大統領選挙データ	176
4.2	2008 年アメリカ大統領選挙世論調査データ	176
4.3	混同行列	183
4.4	顔の見た目実験のデータ	188

4.5	2012年アメリカ大統領選挙データ	201
4.6	フロリダ州の郡レベルでの1996年と2000年の アメリカ大統領選挙データ	211
4.7	政策立案者としての女性データ	218
4.8	イギリス国会議員の個人資産データ	239
4.9	予測市場データ	246
4.10	2012年アメリカ大統領選挙世論調査データ	248
4.11	条件付き現金給付プログラムデータ	250
4.12	ブラジルの政府間移転支出データ	252

目 次

1.1	RStudio のスクリーンショット	15
1.2	RStudio テキストエディタのスクリーンショット	36
2.1	「名前公表」投票推進メッセージ	72
2.2	最低賃金研究における差の差分法	86
3.1	誤った見出しの付いた新聞を掲げるハリー・トルーマン	123
3.2	自然対数	125
3.3	空間投票モデル	132
3.4	ジニ係数とローレンツ曲線	138
4.1	2008 年アメリカ大統領選挙での選挙人団の投票先マップ	166
4.2	実験で使われた候補者の写真の例	187
4.3	散布図におけるデータクラウドの相関係数と形状	190
4.4	ゴルトンの平凡への回帰	199
4.5	パームビーチ郡のチョウ型投票用紙	215

第 1 章 イン트로ダクション

神のことは信頼する。他の者はみな、データをもってこななければならない。

—ウィリアム・エドワーズ・デミング

計量社会科学は、経済学、教育学、政治学、公共政策学、心理学、社会学など幅広い分野を含む学際領域である。計量社会科学では、社会や人間行動に関する問題を理解し、解決するためにデータ分析が行われる。例えば、労働市場における人種差別に関する研究、新しいカリキュラムが生徒の学力に与える影響の評価、選挙結果の予測、ソーシャルメディアの使い方に関する分析などがある。データに基づく同様のアプローチが、隣接分野である衛生学、法学、ジャーナリズム、言語学、そして文学でもとられている。社会科学は現実世界の様々な問題を直接調査するので、その研究結果は社会の個人、政策、ビジネスなどに直接影響を及ぼす大きな可能性を秘めている。

ここ数十年のあいだに、計量社会科学は様々な分野において驚くべきスピードで隆盛をきわめてきた。データ分析によって実証的な証拠を示す学術論文の数は、飛躍的に増加した。学術研究以外の分野でも、多くの組織 — 企業や政治キャンペーン、ニュースメディア、政府機関など — が、データ分析の結果に基づいて意思決定を行うようになってきている。こうした計量社会科学の急速な発展は、環境を一変させる技術面の2つの変化によってもたらされた。第1に、インターネットによってデータ革命(data revolution)が大きく進み、使えるデータの量と多様性が急増した。情報共有により研究者や組織はデジタル形式で膨大な量のデータセットを広めることが可能となったのである。第2に、ソフトウェアとハードウェアの両面における計算革命(computational revolution)のおかげで、誰もが自分のパソコンと好きなデータ分析ソフトウェアとを使えるようになった。

これら技術的変化の直接的な結果として、計量社会科学の研究者が入手可能なデータの量は急速に増大した。ひと昔前であれば、研究者が使えるデータと

例えば大部分は政府機関が発行したもの(例えば、国勢調査、選挙結果、経済指標など)であり、それに加えて研究者グループが集めた少数のデータセット(例えば、国政選挙に関するサーベイ(世論調査)データ、戦争発生や民主主義の諸制度に関して研究者が独自に変数化したデータセットなど)があるだけだった。これらのデータセットは、実証分析を進める上で依然として重要な役割を果たしている。しかし、様々な新しいタイプのデータが登場したことによって、計量社会科学研究が扱うことのできる領域はより広範なものとなった。研究者はいまや、ランダム化実験やサーベイを自分自身でデザインし実行することができる。透明性やアカウントビリティを高める圧力がかかるようになり、政府機関はオンライン上でのデータ公開を進めている。例えばアメリカでは、選挙献金やロビー活動に関する詳細なデータを、誰でも自分のパソコンにダウンロードできる。また、スウェーデンなどの北欧諸国では、所得、税金、教育、健康、職場などの多岐にわたる登録情報を学術研究に使用することができる。

新しいデータセットが、幅広い分野で登場している。電子購買履歴を通じて、消費者の購買情報に関する詳細なデータが入手できるようになった。国際貿易データは、製品ごとに多数の国のペアについて、数十年分にわたり集められている。軍隊もデータ革命に貢献してきた。2000年代のアフغانستان紛争において、アメリカや国際治安支援部隊は、反乱軍攻撃の位置情報、発生時刻、種類などに関するデータを集め、反乱対策戦略の指針を定めるためにデータ分析を行っている。同様に、政府機関や非政府組織が、戦争による民間人死傷者に関するデータを集めている。政治キャンペーンでは有権者動員の戦略を立てるためにデータ分析が用いられ、特定のタイプの有権者層をターゲットとして入念に選んだメッセージを発信している。

これらのデータセットの形態も様々である。計量社会学者は、法案や新聞記事、政治家の演説など様々な電子テキストをデータとして分析している。ウェブサイトやブログ、ツイッター、SMS メッセージ、フェイスブックなどを通じてソーシャルメディア・データを入手できるようになったため、人々がオンライン上でどのようにやりとりしているか調査することが可能となった。地理情報システム(geographical information system; GIS)のデータセットも発達し、空間的位置に注意を払いつつ、選挙区の区割り変更のプロセスや内戦を分析できるようになった。衛星画像データを用いて、途上国の農村部でどの程度電力が普及しているかを調査した研究者もいる。まだ少数ではあるものの、社

会科学的な問いに答えるために、計量的な手法を用いて画像や音声、さらには動画を分析することもある。

情報技術革命に伴って豊富で多様なデータを入手できるようになったことで、学者から実務家まで、企業のアナリストから政策立案者まで、学生から教員まで、あらゆる人がデータに基づいた発見をすることが可能になった。従来は、統計学者や一部の専門家のみがデータ分析を行っていたのに対して、現在では、誰もがパソコンを起動し、インターネットからデータをダウンロードし、自分の好きなソフトウェアを用いて分析を行うことができるのである。こうした変化により、政策の有効性を示す上でも、従来に比べてより一層、説得力のある説明が求められるようになってきている。例えば、非政府組織や政府機関は、政策やプログラム実施の資金を確保し正当性を高めるため、厳格な評価を行い、その有効性を示さなければならなくなった。

透明性の確保やデータに基づいた発見が重視されるようになったことで、社会科学分野の学生は、どのようにデータを分析し、どのように結果を解釈し、そして、そこから得られた知見をどう効果的に公表するかについて学ぶ必要が出てきた。伝統的な統計学の入門コースでは、手計算、あるいはせいぜい関数電卓を用いた計算をさせることで、学生に統計学の基礎的な概念を教えることが多かった。もちろん、そうした概念は依然として重要であり、本書でもそれらを取りあげる。しかし伝統的なアプローチでは、今日の社会的要求に応えられない。一般的な統計学の概念や方法論を学ぶだけでは、「統計リテラシー」が十分に身につかないのである。データから引き出すことができる豊富な情報を余すことなく活用し、データに基づいた発見を通して社会に貢献していくために、社会科学を学ぶ者は皆、基本的なデータ分析のスキルを身につける必要がある。

誰でもデータ分析ができるようになるべきだと考えたことが、本書を書くと思った最大の動機である。この本では、計量社会科学の研究に必要とされるデータ分析の3つの要素について解説していく。すなわち、研究の背景に関する知識、プログラミング技術、そして、統計手法である。これらの要素はどれも単独では不十分である。研究の背景に関する知識なしには、データ分析に必要な仮定の信頼性を評価することができず、また、観察された知見が何を意味しているのか理解することができない。プログラミング技術がなければ、データを分析し、リサーチ・クエスチョンに答えることもできない。統計原理

に則らなければ、シグナルと呼ばれる体系的なパターンとノイズと呼ばれる例外的なパターンを区別することができず、誤った推論をしてしまう可能性がある(ここでは、推論とは、観察データに基づいて、未知の数量について結論を導き出すことを指す)。この3つの要素を組み合わせ、本書はデータ分析の威力を示す。

1.1 | 本書の概観

本書は、データ分析と統計学の初学者向けの本である。社会科学やその他の分野の研究者、大学院生、大学生だけではなく、実務家や意欲的な高校生なども読者として想定している。初歩的な代数を除けば、必要な予備知識はない。微積分や確率論の知識なしに読み進めることができる。プログラミング経験は、あるに越したことはないが、不可欠ではない。また、データ分析をほとんど教えない、伝統的な手計算による統計学入門を履修した人にとっても役に立つ。学生はこの本を通して、データ分析の面白さを体感できるだろう。ここでは、社会科学적인問いに答えるためにRをどう使うかに焦点を当てているが、Rを用いたプログラミングを学びたい人にも本書は役立つだろう。

上で述べたように、本書の特徴は、公開されている計量社会科学研究から直接もってきたデータの分析を通して、プログラミング技術と統計学の概念とを同時に解説しようとする点にある。目指すところは、社会の問題や人間行動に関する重要な問いに答えるために、社会学者がどのようにデータ分析を行うかを示すことである。同時に、本書を読むことで、基礎的な統計学の概念や初歩的なプログラミング技術を身につけることもできる。さらに大事なこととして、約40のデータセットを検討することで、データ分析を経験することができる。

本書は、8つの章から構成されている。この第1章では、本書の最善の活用方法を説明し、広く普及しているオープンソースの統計プログラミング環境であるRの手短な解説を行う。Rは、無料でダウンロードすることができ、Mac, Windows, Linuxのコンピュータで使用できる。読者の皆さんには、データ分析をより簡単に行うことができる無料のソフトウェア・パッケージRStudioの使用を強くお勧めする。この章の最後には、出版された社会科学研究で用いられたデータセットを使って、初歩的なRの機能を練習するための練習問

題を2つ設けている。この本で使用しているデータセットはすべて <http://press.princeton.edu/qss/> から無料でダウンロードできる。また、このウェブサイトには、各章の復習問題など、有用な教材へのリンクもある。なお、第5章を除き、Rの最も基本的なシンタックス(構文)を使うことがほとんどであり、様々な追加パッケージを導入することはない。しかし、この本を読み終える頃には、他のパッケージも使えるRプログラミング技術が身についているだろう。

第2章では、因果関係(causality)の基礎的解説を行う。社会科学の研究において、因果関係は、特定の政策やプログラムが、問題としている結果を変化させるかどうかを調べようとするときには、常に決定的に重要な役割を果たす。しかし、因果関係の分析は、観察することができない反事実を推論しなければならないため、非常に困難であることがよく知られている。例えば、労働市場における人種差別の存在を理解するためには、審査通過の連絡をもらえなかった黒人と同じ人物が、もし白人であれば連絡がもらえたのか、を知らなければならない。この章では、黒人らしい求職者の名前と白人らしい求職者の名前を研究者がランダムに選んで、架空の求職者の履歴書を求人企業に送った、有名な実験研究のデータを分析する。この研究を応用例として、トリートメント(処置)の割り当てをランダム化することでどのようにトリートメントの平均因果効果を特定できるかを説明する。

また、トリートメントの割り当てを研究者がコントロールしない、観察研究における因果推論についても説明する。ここでの主な応用例は、最低賃金の上昇が雇用に与える影響を解明しようとした古典的な研究である。多くの経済学者は、最低賃金を上げると雇用が減少する可能性があると主張する。その根拠は、最低賃金を上げると雇用主は労働者により高い賃金を支払わなければならないために、雇用する労働者の数を減らさざるを得ないからである。残念なことに、最低賃金の引き上げはランダムに決定されているわけではなく、経済成長のように、それ自体が雇用と強く関係している数多くの要因に影響を受けている。そうした複数の要因によって、企業がトリートメントグループになるかが決定されるため、トリートメントを受けたグループと受けていないグループを単純に比較しただけでは、バイアス(偏り)を含んだ推論になってしまう。

そこで、観察研究におけるこの種の選択バイアスを減らそうとする方策をいくつか紹介する。観察研究はトリートメントの効果を不正確に推定してしまう

恐れがあるものの、ランダム化比較試験によって得られた結果よりも一般化しやすい場合が多い。また、この章では、投票推進運動における社会的プレッシャーに関するフィールド実験も取りあげる。第2章の練習問題では、初等教育における少人数クラスの因果効果を調査するランダム化実験や、政治的指導者の暗殺とその効果に関する自然実験を取りあげる。Rのプログラミングに関しては、論理式や部分集合化(一部を抜き出すこと)について説明する。

第3章では、測定(measurement)という基本的な概念を説明する。測定におけるバイアスは誤った結論や誤解に基づく決定につながる恐れがあるため、正確に測定することはどのようなデータに基づく発見においても重要である。まず、標本サーベイ調査でどのように世論を測定すべきか検討する。ここでは、アフガニスタン紛争の期間中に、アフガニスタン市民の間での国際治安支援部隊とタリバン反乱軍への支持がそれぞれどの程度だったのかを測定した研究のデータを分析する。この分析を通じて、サーベイを行う際の抽出において、ランダム化がいかに役立つか説明する。具体的にいえば、母集団から回答者をランダムに抽出することで、母集団を代表した標本を手に入れることができる。その結果、1つの小規模な代表的集団によって、母集団全体の意見を推論することが可能となるのである。また、標本抽出の潜在的なバイアスについても検討する。無回答があると、標本の代表性が損なわれる可能性がある。偽った回答は推論にとって深刻な脅威となる。特に、タリバンの反乱を支持するかどうかというような、デリケートな問題について尋ねられた場合にそのようなことが起こる。

第3章の後半では、計量社会科学において重要な役割を果たすが直接には観察できない概念の測定に焦点を当てる。こうした概念の有名な例には能力やイデオロギーがあるが、この章では政治イデオロギーを取りあげる。まず、議員のイデオロギー位置を点呼投票から推論するためによく用いられるモデルを説明した後、アメリカ議会が時代とともにどのようにして分極化してきたのかを検討する。次に、基本的なクラスタリング計算手法である k 平均法を解説する。この手法により、似たような観察グループを見つけることができる。この手法をデータに適用することで、議会におけるイデオロギーの分極化が近年では政党の違いで特徴づけられていることがわかる。対照的に、それ以前の時期は、それぞれの政党内にイデオロギーの分裂があった。また、本章では、分位数や標準誤差、ジニ係数といったデータのばらつきの尺度も解説する。Rの

プログラミングに関しては、1変量データや2変量データのビジュアル化を解説する。練習問題では同性婚に関する実験を改めて分析する。物議をかもしたこの実験は、この章で取りあげる方法論の例であると同時に、学問的な誠実さに関する問題を突き付けるものである。

第4章では、予測(prediction)について検討する。ある事象が起こるかどうかを予測することは、政策決定や意思決定のプロセスにおいて非常に重要である。例えば、財政計画をたてる際には経済パフォーマンスを予測する必要があるし、対外政策を決定する際には、早い段階で市民の不穏な動きを察知し、事前に対策を講じることが重要となる。本章で主に取りあげる応用例は、選挙前の世論調査を用いたアメリカ大統領選挙の結果予測である。複数の世論調査データを組み合わせることで、比較的簡単に、きわめて正確な予測ができることを明らかにする。加えて、ある心理学実験のデータを分析する。この実験では被験者に知らない候補者の顔写真を見せて、候補者の能力に点数を付けてもらう。この分析によれば、驚くべきことに、一瞬の顔の印象で選挙結果を予測することができる。この例を通して、ある変数から他のある変数の値を予測する際によく用いられる線形回帰モデルを解説する。また、線形回帰と相関関係との関係を説明した上で、「回帰」という言葉の由来となった「平均への回帰」という現象を考察する。

第4章では、さらに、単に予測を行うだけではなく因果効果を推定するために回帰モデルが使える場合について検討する。因果推論は、予測変数(predictor)としてトリートメント変数を使い、観察された結果ではなく反事実の結果を予測するという点で、標準的な予測とは異なる。ここでは、インドにおいてランダムに選ばれた村で村議会のいくつかの議席が女性に割り当てられる事例を用いた、ランダム化自然実験のデータを分析する。この自然実験を通して、女性政治家の存在が政策に影響を及ぼすかどうか、特に、女性の有権者が気にかける政策課題に影響を与えるかどうかを検討する。また、この章では、観察研究で因果推論を行うために用いられる回帰分断デザインについても解説する。これに関しては、イギリスにおいて、政治家の資産がどの程度、政治的役職に就いていることによって得られたものなのかを考える。この問いに答えるため、選挙に辛うじて当選した候補者と惜しいところで落選した候補者の比較を行う。さらに、この章では、少々難しいが役に立つRプログラミングの概念である、ループ(繰り返し)と条件文について解説する。章末の練習問題で

は、選挙予測の賭博市場が選挙結果を正確に予測することができるのかという分析などを行う。

第5章では、様々な種類のデータからのパターンの発見(discovery)について解説する。「ビッグデータ」を分析する際には、データにある体系的なパターンを特定するために、機械的なデータ処理の手法やビジュアル化のためのツールが必要である。まず、テキスト(文章)をデータとして分析する。ここでは主に、合衆国憲法の基となった論文集『ザ・フェデラリスト』の著者が誰なのかを予測する。この論文集の中には著者が明らかなものがある一方で、著者不明のものもある。著者が明らかな文章に出てくる特定の単語の頻度を分析することで、著者不明とされてきた論文を執筆したのはアレクサンダー・ハミルトンとジェームズ・マディソンのどちらなのかを予測できる。次に、分析ユニット間の関係に焦点を当て、ネットワーク・データの分析方法を示す。ルネサンス期フィレンツェにおける姻戚関係ネットワークの中で、メディチ家が重要な役割を果たしていたことを定量的に示す。より現代的な例として、中心性に関する様々な尺度を紹介し、アメリカの上院議員がツイッターにおいて生成するソーシャルメディア・データに応用する。

第5章では、地理空間データを最後に解説する。まず、空間データ分析の古典であるジョン・スノーの研究を検討する。これは、1854年にロンドンでコレラが大流行した原因の研究である。次に、アメリカの選挙データを例に、地図を作成しながら空間データをビジュアル化する方法を示す。時空間データを扱う際には、時間とともに変化する空間のパターンをビジュアル化するために、いくつものマップを連続的に重ねてアニメーションを作成する。そのためこの章では、ビジュアル化専用のいくつかのRパッケージを用いながら、様々なデータビジュアル化の技術を応用する。

第6章では、データ分析の話から、不確実性に関する統一的な数学モデルの1つである確率(probability)に焦点を移す。これまでの章では、どのようにパラメーター(母集団についての関心のある値)を推定し予測を立てるか検討してきたが、実証的知見に含まれる不確実性に関しては論じていなかった。この点については第7章で説明する。確率は推論の不確実性を定量的に示す統計的推論の基礎であるため重要である。ここではまず、頻度論者とベイズ論者という2つの主要な観点から、確率をどのように解釈すべきかという問いについて考える。次に、確率と条件付き確率の数学的定義を与えた後、確率論の基

本的な法則について解説する。その1つが「ベイズの公式」である。ベイズの公式を用いることで、サーベイ・データが入手できなくても、名字と住所から個人のエスニシティを正確に予測できることを示す。

また、この章では、確率変数と確率分布という重要な概念についても説明する。これらの概念ツールを用いて、第4章の選挙前の世論調査データを用いた選挙結果の予測に不確実性の尺度を付け加える。また、賭博市場のデータを用いた選挙結果予測の不確実性に関する練習問題もある。最後に、大数の法則と中心極限定理という2つの基本的な定理の説明を行う。この2つの定理は、様々な場面で使うことができ、標本サイズが大きくなり標本抽出を何回もするにつれて推定結果がどのように変わるかを理解する上で役に立つ。章末の練習問題では、第2次世界大戦中にドイツで使われていた暗号機エニグマと、ロシアにおける選挙不正の検出を扱う。

第7章では、推定結果や予測の不確実性(uncertainty)をどのように定量的に示すのかを検討する。これまでの章では、データに存在するパターンを見つけるために使われるデータ分析の手法をいくつも紹介してきた。第6章の知識を土台として、第7章では、そのようなパターンに対し、どの程度の確実性を見込めるかについて詳細に解説する。この章では、標準誤差や信頼区間の計算、仮説検定の使用などを通してノイズとシグナルを区別する方法を紹介する。つまり、この章のテーマは統計的推論である。これまでの章で出てきた例を用いているが、ここでの焦点は、計算された推定値がどの程度不確実なのか、である。扱うのは、選挙前に行われた世論調査の分析や、初等教育におけるクラスの人数が生徒の学力に及ぼす影響に関するランダム化実験、最低賃金の上昇が雇用に与える影響を評価した観察研究などである。また、統計的仮説検定について検討する際には、多重検定(multiple testing)や「出版バイアス」の危険性についても注意を喚起する。最後に、線形回帰モデルから導かれた推定値の不確実性がどの程度であるかをどのように定量化すべきか検討する。ここでは再び、インドの女性政治家に関するランダム化自然実験や、イギリスの政治家が政治的地位に就くことによって蓄積できる資産量を推定するための回帰分断デザインについて検討する。

最終章では、この本を読み終えた後で次に何をすべきか手短かに説明する。また、計量社会科学の研究においてデータ分析が果たす役割についても検討する。

1.2 | 本書の使い方

この節では本書の使い方を説明するが、これは次の原則に則っている。

データ分析はやってみて身につくもので、読むだけでは身につかない。

この本はただ読むようには書かれていない。実際にデータ分析を経験することが何よりも重要である。それには本書に出てくるコードを自分で試し、あれこれ変えてみて、各章の章末にある練習問題に取り組むのが一番である。この本で使うコードやデータセットはどれも <http://press.princeton.edu/qss/>にあるリンクから無料でダウンロードできる。

この本は段階を踏みながら解説を行う。後の章は、前の方で取りあげた題材の多くに読者がなじんでいることを前提としている。そのため、途中の章を飛ばすことはお勧めしない。例外は第5章「発見」で、ここの内容は後の章で使わない。とはいえこの章は、本書の中でも特に興味深いデータ分析の例を扱っており、取り組むことを勧めたい。

この本は、様々な形で使うことができる。例えば、伝統的な統計学の入門コースで、データ分析演習の副読本として、この本の全体または一部を課すことができる。この本の一番よい使い方は、一方的な講義にあまり時間を費やすことなく、クラス内で教員が学生とやりとりしながらデータ分析の演習を行うようなコースで使うというものである。そうしたクラスでは、授業の前に学生は本書の該当箇所を予習しておき、授業では、新しく出てきた手法やプログラミング技術について教員が解説を行った上で、本書の練習問題や自分たちで選んだ他の同様の応用例にそれらを使ってみる。このようなプロセスを経ることで、クラス全員が答えにたどり着くまでに、場合によってはソクラテス式問答法も使って、教員と学生が練習問題に関して双方向的に議論できる。クラスでのディスカッションの後で、少数の学生が教員と共に別の練習問題に取り組むようなコンピュータ実習室でのフォローアップを行えると理想的である。

この教え方は「個別・一般・個別」原則に沿っている¹⁾。この原則によれ

1) Frederick Mosteller (1980) "Classroom and platform performance." *American Statistician*, vol. 34, no. 1 (February), pp. 11-17.