

無断使用をお断りします。日科技連出版社

統計的機械学習 ことはじめ

データ分析のセンスを磨く
ケーススタディと数値例

廣野元久 [著]



日科技連

まえがき

最近のデータ分析の進歩はめざましいものがあり、AI(人工知能)や機械学習・ディープラーニングといった言葉が踊っている。ここでいう進歩とは、「実務への応用が指数関数的な速度で展開されている」という意味である。ビジネスの世界で市民権を得た機械学習であるが、中身を知りたいければ理屈を学ぶ必要がある。それには数理的な素養が不可欠であり、演算プロセスをトレースできる計算環境が必要となる。成書では演算プロセスの説明にPythonやRなどのコマンドを使った短いプログラムが示されているものが多い。例題のプログラムを参考にしてキーボードを打てば何かしらの計算結果が返ってくる。それが正しい処理かどうかは学習しないとわからない。正しい処理であっても実務に役立つとは限らない。機械に何を実現させるかは対峙している問題に対する深い洞察と専門性が需要で、データの分析ではどのような方法を選択すれば良いかを見極めることが大切である。このため、問題を解決しようとする専門家とデータ分析の専門家、加えてソフトウェア技術者の協力が欠かせない。

機械学習では主にビッグデータを対象とするが、ビッグデータとは単にデータ量が多い(Volume)だけでなく、動画や会話・文章などさまざまな種類・形式(Variety)が含まれる非構造化データ・非定型的データを対象にしたものであり、さらに、日々膨大に生成・記録される時系列性・リアルタイム性のあるようなもの(Velocity)を指すことが多い。この性質(3V)から、ビッグデータの多くは観察データである。ビッグデータを分析する目的は、今までは管理しきれなかったために見過ごされてきた情報を記録・保管して即座に分析することで、ビジネスや社会に有用な知見を得たり、これまでになような新たな仕組みやシステムを生み出したりして問題解決を行うことである。この目的を達成するための方法の一つが機械学習である。

機械学習では、最初にさまざまなデータから数値化する作業である特

特徴抽出が行われることが多い。特徴抽出には最新のソフトウェア技術が使われる。数値化されたデータは**特徴量**と呼ばれ、得られた無数の特徴量からデータ分析に必要とされる複数の特徴量を選択するプロセスは**特徴量選択**と呼ばれる。

こうして選ばれた特徴量を使ってデータ分析を行う。データ分析では目的に合わせてデータの分類や予測などが行われる。機械学習で得られたモデルはアルゴリズムと呼ばれるが、アルゴリズムや推定されたパラメータの解釈は行われず、統計モデルとは異なり、得られたアルゴリズムは人が解釈できるような式の形では求まらない。

ところで、問題解決には「原因を究明して、そこに手を打ち悪い結果を発生させない」ようにする原因追求型と「現象を深く観察して結果を予測することで、悪い結果を回避する」結果予測型の活動がある。

機械学習で扱う特徴量の多くは研究対象の状態を表す結果を数値化したものであるから、そもそも機械学習は結果に影響を与える主要な要因を探し出し、要因を制御することを目的にしている。機械学習の目的は分類や予測に役立つアルゴリズムを使って、リアルタイムに演算を行い、意思決定を支援する学習システムを構築することである。

例えば、万引き防止を考えると、万引き犯の心理を解明して万引きという行為を撲滅するのではなく、店舗に来る人の行動や表情から特徴量を見つけ、学習を繰り返して万引き犯を事前に特定して、万引きという行為を留まらせる予測システムを構築するのである。

本書は機械学習全体についての概要を説明することを目的にしている。機械学習で利用される基本的なアルゴリズムについて、統計学、特に多変量解析の延長線上にある非線形な手法であると位置づけて、その概要と使い方を数値例やケーススタディを使って説明するものである。

執筆にあたり、読者層をQC検定(品質管理検定)で2級合格を目指す技術者(実務経験があり統計学の基礎知識を有している)とした。その理由は、次世代のものづくりの中核をなす方々がもつ疑問、例えば、「機械学習で何ができるのか」「どのようなことが伝統的な統計的方法に比べて新しいのか」「統計的方法が苦手とする問題にどのように適用でき

るのか」といった点に答えることが必要と考えたからである。多くの図表(カラーの図は日科技連出版社のウェブサイトからダウンロードできる)を使って、読者が以下の①～③を達成できるよう、配慮した。

- ① 17の数値例を使って機械学習の前提条件や考え方の基本を理解できる(表1)
- ② 32のケーススタディを通じて機械学習が何を叶えてくれるのかを擬似体験できる(本書ではデータファイルを《……》で、特微量(変数)を『……』で表している)(表1)
- ③ 統計的方法と機械学習の両者を比較することで各方法の特色を明らかにする

ただし、本書を執筆するにあたり意識したことが2つある。1つは「PythonやRなどを使った計算例を示さない(成書やウェブサイトの記事に譲る)」ことであり、もう1つは「詳しい数理的説明をしない(研究者の成書や論文に譲る)」ことである。これらは機械学習を提供する研究者の領域であり、すでに良質な文献が多く出回っている(インターネット上には無料で読めるものさえある)からである。なお、PythonやRなどの使い手であれば、日科技連出版社のウェブサイトからケーススタディのデータ(Excelファイル)をダウンロードしてデータ分析を楽しむことができる。

本書ではSAS InstituteのJMP Pro 15と日本科学技術研修所のJUSE-StatWorksV5 機械学習編の2つを使ってデータ分析を行っている。手法名やハイパーパラメータの設定、扱える手法は2つのソフトウェアで異なるところがある(表2)。詳しくはソフトウェア開発先が主催するセミナーやウェブページを参照してほしい。

最後に、早稲田大学の小島隆矢先生には草稿の段階から多くのアイデアを頂戴した。日本科学技術研修所の機械学習セミナー講師をされている立教大学の山口和範先生、青山学院大学の西垣貴央先生の貴重な資料からは多くを勉強させていただき、アドバイスも数多く頂戴した。友人であるデンソーの吉野陸氏と産業技術総合研究所の遠藤幸一氏には内容を詳細に添削していただき、多くの助言をいただいた。また、日本科学

表 1 本書のケーススタディ (32 個) と数値例 (17 個) の一覧表

章	題名	章	題名
1	数値例①：ヒストグラムに潜むクラスの発見	6	数値例⑩：判別分析の弱点
	数値例②：正規混合分布を使った 2 クラス分類		数値例⑪：必要な個体だけを使った分類
	ケーススタディ①：誘発磁場の分布		ケーススタディ⑫：鮭と鱈の線形 SVM
	数値例③：時系列データに潜むパターンの発見		数値例⑬：2 つの特徴量の高次元化
	ケーススタディ②：誘発磁場の傾向		ケーススタディ⑭：鮭と鱈の非線形 SVM
	数値例④：散布図に潜むクラスの発見		ケーススタディ⑮：誘発磁場の非線形 SVM
2	ケーススタディ③：誘発磁場の等高線図	7	ケーススタディ⑯：統計的な見方・考え方
	ケーススタディ④：投手成績の等高線図		ケーススタディ⑰：良否判定を予測するロジスティック判別
	ケーススタディ⑤：投手成績のカラーマップ		ケーススタディ⑱：ROC 曲線を使ったカットオフ値の算出
	ケーススタディ⑥：投手成績の相関係数のカラーマップ		数値例⑬：2 次ロジスティック判別
	ケーススタディ⑦：誘発磁場のカラーマップ		ケーススタディ⑲：原油成分の隠れ層とノード数の設定
	ケーススタディ⑧：杉花粉の重回帰分析		数値例⑭：グラフを見ればわかるニューロ判別境界
3	ケーススタディ⑨：投手成績のホールアウト検証	8	ケーススタディ⑳：誘発磁場のニューロ判別
	ケーススタディ⑩：杉花粉データの正則化		ケーススタディ㉑：プリンタのローラ径のニューロ判別
	数値例⑤：特徴量の次元圧縮		ケーススタディ㉒：カラー画像のニューロ判別
	数値例⑥：カーネル主成分によるクラス発見		数値例⑮：簡単な予測
	数値例⑦：重要なカーネル主成分の探索		ケーススタディ㉓：住宅価格予測のニューラルネットワーク
	ケーススタディ⑪：誘発磁場の主成分分析		ケーススタディ㉔：時系列データを予測するニューラルネットワーク
4	数値例⑧：非階層クラスター分析によるクラス発見	9	数値例⑯：木モデルを使った判別境界
	ケーススタディ⑫：誘発磁場計測の k -平均法		数値例⑰：決定木を使った予測
	数値例⑨：出現可能性を使ったクラス発見		ケーススタディ㉕：糖尿病患者の症状進行予測
	ケーススタディ⑬：死因の分類		ケーススタディ㉖：クレジットリスクを求めるランダムフォレスト
	ケーススタディ⑭：鮭と鱈 (すずぎ) の判別分析		ケーススタディ㉗：アルゴリズムのコンテスト
	ケーススタディ⑮：鮭と鱈の重判別分析		ケーススタディ㉘：アンサンブル学習による体脂肪率の予測
5	ケーススタディ⑯：金属部品の重判別分析	10	——

表2 本書の表記と使用したソフトソフトウェアの表記の対応表

章	本書の表記	JMPの表記	JUSE-StatWorksの表記
1	正規混合分布	正規混合	密度プロット
	$AICc$	—	—
	カーネル平滑化	カーネル平滑化	—
	等高線図	密度等高線図	等高線図
	メッシュプロット カラーマップ	メッシュプロット カラーマップ	— —
2	重回帰分析	標準最小2乗 ステップワイズ法	重回帰分析
	予測判定グラフ	予測値と実測値のプロット	予測判定グラフ
	評価(用)データ	テスト(用)データ	評価(用)データ
3	正規化回帰分析	一般化回帰	正規化回帰
	主成分分析	主成分分析	主成分分析
4	カーネル主成分分析	—	カーネル主成分分析
	k -平均法	k means クラスタ分析	k -means 法 (マハラノビス汎距離)
5	正規混合法	正規混合	混合ガウス分布
	判別分析	判別分析	判別分析
6	サポートベクターマシン	サポートベクトルマシン	サポートベクターマシン
7	ロジスティック判別	名義ロジスティック	ロジスティック回帰分析
	混同行列	混同行列	誤判別表
	的中率	—	再現率
8	感度/特異度	混同行列	適合率
	ニューロ判別	ニューラル	—
9	ニューラルネットワーク	ニューラル	—
10	決定木分析	パーティション ¹⁾ (G^2 を表示)	多段層別分析
	ランダムフォレスト	ブートストラップ森	ランダムフォレスト (誤分類率を表示)

1) パーティションは厳密には CHAID のアルゴリズムではない。

技術研究所の犬伏秀生氏と SAS Institute Japan の岡田雅一氏にはソフトウェアを提供する立場からソフトウェアの使い方などを指導していただいた。ここでは名前を挙げなかったが、筆者と縁ある多くの方々の叱咤激励がなければ到底、完成に辿り着くことはできなかったと思う。心より感謝を申し上げたい。

出版に際しては、予定より大幅に遅れても励まし続けていただいた日科技連出版社の田中延志氏にこの場を借りて御礼を申し上げたい。また、PC を前に渋い顔をして原稿を書いている筆者を終始和ませてくれた妻、峰子にはいつもながら頭の下がる思いである。

本書は伝統的な統計的データ分析の視点から対岸にある機械学習を眺めた内容になっているが、統計的品質管理と機械学習の架け橋になればという思いとともに読者が機械学習の知識を知恵に変えて、ビジネスに活用していただく一助となれば幸いである。

2021年8月 コロナ禍の中だからこそ、愛で世界が満たされますように

廣野 元久

目 次

まえがき iii

第 1 章	ビッグデータの可視化	1
1.1	統計学と機械学習	1
1.2	データの可視化	5
1.3	正規混合分布	6
1.4	カーネル平滑化	14
1.5	(密度)等高線図	19
1.6	カラーマップ(ヒートマップ)	26
第 2 章	モデル検証	31
2.1	重回帰分析	31
2.2	クロスバリデーション	37
2.3	正則化回帰分析	44
第 3 章	カーネル主成分分析	49
3.1	カーネル法と主成分分析	49
3.2	カーネル主成分分析の基礎	53
第 4 章	クラスター分析	61
4.1	k -平均(k -Means)法	61
4.2	正規混合法(混合ガウス法)	68
第 5 章	判別分析	75
5.1	判別分析の基礎	75
5.2	判別関数	81
5.3	判別関数と外れ値	87

第6章	サポートベクターマシン	91
6.1	線形SVMの基礎	91
6.2	非線形SVMの基礎	101
第7章	ロジスティック判別分析	109
7.1	独立性の検定	109
7.2	ロジスティック判別分析の基礎	115
7.3	ロジスティック判別分析の適用	119
第8章	ニューロ判別分析	125
8.1	ニューロ判別分析の基礎	125
8.2	ニューロ判別分析の適用	137
第9章	ニューラルネットワーク	151
9.1	ニューラルネットワークの基礎	151
9.2	ニューラルネットワークの適用	159
第10章	ランダムフォレスト	171
10.1	決定木分析	171
10.2	ランダムフォレストの基礎	183
参考文献		197
索引		199

コラム一覧

コラム 1	機械学習に標本抽出の考え方は必要か	4
コラム 2	ヒストグラムは級の数と幅で形が変わる	8
コラム 3	$AICc$ とは何か	13
コラム 4	時系列の分析はデータ範囲で見た目が変わる	19
コラム 5	データ分析の前処理に手を抜かない	30
コラム 6	k -平均法の盲点	67
コラム 7	初期値で結果が変わる	72



第6章 サポートベクターマシン

判別の目的は境界を決めることでもあるから、境界付近にある点を重視し、境界から遠い点を無視するほうが都合のよい場合がある。その要求に答えてくれるのが本章で紹介する**サポートベクターマシン(SVM)**である。SVMは機械学習の主要なアルゴリズムの1つで、線形な方法とカーネル法を用いる非線形な方法とがある。

6.1 線形 SVM の基礎

判別分析ではすべての個体を使って判別境界を求めている。その際、観測値は多変量正規分布に従う確率変数であること、また、各群の母分散と母共分散は等しいという前提にもとづく。この前提はある意味で行儀が良すぎるものである。**線形 SVM**は判別に役立つ個体を数理的に選んで、少数の観測点(**サポートベクター**)で判別境界を定める方法である。

6.1.1 【数値例⑩：判別分析の弱点】

データの背後に正規性・等分散共分散性といった統計仮説が成り立っている条件で判別分析が行われる。もし、その仮説が崩れた場合にはどのようなことが起きるのかを数値例で確認する。図 6.1-①は《判別分析の弱点》のデータファイルから『 x 』と『 y 』を『クラス』で分けた確率楕円を追記した層別散布図である。『クラス』は C1 がマーカーの△で、C2 がマーカーの○で識別されている。2つの確率の様子(長軸の方向や楕円の短軸方向の幅)から、「2つのクラスに等分散共分散性が成り立っている」とグラフからは感じられない。統計仮説が崩れている状況で判別境界を求めると何が起ころうか。判別分析で得られた判別境界を追記したグラフを図 6.1-②に示す。C1 を C2 と誤分類することはないが、C2 を C1 と誤分類することが多い結果となった。グラフ

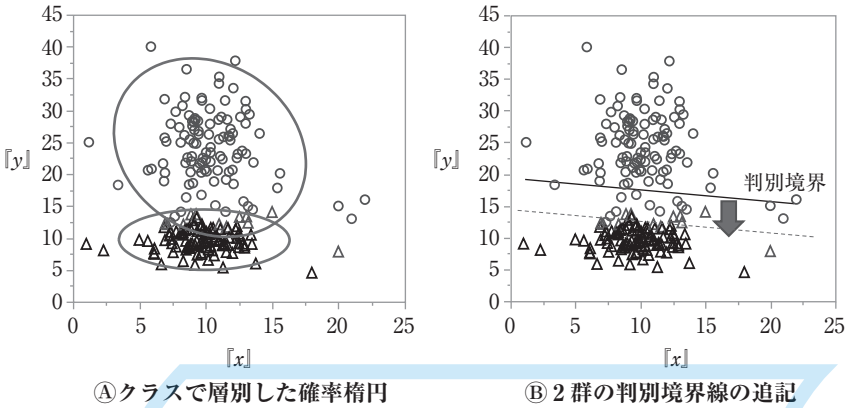


図 6.1 《判別分析の弱点》の $[x]$ と $[y]$ での判別 (\triangle はC1, \circ はC2を表す)

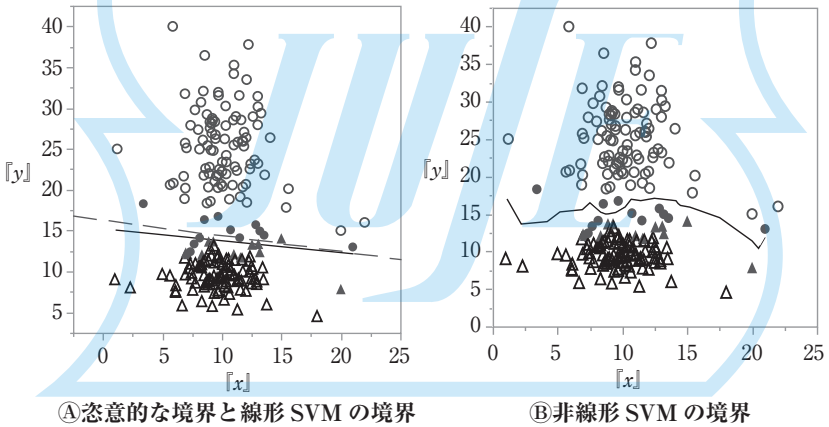


図 6.2 線形判別分析とは異なる方法で求めた判別境界

から直感的に、「判別境界線を垂直方向にもっと下げたほうがよい」と感じるであろう。

そこで、主観的に境界に近い観測点だけを使って判別分析を行い、新たな境界線を引くことを考える。図 6.2-①の実線は●と▲以外の観測点の重みを0にした場合の判別境界である。すべての観測点を使った判別境界線よりもよい判別結果が得られている。しかし、このような恣意的な方法は分析者の主観が色濃く反映されるので、客観性に乏しく好ま

しい方法とはいえない。

ところで、図中の破線は線形 SVM により求めた判別境界線である。こちらの直線は客観的な基準で選定された観測点を使って得られたものである。また、図 6.2-⑧の複雑な境界は非線形 SVM により求めた判別境界線である。非線形 SVM のほうがよい判別ができていのように感じられるであろう。以降では、判別境界を求めるためにどのような基準で観測点が選ばれるのか、および、SVM の概念を説明する。

6.1.2 【数値例⑪：必要な個体だけを使った分類】

SVM は広く使われている学習アルゴリズムの 1 つである。判別分析は、すべての個体の情報を使って判別境界を計算するが、SVM では分類に必要な個体の情報(分類をサポートするベクトルとして)を使って判別境界を計算する。また、カーネル法と組み合わせることで、非線形な判別を行うことができ、高い判別性能が得られる。

図 6.3 は《SVM》の『X』『Y』に線形な判別モデルをあてはめた結果を示したものである。SVM では判別境界をサポートするベクトルは境界に面した個体を利用する。《SVM》では、クラス 1 から 2 つ(◆のマーカ)とクラス 2 から 1 つ(◇のマーカ)を使って判別境界を作る。

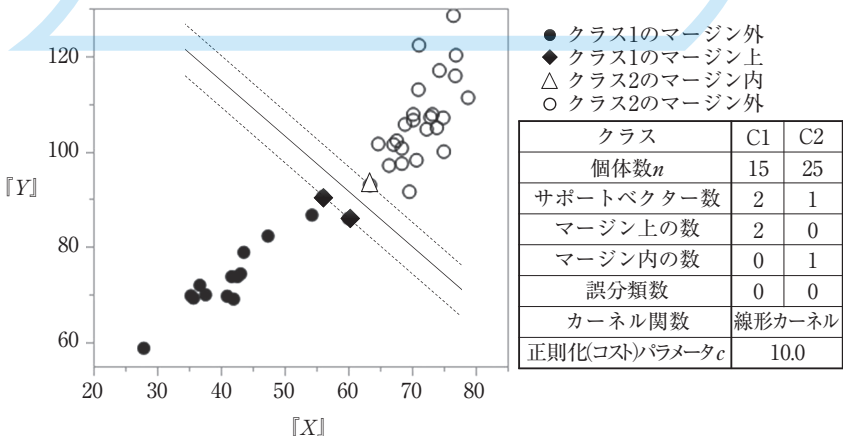


図 6.3 《SVM》の線形な判別境界

図中の実線は判別境界を、破線はマージンを加えた境界線を表している。

6.1.3 SVMの特徴

判別分析もSVMも判別関数の符号で、どちらのクラスに所属するかを決定する。判別分析では図6.4-①のように正規分布を仮定しており、マハラノビス距離を基準に分類する。判別境界は各クラスの平均からのマハラノビス距離が等しい点になる。SVMでは図6.4-②のように判別境界は各クラスの1番近い個体からできるだけ距離をとるように作られる。これをマージン最大化という。この判別境界線に1番近い観測点をサポートベクターという。この名前の由来は、観測点が境界を支える(サポートしている)ベクトルであることから来ている。図6.4はサポートベクターの数が2つの例を示している。一般に、サポートベクターの数が少ないと小さな次元で分類できており、判別境界は単純である。一方、サポートベクターの数が多いと誤分類の数が見た目で少なくなるが、高い次元で分類されているので判別境界が複雑になる。サポートベクターの数が必要以上に多い場合は、モデルが過学習している可能性が高く、汎化性能が劣っている可能性がある。また、SVMはモデルの再構築が容易であることが知られている。それは判別境界に近い個体と新しい個体だけを使ってモデルの再構築を行うからである。

ところで、SVMにも弱点がある。それは、計算時間が個体数 n に大きく依存することである。後述するカーネルトリックによって特徴量の

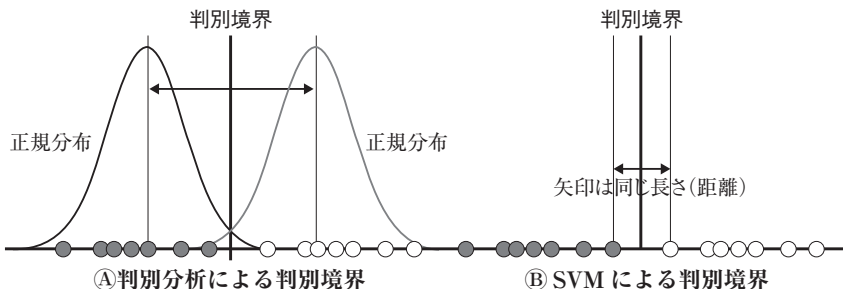


図6.4 2つの判別境界の考え方の比較

●著者紹介

廣野元久 (ひろの もとひさ)

1984年、㈱リコー入社。以来、社内の品質マネジメント・信頼性管理の業務、SQCの啓蒙普及に従事、品質本部QM推進室長、NA事業部SF事業センター所長を経て、現在、㈱リコー 倫理審査委員会委員。

東京理科大学工学部経営工学科 非常勤講師(1997～1998年)、慶應義塾大学総合政策学部 非常勤講師(2000～2004年)。

主な専門分野はSQC、信頼性工学。主著に『グラフィカルモデリングの実際』(共著、日科技連出版社、1999年)、『JMPによる多変量データ活用術』(海文堂出版、2004年)、『SEM因果分析入門』(共著、日科技連出版社、2011年)、『アンスコムのな数値例で学ぶ統計的方法23講』(共著、日科技連出版社、2013年)、『目からウロコの統計学』(日科技連出版社、2017年)、『JMPによる技術者のための多変量解析』(日本規格協会、2018年)、『目からウロコの多変量解析』(日科技連出版社、2019年)など。

統計的機械学習ことはじめ

データ分析のセンスを磨くケーススタディと数値例

2021年9月28日 第1刷発行

著者 廣野元久
発行人 戸羽節文

検 印
省 略

発行所 株式会社 日科技連出版社
〒151-0051 東京都渋谷区千駄ヶ谷5-15-5
DSビル
電 話 出版 03-5379-1244
営業 03-5379-1238

Printed in Japan

印刷・製本 東港出版印刷(株)

© Motohisa Hirono 2021

ISBN 978-4-8171-9740-5

URL <https://www.juse-p.co.jp/>

本書の全部または一部を無断でコピー、スキャン、デジタル化などの複製をすることは、著作権法上での例外を除き禁じられています。本書を代行業者等の第三者に依頼してスキャンやデジタル化することは、たとえ個人や家庭内での利用でも著作権法違反です。