

機械学習の数理 100 問シリーズ
統計的機械学習の数理 100 問 with R

鈴木 譲

平成 32 年 2 月 20 日

シリーズ序文

機械学習の書籍としておびただしい数の書籍が出版されているが、ななめ読みで終わる、もしくは難しすぎて読めないものが多く、「身につける」という視点で書かれたものは非常に少ないと言ってよい。本シリーズは、100の問題を解くという、演習というスタイルをとりながら、数式を導き、R言語もしくはPythonのソースプログラムを追ひ、具体的に手を動かしてみて、読者が自分のスキルにしていくことを目的としている。

各巻では、各章でまず解説があり、そのあとに問題を掲載している。解説を読んでから問題を解くこともできるが、すぐに問題から取り組む読み方もできる。その場合、数学の問題において導出の細部がわからなくても、解説に戻ればわかるようになっている。

「機械学習の数理100問シリーズ」は、2018年以降に大阪大学基礎工学部情報科学科数理科学コース、大学院基礎工学研究科の講義でも使われ、また公開講座「機械学習・データ科学スプリングキャンプ」2018, 2019でも多くの参加者に解かれ、高い評価を得ている。また、その間に改良を重ねている。講義やセミナーでフィードバックを受け、洗練されたものだけを書籍のかたちにしていく。

プログラム言語も、大学やデータサイエンスで用いられているR言語と企業や機械学習で用いられているPythonの2種類のバージョンを出す。これも本シリーズの特徴の一つである。

本シリーズのそれぞれの書籍を読むことで、機械学習に関する知識が得られることはもちろんだが、脳裏に数学的ロジックを構築し、プログラムを構成して具体的に検証していくという、データサイエンス業界で活躍するための資質が得られる。「数理」「情報」「データ」といった人工知能時代を勝ち抜くための、必須のスキルを身につけるためにうってつけのシリーズ、それが本シリーズである。

まえがき

人工知能時代をエンジョイするのか，人工知能のエジキになるのか

2016年に、大阪大学理学部数学科の学生を対象に（統計的）機械学習に関する講義をはじめたことが、本書を執筆するきっかけになったと思います。当時、テキストとして適当なものがなく、2018年に日本語に翻訳された『Rによる統計的学習入門』“*Introduction to Statistical Learning with Application in R*”を参考に講義を組み立てました。説明が丁寧で、わかりやすく、非常に気に入って、日本語の翻訳をしてみたいと思うぐらいになりました（その時点で他の方が権利をすでに獲得していて、翻訳は実現しませんでした）。本書の扱っている単元や、章立てもその著述と似ているのも、そうしたところによっています。

2017年に現在の所属、基礎工学部情報科学科数理科学コースに異動となりました。早速、3年生後期の計算数理Bという科目で、2016年と同様の内容の講義を試みました。数理科学コースは2年生でR言語を習熟していて、当初は違和感なく首尾よく講義を行えました。しかし、よりどころとしていた『Rによる統計的学習入門』については、2018年、2019年というように講義を重ねるごとに、問題意識を抱くようになりました。初学者の方が機械学習の概要を把握するには、最適な書籍とは思われますが、

1. 本質的な理解というよりは、感覚的な理解だけで十分
2. R言語の処理のステップを見ないで、パッケージにデータを放り込めば十分

という姿勢には、どうしても合点がいきませんでした。

現在は、人工知能の時代と言われます。言うまでもありませんが、インターネットのおかげで、必要な情報を即座に得ることができ、職場の業務が効率化され、生活も豊かになりました。その一方で、人間が行ってきた業務の多くが人工知能に置き換えられるのではないかという危惧が生じています。私自身は、データサイエンスや機械学習の業務に携わる人との付き合いが多いのですが、「業務でどのような資質が求められますか」と聞くと、知識や経験というよりは、「ロジック」という答える人が多く、活躍している人ほどその傾向が強いように思われます。情報の真偽を吟味する、人が見えない本質やチャンスを見る、制約にとらわれない発想などが、「ロジック」の結果として生まれているという視点です。逆に、そういう「ロジック」が欠如していて、ヤマカンに頼るというのであれば、人工知能のエジキになる可能性が高くなるように思われます。

もっとも、そうした「ロジック」が、数学やプログラミングをやらないと身につかない、というのは真実ではないかもしれませんが、大学教員30年間の中で、多くの学生をみてきた経験からす

ると、それらには非常に強い相関があるように思われます。本書は、機械学習に関する知識も提供しますが、それと同時に、数学的に本質を理解して、プログラミングで処理を構成して、検証するという経験を通じて、読者の方々の「ロジック」を脳裏に構築することを、目標の一つに掲げています。

『Rによる統計的学習入門』は、講義を組み立てる上で大変参考にはなりましたが、人工知能時代をエンジョイするために不可欠な「ロジック」を構築するという視点が十分ではなく、どうしても本書を執筆せざるを得なかった、というのが私の正直な気持ちです。

また、数年前に日本語に翻訳された“*Elements of Statistical Learning*”『統計的学習の基礎』（共立出版）は、分量が多く、輪講して挫折したという話をよく聞きます。困ったときに参考にする百科事典として使う方が多いように思われます（私は頻繁に利用しています）。『統計的学習の基礎』を大学の講義半期に圧縮できたら、という気持ちも本書を執筆するきっかけとなりました。また、『統計的学習の基礎』は信頼できる知識を提供していますが、ソースコードをおいたり、読者が具体的にスキルを身につけるような誘導があってもよいように思いました。

本書の100問は、2017～2019年の講義で学生に課した演習問題、特に『Rによる統計的学習入門』を漫然と読んでいる学生にツッコミを入れるために作成したものです。数式を正しく導出し、プログラムを組んで実行結果を見るなど、手を動かして自分のものにすべきと伝えました。ただ、問題だけを与えても、優秀な学生以外は自力では解けないので、全員が課題を提出できるよう、講義で解答に限りなく近い丁寧なヒントを与えました。それが本書の本文です。問題1～100と本文で重複している記述が若干あるのは、そのためです。また、講義で提供している10～15分間の復習ビデオも提供しています。

本シリーズの特徴

本書というよりは、本シリーズの特徴を以下のようにまとめてみました。

1. 身につける：ロジックを構築する

数学で本質を把握し、プログラムで処理を構成して、データを処理していきます。読者の皆さんの脳裏に、「ロジック」を構築していきます。機械学習の知識だけではなく、視点が身につきますので、新しい機械学習の技術が出現しても追従できます。100問を解いてから、「大変勉強になりました」と言う学生がほとんどです。

2. お話だけで終わらない：コードがあるのですぐにコード（行動）に移せる

機械学習の書籍でソースプログラムがないと、非常に不便です。また、パッケージがあっても、ソースプログラムがないとアルゴリズムの改良ができません。Gitなどでソースが公開されている場合もありますが、MATLABやPythonしかなかったり、十分でない場合もあります。本書では、ほとんどの処理にプログラムのコードが書かれていて、数学がわからなくても、それが何を意味するかを理解できます。

3. 使い方で終わらない：大学教授が書いた学術書

パッケージの使い方，実行例ばかりからなる書籍も，よく知らない人がきっかけを掴めるなど，存在価値はありますが，手順にしたがって機械学習の処理を実行できても，どのような動作をしているかを理解できないので，満足感として限界があります。本書では，機械学習の各処理の数学的原理とそれを実現するコードを提示しているので，疑問の生じる余地がありません。本書はどちらかというと，アカデミックで本格的な書籍に属します。

4. 100 問を解く：学生からのフィードバックで改善を重ねた大学の演習問題

本書の演習問題は，大学の講義で使われ，学生からのフィードバックで改良を重ね，選びぬかれた最適な 100 問になっています。そして，各章の本文はその解説になっていて，本文を読めば，演習問題はすべて解けるようになっています。

5. 書籍内で話が閉じている (self-contained)

定理の証明などで，詳細は文献〇〇を参照してください，というように書いてあって落胆した経験はないでしょうか。よほど興味のある読者（研究者など）でない限り，その参考文献をたどって調査する人はいないと思います。本書では，外部の文献を引用するような状況を避けるように，題材の選び方を工夫しています。また，証明は平易な導出にし，難しい証明は各章末の付録においています。本書では，付録まで含めれば，すべての議論が完結しています。

6. 読者ページ：章ごとのビデオ，オンラインの質疑応答，プログラムファイル

大学の講義では，slack で 24/365 体制で学生からの質問に回答していますが，本書では Disqus という質疑応答のシステムを採用しています。また，講義で利用した各章 10～15 分のビデオを公開しています。また，本書にあるプログラムをすべて手で打ち込むのは難しいので，プログラムのリストのファイルをサイトにおいています。

7. 線形代数

機械学習や統計学を学習するうえでネックになるのが，線形代数です。研究者向きのものを除くと，線形代数の知識を仮定しているものは少なく，本質に踏み込めない書籍がほとんどです。そのため，シリーズ第 1 号の『統計的機械学習 100 問 with R』と第 2 号の『同 with Python』では，第 0 章として，線形代数という章を用意しています。14 ページしかありませんが，例だけでなく，証明もすべて掲載しています。ご存知の方はスキップしていただいて結構ですが，自信のない方は休みの日を 1 日使って読まれてもよいかと思います。

本書の使い方

各章は，問題，その解説（本文），付録（証明，プログラム）からなっています。問題を解いてみて，わからないときだけ本文を読める方（上級者）もいらっしゃるでしょうが，本文から読み始めて最後に問題をとくという形でも問題ありません。最後まで読破して，身につけることを優先してください。

講義で使われる場合、第1章を3回、第4章を1～2回、第6章を2回、それ以外の章で各1回で、合計12～13回程度の講義が組めるかと思います。線形代数の章（第0章）に2～3回かけて15回としてもよいですが、第1章～第9章をやりながら適宜戻ってくるという進め方でもよいかと思います。かなりできる学生なら、最初に100問を課してもよいでしょう。本文をじっくり読めば、回答できるようになっています。輪講（担当を決めて交代で読み続けていく）でも、同じ程度の回数で読めるかと思います。

謝辞

機械学習の数理100問の執筆をご提案いただいた共立出版の皆様、特に本書の担当編集者で、シリーズ化を受け入れ、本書の出版に関して細かい点までチェックしていただいた大谷早紀氏に感謝します。また、大阪大学大学院の学生である稲岡雄介君には、提出前の原稿に目を通して、プログラムや数式の誤りを指摘してもらいました。2018～2019年の講義の学生（大阪大学基礎工学部）からのフィードバックが、本書にとって有益な情報となったことも、特記すべきと思っています。

目 次

第0章 線形代数	1
0.1 逆行列	1
0.2 行列式	3
0.3 一次独立性	5
0.4 ベクトル空間とその次元	7
0.5 固有値と固有ベクトル	9
0.6 正規直交基底と直交行列	10
0.7 対称行列の対角化	11
付録 命題 2,4,6 の証明	12
第1章 線形回帰	15
1.1 最小二乗法	15
1.2 重回帰	18
1.3 $\hat{\beta}$ の分布	20
1.4 RSS の分布	21
1.5 $\hat{\beta}_j \neq 0$ の仮説検定	22
1.6 決定係数と共線形性の検出	28
1.7 信頼区間と予測区間	30
付録 命題 12,13 の証明	33
問題 1~18	35
第2章 分類	43
2.1 ロジスティック回帰	43
2.2 Newton-Raphson 法の適用	45
2.3 線形判別と二次判別	49
2.4 K 近傍法	51
2.5 ROC 曲線	53
問題 19~31	55
第3章 リサンプリング	61

x

目 次

3.1	クロスバリデーション	61
3.2	線形回帰の場合の公式	65
3.3	ブートストラップ	68
付録	命題 15,16 の証明	72

問題 32～39 74

第 4 章 情報量基準 79

4.1	情報量基準	79
4.2	有効推定量と Fisher 情報量行列	83
4.3	Kullback-Leibler 情報量	85
4.4	赤池の情報量基準 (AIC) の導出	87
付録	命題 18,19 の証明	88

問題 40～48 91

第 5 章 正則化 95

5.1	Ridge	95
5.2	劣微分	97
5.3	Lasso	100
5.4	Ridge と Lasso を比較して	102
5.5	λ の値の設定	104

問題 49～56 106

第 6 章 非線形回帰 109

6.1	多項式回帰	109
6.2	スプライン回帰	112
6.3	自然なスプライン関数への回帰	114
6.4	平滑化スプライン	118
6.5	局所回帰	121
6.6	一般化加法モデル	125
付録	命題 20, 21, 22 の証明	126

問題 57～68 131

第 7 章 決定木 139

7.1	回帰の決定木	139
7.2	分類の決定木	147
7.3	バギング	150
7.4	ランダムフォレスト	151

7.5 ブースティング 154

問題 69～74 157

第 8 章 サポートベクトルマシン 161

8.1 最適な境界 161

8.2 最適化の理論 163

8.3 サポートベクトルマシンの解 166

8.4 カーネルを用いたサポートベクトルマシンの拡張 169

付録 命題の証明 174

問題 75～87 176

第 9 章 教師なし学習 183

9.1 K -means クラスタリング 183

9.2 階層的クラスタリング 187

9.3 主成分分析 193

付録 樹形図のプログラム 198

問題 88～100 200

索 引 205

以下では、実数全体を \mathbb{R} とかき（複素数全体は \mathbb{C} ）、 $n \times m$ の実数成分の行列の集合を $\mathbb{R}^{n \times m}$ 、 $n \times 1$ の実数成分の行列（列ベクトル）の集合を \mathbb{R}^n と書くものとする。また、行列やベクトルの転置は、 A^T, b^T のように右上に T を上付きで表記する。