

● はしがき

データ分析のツールは日々進化しており、次々と新しい分析手法が経済学の実証研究に取り入れられています。特に最近の大きな変化は、機械学習の隆盛でしょう。機械学習の手法を用いた実証研究は、この数年で急速に増加しています。ところが、教育の側はそのような流行に追いついているかというところ、どうでしょうか。特定の教員によって単発的に授業が行われることはあっても、計量経済学の標準的なトピックとして機械学習の手法が扱われるまでには至っていないのではないかと思います。

もちろん、最近では機械学習関連の書籍は日本語でも数多く出版されていますので、それらを読めば済むのかもしれませんが。しかし、私の知っている範囲では、それらの本は一部を除いては、経済学者の興味とは異なるところに重点を置いていたり、あるいは、基礎的過ぎて物足りなかったりと、計量経済学を学ぶ人たちにとって丁度良い塩梅のものは、依然として多くはない印象です。

本書はこのような昨今の状況を鑑みて書かれたものです。経済学をバックグラウンドとする人たちを主たるターゲットとして、従来の計量経済学の手法の拡張という視点から機械学習の手法を解説しています。ですが、経済学の知識は特段必要としません。

本書は7つの章から成り、大きく前半（1～4章）と後半（5～7章）に分けられます。前半は伝統的な計量経済学・統計学の手法に関するパートです。ただし、伝統的とは言っても、通常の授業では詳しく取り上げられない内容に焦点を当てており、最近の研究についても触れています。1章で本書の問題意

識を述べた後、2章で線形回帰モデルの変数選択、3章でノンパラメトリック回帰、4章でセミパラメトリック回帰について解説しています。2～4章で扱う手法は、計量経済学や統計学で以前から用いられてきたデータ駆動型の分析手法です。これらはそれ自体で重要であるとともに、後半の内容の理解の助けにもなると考えて、トピックとして採用しました。

後半は機械学習の手法に関するパートです。機械学習が対象とする範囲は膨大で、私の理解が及ぶ範囲にも限りがありますので、教師あり学習と呼ばれるカテゴリーに属するものの一部のみを扱っています。教師あり学習の主たる目的は予測ですが、計量経済学では区間推定や検定などの統計的推測も重要ですので、機械学習の手法を基にした統計的推測についても解説しています。具体的には、5章で回帰木などのツリーベースの手法、6章で線形回帰モデルの正則化推定量を紹介し、7章ではそれらの手法を応用して、興味のあるパラメータの妥当な信頼区間を求める方法を考察します。

私がはじめて機械学習と接点を持ったのは2010年で、Ph.D.を取得した直後でした。その年に開催された情報論的学習理論ワークショップ (IBIS 2010) で、企画セッション「計量経済と機械学習」に登壇者として招待していただいたのがきっかけです。そのときに機械学習に対して抱いた印象は、自分の関心とかなり近いことをやっているというものでした。機械学習の主眼はデータ駆動的なモデリングにあります。元々私は変数選択やノンパラメトリック・セミパラメトリックな分析手法の研究をしていたこともあり、データ駆動的という考え方はすでに馴染みのあるものだったからです。最近では機械学習の手法を自分の研究にも取り入れるようになりましたが、私にとっての機械学習はデータ駆動型の分析を行うためのいくつかある手段のうちの一つであり、従来の計量経済学や統計学を無用にしてしまうようなものであるとは考えていません。とは言え、機械学習の有用性を過小評価しているわけでもありません。

「データ駆動型」と並ぶ本書のもうひとつのキーワードは、「一様性」です。統計分析の手法が満たすべき望ましい性質のひとつとして、分析者が想定するどのような分布からデータが発生したとしても、その手法がうまく機能するというものがあります。そのような性質は一様妥当性 (uniformly validity) と言われ、近年の計量経済学の理論研究では重要視されている性質なのですが、大

抵の教科書では触れられていないのではないかと思います。理論の詳細までは踏み込みませんが、雰囲気だけでも伝わればと思い、本書のいくつかの箇所での問題を取り上げています。

本書では学部上級レベルの計量経済学の知識を前提としています。『計量経済学』（西山慶彦他、有斐閣、2019年）くらいを想定していますが、それにプラスして、線形代数の知識も多少必要です。一部ではそれだけでは理解が困難かもしれない発展的な内容も扱っていますが、それらは主として理論研究に関心がある人やプロの研究者向けに書いています。星印（*）の付いている節は難易度が高めで、星印が2つ（**）付いている節はかなり難しいです（星印が付いていない部分が決して易しいわけでもないですが）。星印の有無にかかわらず、難しそうな部分は適宜読み飛ばしてもらってもよいでしょう。

本書のスタイルは、基本的には拙著『計量経済学：マイクロデータ分析へのいざない』（日本評論社、2015年）のそれを踏襲しています。証明はほぼ省略し、興味のある人には参考文献を挙げるにとどめています。その理由は、話の全体像が見えづらくなることを避けるためであり、厳密性を犠牲にしているというわけではありません。証明をしない代わりに、説明は丁寧にしています。新しく登場する事柄については、可能な限り数式と言葉の両方を用いて解説をするようにしています。また、どのような分析手法であれ、特定の目的には有用であっても、その他の目的にとっては有害にもなりうるので、分析手法の優れた点だけでなく、問題となりうる点についても極力言及するようにしています。

本書の執筆にあたり、多くの方々にお世話になりました。東京大学の奥井亮氏と坂口翔政氏、慶應義塾大学の岡達志氏、滋賀大学の新久章氏、ニューヨーク大学博士課程の池上慧氏、カリフォルニア大学サンディエゴ校博士課程の林田光平氏には、本書の草稿に目を通していただき、多くの有益なコメントをいただきました。また、本書は私のサバティカル中に執筆されたものです。サバティカル中に一緒に仕事をさせていただいた株式会社サイバーエージェント AI Lab の安井翔太氏、竹浪良寛氏、松木一永氏からも有益なコメントをいただきました。紙幅の都合上、あるいは、私の能力の問題で、いただいたコメントのすべてを最終稿に反映させることは叶いませんでしたが、貴重な時間を割

いて原稿を読んでいただいた皆様に心から感謝申し上げます。

また、前著に続いて執筆の機会を与えてくださった日本評論社の吉田素規氏には、今回も大変お世話になりました。原稿の数々の間違いを吉田氏に発見していただかなければ、大変なことになっていたでしょう。もちろん、まだありうる本書の間違いについては、すべて私の責任です。

回帰分析と言えば、佐和隆光先生の『回帰分析』（朝倉書店）が有名です。初版が出版された1979年は、私が生まれた年でもあります。それから45年が経過し、見識は佐和先生にはるかに及ばないものの、私なりに現代の回帰分析のあり方を模索し苦闘した跡を、本書に読み取っていただければ幸いです。

2024年3月

末石直也

● 目 次

はしがき i

第1章 回帰分析の課題 1

- 1.1 回帰分析 1
 - 1.1.1 構造モデルと回帰モデル 1
 - 1.1.2 回帰分析の目的 4
- 1.2 線形回帰モデル 7
 - 1.2.1 線形回帰モデル 7
 - 1.2.2 OLS 推定量とその性質 8
 - 1.2.3 線形射影モデル 11
 - 1.2.4 欠落変数バイアス 13
- 1.3 本書の課題と構成 14
- 1.4 補論 17

第2章 変数選択 23

- 2.1 設定 24
- 2.2 推定された予測誤差に基づく変数選択 27
 - 2.2.1 Mallows の C_p と Stein の SURE 28
 - 2.2.2 交差検証法 30
- 2.3 情報量規準 33
 - 2.3.1 一般のモデルの AIC と BIC 33
 - 2.3.2 線形回帰モデルの AIC と BIC 35
- 2.4 変数選択の一致性と漸近最適性 37
- 2.5 その他のモデル評価基準 39
- 2.6 変数選択後の統計的推測の問題 40
 - 2.6.1 被覆確率に与える影響 41
 - 2.6.2 一致性を満たす変数選択法の問題点 42

第3章 ノンパラメトリック回帰 45

- 3.1 カーネル推定 46
 - 3.1.1 Nadaraya-Watson 推定量と局所線形推定量 46
 - 3.1.2 推定量の性質とバンド幅の選択 49
 - 3.1.3 次元の呪い 54
 - 3.1.4 サポートの境界における回帰関数の推定 55
 - 3.1.5 信頼区間 57
- 3.2 シリーズ推定 59
 - 3.2.1 推定方法 59
 - 3.2.2 推定量の性質とシリーズの長さの選択 61
 - 3.2.3 信頼区間 63
- 3.3 回帰不連続デザイン 64
 - 3.3.1 識別 64
 - 3.3.2 推定 66
- 3.4 補論 70

第4章 セミパラメトリック回帰 73

- 4.1 部分線形モデル 74
 - 4.1.1 モデル 74
 - 4.1.2 Robinson (1988) の推定量 75
- 4.2 シングルインデックスモデル 76
 - 4.2.1 モデル 76
 - 4.2.2 Ichimura (1993) の推定量 77
- 4.3 平均処置効果 79
 - 4.3.1 強い無視可能性の仮定の下での識別 79
 - 4.3.2 推定 80
- 4.4 プラグイン推定量の漸近分布* 81
 - 4.4.1 プラグイン推定量 81
 - 4.4.2 η_0 の推定が $\hat{\theta}$ の漸近分散に影響を与えない場合 83
 - 4.4.3 η_0 の推定が $\hat{\theta}$ の漸近分散に影響を与える場合 85
 - 4.4.4 η_0 が回帰関数の場合 89
- 4.5 補論** 91
 - 4.5.1 正則な推定量と微分可能なパラメータ 91
 - 4.5.2 なぜ正則な推定量か 94

第5章 回帰木とアンサンブル学習 99

- 5.1 決定木 100
 - 5.1.1 分類木 100
 - 5.1.2 回帰木 104
- 5.2 バギングとランダムフォレスト 107
- 5.3 ブースティング 110
 - 5.3.1 ブースティングのアイデア 111
 - 5.3.2 勾配ブースティング 112
 - 5.3.3 勾配ブースティング木 113
- 5.4 因果木と因果フォレスト* 115
- 5.5 補論* 119

第6章 正則化法 121

- 6.1 リッジ回帰 122
 - 6.1.1 リッジ推定量の解釈 122
 - 6.1.2 正則化パラメータの選択 125
- 6.2 Lasso 回帰 127
 - 6.2.1 スパース性 127
 - 6.2.2 最小角回帰と Lasso 129
 - 6.2.3 Lasso の性質 133
 - 6.2.4 正則化パラメータの選択 135
- 6.3 エラスティックネット 136
- 6.4 オラクル性を満たす正則化法 137
- 6.5 オーバーフィッティングは本当に問題か* 140
- 6.6 補論* 143

第7章 変数選択後の統計的推測 147

- 7.1 一様に妥当な信頼区間 148
- 7.2 debiased Lasso 150
- 7.3 post-double-selection 154
- 7.4 double/debiased machine learning* 156
 - 7.4.1 DML 推定量とその性質 157
 - 7.4.2 Neyman 直交化 160
 - 7.4.3 Neyman 直交性が満たされないときの問題 162

7.5	選択的推測*	163
7.5.1	PoSI 信頼区間	164
7.5.2	正確な被覆確率を持つ信頼区間	167
	参考文献	171