

まえがき

Rは統計計算、データ分析、そしてその結果を可視化するために設計されたプログラミング言語であり、近年、データサイエンスおよび統計学の分野においては最も有名なものとなっている。Rによるプログラミングにはデータプロセッシングが大きな比重を占めており、これはRに不慣れなユーザにとってはかなり骨の折れる作業でもある。

Rは動的言語 (dynamic language) であり、C++ や Java, C# といった厳密な型をもつ言語に比べると非常に柔軟なデータ構造を有する。実際、そのような言語に慣れていた筆者がRを利用し始めた時には、Rの挙動は非常に奇妙かつ予測不能で、一貫性のないものに思えたものである。Rを利用するようなデータ分析プロジェクトにおいて、多くの時間はモデルの作成ではなく、データクリーニングや探索的分析、その可視化に割かれる。そして、おかしな結果やエラーを吐き出すコードの間違い探しも、我々から多くの時間を奪う。本質的な問題解決とは程遠いプログラミングそのものの問題と戦っている時、とりわけヒントなしに何時間もバグと戦っている時などはイライラして仕方がない。しかし、データ分析プロジェクトにおける経験を積み、Rにおけるオブジェクトや関数の挙動を理解するにつれ、思ったよりもこの言語は一貫性があり美しく設計されているのではないかと思うようになった。この観点を多くの人と共有したい、これがこの本を執筆するに至った動機である。

本書を通じて、あなたはRをプログラミング言語として広く一貫性をもった形で理解でき、周辺の有用なツールについても知ることができる。また、あなたの生産性を向上させるベストプラクティスやデータ操作についての深い理解も得られるだろう。そして、Rを用いてプログラミングを行い、適切な問題解決を実行するために十分な自信も獲得できるものと思われる。

本書がカバーする内容

第1章「クイックスタート」では、Rに関する基本的な知識やRの環境構築方法、そしてRStudioを用いたプログラミングについて学ぶ。第2章「基本的なオブジェクト」では、Rにおける基本的なオブジェクトの内容および挙動について学ぶ。第3章「作業スペースの管理」では、作業ディレクトリやR内の環境、そして拡張パッケージの管理方法について学ぶ。第4章「基本的な表現式」では、代入式、条件式、ループといった基本的な条件式について学ぶ。第

iv まえがき

5章「基本的なオブジェクトを扱う」では、Rにおけるオブジェクトを操作するための基本的な関数について学ぶ。第6章「文字列を扱う」では、文字列に関連するRのオブジェクトおよび文字列操作のテクニックについて学ぶ。第7章「データを扱う」では、実例を交えた形でシンプルなる入出力関数について学ぶ。第8章「Rの内部を覗く」では、遅延評価や環境、関数、レキシカルスコーピングといったトピックを取り上げながら、Rにおける評価の仕組みについて学ぶ。第9章「メタプログラミング」では、Rの言語オブジェクト (language object) や非標準評価 (non-standard evaluation) を理解する際に役立つメタプログラミングについて学ぶ。第10章「オブジェクト指向プログラミング」では、S3, S4, 参照クラス, 拡張パッケージとして提供されているR6といったRにおけるクラスシステムについて学ぶ。第11章「データベース操作」では、SQLiteやMySQLといったポピュラーなリレーショナルデータベースをRでどう扱うかに始まり、MongoDBやRedisといったNoSQL型データベースの操作方法についても学ぶ。第12章「データ操作」では、`data.table` パッケージおよび `dplyr` パッケージを用いたデータフレームの操作、そして `rlist` パッケージを用いたリストの操作について学ぶ。第13章「ハイパフォーマンスコンピューティング」では、パフォーマンスの話題を扱い、`Rcpp` パッケージの利用といったRにおける計算速度を高める方法について紹介する。第14章「ウェブスクレイピング」では、ウェブページ, CSSセレクタ, XPathの基礎および `rvest` パッケージを用いたウェブスクレイピングについて学ぶ。第15章「生産性を高める」では、分析結果のレポートやプレゼンテーションを効率よく実行できるRMarkdown, Shinyについて学ぶ。

本書のコードを実行するために必要な環境

本書のコードはR3.4.0日本語版で実行を確認している¹⁾。また、開発環境はRStudioが望ましい。

第11章のコードを実行するためにはMongoDB, Redisの実行環境が必要である。第13章におけるRcppのコードを実行する際、Windowsの場合は利用しているRのバージョンに合わせたRtoolsのインストールが、Linux/macOSの場合はgccのインストールがそれぞれ必要である。

本書の想定読者

本書の読者としては、データ分析プロジェクトにかかわっており生産性を上げたいと思っているものの、プログラミングには不慣れな人を想定している。また同時に、周辺技術やツール、拡張パッケージも交えた形でRを言語として体系立てて理解したい上級者も想定読者に含まれる。

いくつかの章は初心者にとっては難しいため、読み飛ばしてもらっても構わない。だが、これらの章を読んでおくと、Rに限らないプログラミングの基本概念やデータ分析における基本

¹⁾ 訳注：エラーメッセージ、警告メッセージについては和英が混在している形になっているが、これは本実行環境を反映したものとなっている。

概念を知ることができるので、上級者を志すなら一読をお勧めしたい。

本書の表記について

本書では用途に合わせて表記を変えている。データベースのテーブル名、フォルダ（ディレクトリ）名、ファイル拡張子、Twitter のハンドル名は、コードと同様の書体で表記している。また、関数名については `apply()` のように末尾に括弧を付与している。以下にその表記例を示す。

例：`apply()` は配列による入力を受け付け、行列の形で出力する。

本文中のコードは以下のようにインデントを加えて表記しており、出力結果については#を2つ並べた形で区別している。

```
x <- c(1, 2, 3)
class(x)
## [1] "numeric"
typeof(x)
## [1] "double"
str(x)
## num [1:3] 1 2 3
```

一部のコードでは、重要なポイントを表すために太字で表記している。

```
x <- rnorm(100)
y <- 2 * x + rnorm(100) * 0.5
m <- lm(y ~ x)
coef(m)
```

初出の単語および重要な用語についても太字で表記している。

²⁾ 訳注：本書の正誤表およびコードに関する補足事項は、GitHub 上のサポートサイトに掲載している。
https://github.com/HOXOMInc/Learning_R_Programming