

データサイエンスをはじめましょう - Data Science for All -

John Doe

2023-02-18

目次

この文書について	4
著者について	5
コンピュータ言語について	5
言語について	5
PDF、ePub 版について	6
0.1 はじめに	6
 第 I 部 PART I PUBLIC DATA	 7
0.2 Public Data	9
0.2.1 World Bank	9
0.2.2 Worldbank Data	9
0.2.3 World Bank: WDI - World Development Indicators	9
0.2.4 OECD	10
0.2.5 UN Data	10
0.2.6 Our World in Data	10
0.2.7 Eurostat	10
 第 II 部 PART II BASICS	 11
0.3 R on R Studio	13
0.3.1 はじめに	13
0.3.2 R と R Studio	13
0.3.3 R と R Studio のインストール	13
0.3.4 プロジェクト - Project	14
0.3.5 コンソールで実行 - Run in Console	15
0.3.6 RStudio について	19
0.3.7 R Script 実行記録	19
0.3.8 パッケージ - Packages	28
0.3.9 クラウド - Posit Cloud	29
0.3.10 練習問題 Posit Primers	29
0.3.11 参考文献 References	30

0.3.12	YouTube Video - getstarted	30
0.4	R Markdown	30
0.4.1	Reproducible and Literate Programming	30
0.4.2	準備: パッケージのインストール	31
0.4.3	R Notebook	32
0.4.4	日本語のテンプレート	32
0.4.5	R Markdown いくつかの Output	32
0.4.6	YouTube Video - rmarkdown	33
第 III 部	PART III INSTITUTIONAL DATA	35
0.5	World Bank	37
第 IV 部	PART IV EDA	39
0.6	What is EDA?	41
第 V 部	PART V EXAMPLES	43
0.7	Example 1	45
.1	日本語の扱いについて	45
.1.1	日本語□中国語□韓国語	45
.1.2	Base R でタイトルに日本語	45
.1.3	列名や、データに日本語	46
.1.4	kable で表示	46
.1.5	ggplot でグラフを作成	46
.1.6	備考:	47
.1.7	参考: 日本語の表示について	47
.2	Bookdown	49
.2.1	About	49
.2.2	Hello bookdown	50
.2.3	Cross-references	51
.2.4	Parts	52
.2.5	Footnotes and citations	53
.2.6	Blocks	53
.2.7	Sharing your book	54

この文書について

データサイエンスを始めてみませんか。

データサイエンスは、広い意味をもったことばで、一口に、まなび始めると言っ

でも、さまざまな始め方があると思います。本書では、そのひとつを提案するとともに、共に学んでいきたいと願って、書き始めました。

みなさんも一緒にデータサイエンスを学んでみませんか。

著者について

著者は、大学の学生の時以来、数学を学び、大学で教え、2019 年春に退職。それ以来、少しずつ、データサイエンスを学んでいます。

幸運にも、2019 年 9 月の日本数学会教育委員会主催教育シンポジウムで、「文理共通して行う数理口データサイエンス教育」という題で、話す機会が与えられ、その後、あることが契機となり、2020 年度から、毎年、冬に、大学院一般向け（分野の指定なし）の授業、「研究者のためのデータ分析 (Data Analysis for Researchers)」を担当しています。複数の教員で担当しますが、基本的な部分は、わたしが教えています。受講生は 20 人程度ですが、殆どが、外国人。それも、多国籍で、多くても一国から三人程度。英語で教えています。

コンピュータ言語について

統計解析のために開発された R を使います。いずれは、python についても触れたいと思いますが、プログラミングの経験がない方も含めて、最初にデータサイエンスを学ぶには、R は最適だと考えています。特に、R Studio IDE (integrated development environment, 統合開発環境) で、R を使うことがとても、簡単になっています。さらに、簡単なものであれば、Posit Cloud で試したり、共有することも可能です。また、再現性 (Reproducibility) や、なにを実行しているのかの説明を同時に記述すること (Literate Programming) は、非常に重要ですが、その記述も、R Markdown によって、可能になっています。この文書も、R Markdown の一つの形式の、bookdown を利用しています。最後に、Bookdown に関連して、膨大な数の、参考書も、無償で提供されており、オンラインで読むことができることも、R をお勧めする理由です。

ただし、日本語のものは、まだ十分とは言えない状況です。この文書を書き始めたのも、すこしでも、お役に立つことができればとの、気持ちが背景にあります。

言語について

ご覧の通り、本書は、日本語で書かれています。用語は、英語、あるいは、英語を追記、または、英語をカタカナにしただけのものを使用する可能性が大きいですが、説明は、極力、日本語で書いていく予定です。

しかし、基本的に、コード (プログラムの記述) には、日本語を使わないで書いて

いく予定です。とくに、初心者にとっては、日本語の扱いは、負担になることが多いからです。最近では、コードの中で日本語を使用しても、ほとんど、問題は起きないように思います。そうであっても、世界の人の共通言語として、プログラム言語を学んでいくときには、日本語を使わないことは意義があると思います。

少し慣れてきて、日本語のデータなどを扱うときには、コードにも日本語を使う必要ができていますから、日本語の利用についても、追って説明していきます。.1を参照してください。

最初は、みなさんも、変数 (variable) や、オブジェクト (object) に名前をつけるときは、半角英数を使い、日本語は、使わないようにすることをお勧めします。

PDF、ePub 版について

実は、PDF 版と、ePub 版も作成しています。しかし、扱いが異なるので、ある程度完成するまでは、ほとんど更新しない予定です。いずれ、これらも、更新したものを公開できると良いのですが。試験公開版は、下のリンクにあります。

- PDF 版
- ePub 版

0.1 はじめに

Data Science: データ (Data) を活用して課題を発見□探求し、適切な解決策を探る意思決定のための科学(Decision Science)で、エンピリカル(Empirical Study) すなわち、理論ではなく、実証性を特徴とする。データから得られる特徴を表示するとともに、数理モデルを適用し□機械学習などで評価し□アルゴリズムを策定する数理的思考を通して得られた結果を、可視化などによってコミュニケーションをおこない、共有し、他者の意見を聞き理解する努力をしながら、さらに課題について、あらたにデータを活用して考え、検証し、適切な解決策がもたらす新たな課題も予測しながら、調整をはかる。

第 I 部

PART I PUBLIC DATA

0.2 Public Data

まずは、パブリックデータを見てみましょう。大きな機関のパブリックデータには、ダッシュボード (dashboard) と呼ばれている、パラメタを変更して、そのグラフを描くなどの機能が付いているものもあります。

0.2.1 World Bank

0.2.1.1 Open Government Data Toolkit: Open Data Defined

The term **Open Data** has a very precise meaning. Data or content is open if anyone is free to use, re-use or redistribute it, subject at most to measures that preserve provenance and openness.

1. The data must be *legally open*, which means they must be placed in the public domain or under liberal terms of use with minimal restrictions.
2. The data must be *technically open*, which means they must be published in electronic formats that are machine readable and non-proprietary, so that anyone can access and use the data using common, freely available software tools. Data must also be publicly available and accessible on a public server, without password or firewall restrictions. To make Open Data easier to find, most organizations create and manage Open Data catalogs.

0.2.2 Worldbank Data

- Climate Change Knowledge Portal: <https://climateknowledgeportal.worldbank.org>
 - country summary

0.2.3 World Bank: WDI - World Development Indicators

- World Bank: <https://www.worldbank.org>
- Who we are:
 - To end extreme poverty: By reducing the share of the global population that lives in extreme poverty to 3 percent by 2030.
 - To promote shared prosperity: By increasing the incomes of the poorest 40 percent of people in every country.
- World Bank Open Data: <https://data.worldbank.org>
 - Data Bank, World Development Indicators, etc.

0.2.3.1 World Development Indicator

- World Development Indicators (WDI) : the World Bank's premier compilation of cross-country comparable data on development; 1400 time series indicators
 - Themes: Poverty and Inequality, People, Environment, Economy, States and Markets, Global Links
 - Open Data & DataBank: Explore data, Query database
 - Bulk Download: Excel, CSV
 - API Documentation

0.2.4 OECD

OECD Data: <https://data.oecd.org/>

0.2.5 UN Data

UNdata: <https://data.un.org>

0.2.6 Our World in Data

owid: <https://ourworldindata.org/>

0.2.7 Eurostat

eurostat: <https://ec.europa.eu/eurostat>

第 II 部

PART II BASICS

0.3 R on R Studio

0.3.1 はじめに

R Studio で R を使うことを始めましょう。

このページの一番下に、簡単な解説ビデオがついています。

0.3.2 R と R Studio

R is a free software environment for statistical computing and graphics. It compiles and runs on a wide variety of UNIX platforms, Windows and MacOS. <https://www.r-project.org>

R は、無償で提供されている、統計解析とグラフを描写する環境です。Windows、MacOS や、Linux で利用することが可能です。

RStudio is an integrated development environment (IDE) for R and Python. It includes a console, syntax-highlighting editor that supports direct code execution, and tools for plotting, history, debugging, and workspace management. RStudio is available in open source and commercial editions and runs on the desktop (Windows, Mac, and Linux). <https://posit.co/products/open-source/rstudio/>

RStudio は、R と Python のための、総合開発環境です。RStudio には、プログラムを実行したり、制御やジョブ管理のための、コンソール (console)、コードを書いたり、実行したりする、文書の編集をする、エディター (Editor) とともに、グラフを表示したり、履歴や、プログラムを修正するなどのための、さまざまなツールが付属しています。RStudio はオープンソースで提供され、Windows、Mac および、Linux で利用可能で、有償版のサービスと無償版を提供しています。

R は、統計解析のためのシステムで、R Studio は、R (および Python) を利用するための、総合開発環境です。そこで、「R Studio で R を利用する」という表現をします。

0.3.3 R と R Studio のインストール

R と R Studio をインストールします。

両方とも、インストールする必要があります。

0.3.3.1 R のインストール

<https://cloud.r-project.org>

上のリンクから、Windows、macOS または、Linux を選択して、インストールしてください。

macOS の場合は、M1, M2 など、最近の Apple Silicon の CPU で動くコンピュータか、以前の、Intel の CPU で動くものか、選択してください。Mac の左上の、りんごマークの、このコンピュータについてから、確認できます。

不明の場合は、「R のインストール」と検索してみてください。

0.3.3.2 R Studio のインストール

<http://www.rstudio.com/download>

上のリンクから、Windows 10/11 または、macOS 11+ を選択してください。これら以外の、古いシステムのコンピュータの場合は、下のサイトから、探してください。

<https://docs.posit.co/previous-versions/>

不明の場合は、「RStudio のインストール」と検索してみてください。

0.3.4 プロジェクト - Project

RStudio で R を利用する場合には、プロジェクトを作成することを強く勧めます。

1. まず、R Studio を起動します。
2. 上のメニューの、File から、New Project を選択します。New Directory (新しいディレクトリー) を選択し、プロジェクトを作成する Directory を決めて、名前をつけます。その名前が、プロジェクト名になります。
 - Directory (フォルダー) を指定してその名前をつけて、プロジェクトを作成します。
 - Directory が階層に分かれているときは、どこに作成するかを選択してから、名前をつけて、作成します。
3. 一旦、R Studio を終了してみましよう。
4. プロジェクトの起動には、いくつかの方法があります。
 - まず、R Studio を起動。一つしかプロジェクトがない場合は、そのプロジェクトが起動すると思います。上に、プロジェクト名が掲載されていれば、

問題ありません。

- File から、Open Project を選択し、起動したい、プロジェクトの Directory (フォルダー) を選択して起動します。
- File から、Recent Project (最近使ったプロジェクト) を選択すると、プロジェクト名が表示されますから、選択すると起動することができます。
- コンピュータのプロジェクト入っているディレクトリー (フォルダー) をさがし、そこに、プロジェクト名.Rproj とあるものを見つけて、それを開くと、そのプロジェクトが起動します。

5. 作業後は、保存しますかと聞かれますから、保存して終了してください。

0.3.5 コンソールで実行 - Run in Console

プログラム (コード) の実行は、いくつかの方法がありますが、一番、基本的な、コンソール (Console) での実行にすいて、説明します。Console は、R Studio の左下にあります。(左の枠が一つになっているかもしれません。その一番左のタブが Console です。選択されていない場合は、Console を選択してください。)

0.3.5.1 最初の四つ

下の、四つを、一つずつ、一番下の、> マークの次に関書き (または、コピー□ペーストして) Return または、Enter キーを押してください。実行結果が、その下に出ます。最後の、`plot(cars)` は、`cars` というデータの、散布図が右下の、Plots タブに表示されます。

- `head(cars)`
- `str(cars)`
- `summary(cars)`
- `plot(cars)`

エラーが表示されたら、もう一度、スペルを確認して、入力してみてください。

次のような、結果が表示されると思います。簡単な説明をつけます。

```
head(cars)
#>   speed dist
#> 1     4    2
#> 2     4   10
#> 3     7    4
#> 4     7   22
#> 5     8   16
#> 6     9   10
```

`head(cars)` は、`cars` という、R に付属している、データの、最初 (頭 head) の

6 行を、表示します。

```
str(cars)
#> 'data.frame':   50 obs. of  2 variables:
#> $ speed: num  4 4 7 7 8 9 10 10 10 11 ...
#> $ dist : num  2 10 4 22 16 10 18 26 34 17 ...
```

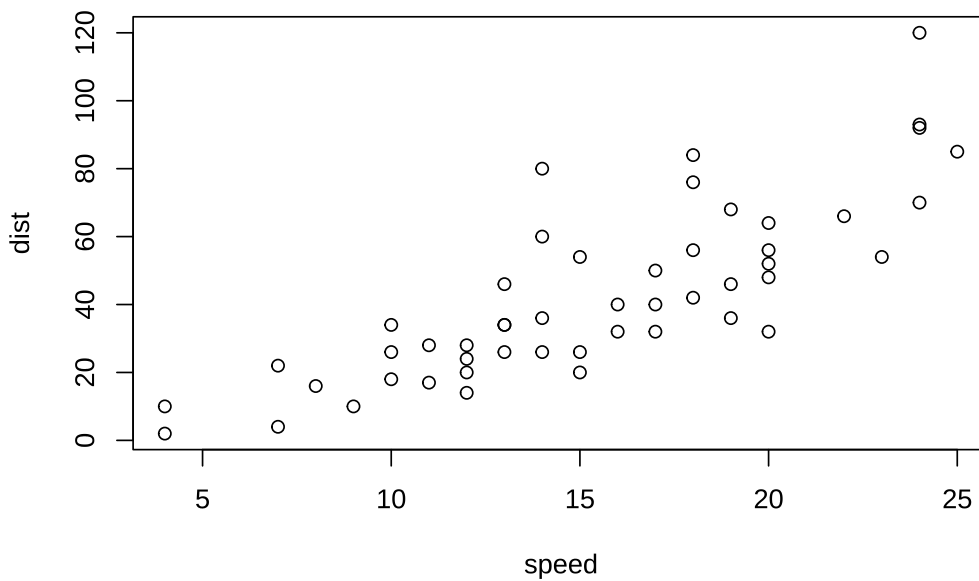
`str(cars)` は、`cars` という、R に付属している、データの構造 (structure) を表示します。`data.frame` とありますが、これは、矩形になったデータ (各列の長さがおなじ) の一番簡単なクラスの名前で、2 変数、それぞれが、50 個の数値データ (numerical data) からなっていることがわかります。

`head(cars)` では、縦に表示されていたものが、横に表示されています。`speed` `dist` とありますが、`cars$`speed``,`cars$`dist`` は、`cars`` データの、それぞれの列を意味します。

```
summary(cars)
#>      speed      dist
#> Min.   : 4.0   Min.   :  2.00
#> 1st Qu.:12.0   1st Qu.: 26.00
#> Median :15.0   Median : 36.00
#> Mean   :15.4   Mean    : 42.98
#> 3rd Qu.:19.0   3rd Qu.: 56.00
#> Max.   :25.0   Max.    :120.00
```

`cars` データの概要 (summary) が表示されます。各列 (変数) について、最小値 (Minimum)、小さい方から、4 分の 1 を切り捨てたときの最小の値 (1st Quadrant)、中央値 (Median)、平均 (Mean)、大きい方から、4 分の 1 を切り捨てたときの最大の値 (3rd Quadrant)、最大値 (Maximum) が表示されます。

```
plot(cars)
```

0.3.5.2 アサインメント、ヘルプ

コンソールで次のそれぞれを、試してみてください。

- `df <- cars`

`df` に、`cars` をアサインします。すなわち、`df` が、`cars` の内容に置き換わります。`cars` はデータですが、データを含む、オブジェクトの名前を設定するためにも使います。オブジェクト名は、英文字から始まれば、かなりの自由度がありますが、わたしは、英文字と数字と `_` (underscore) 程度しか使わないようにしています。

- `head(df)`

`head(df)` は、`head(cars)` と同じ出力が得られます。

- `View(cars)`

左上の、窓枠が開き、`cars` の内容が表示されます。列名のところには、三角形も表示され、それを用いると、大きい順、小さい順などに、並び替えることも可能です。

- `?cars`

右下の、窓枠の `Help` タブに、`cars` の情報が表示されます。`Help` タブにある、虫眼鏡がついた、検索窓 (search window) に、`cars` といっても、同じ結果が得られます。内容を確認してください。

一番上には `cars {datasets}` とありますが、これは、`datasets` というパッケージの、`cars` だという意味です。そこで、`datasets` を調べてみましょう。

- `?datasets`

“The R Datasets Package” だと書かれていて、さらに、

This package contains a variety of datasets. For a complete list, use `library(help = "datasets")`.

さまざまなデータが含まれています。全てのリストをみるには、`library(help = "datasets")` を使ってください。

とありますから、`library(help = "datasets")` をコンソールに入力してみてください。

- `library(help = "datasets")`

左上の窓枠に、リストが表示されます。古いデータばかりですが、例として使うには、十分すぎるぐらいの、数のデータがあります。これらは、Toy Data（おもちゃのデータ）と呼ばれることもあります。

`cars` も見つかりましたか。

0.3.5.3 おすすめ

コンピュータのシステムが、日本語であると、R の言語も日本語になっているはずですが、そこで、エラーが発生すると、一部、日本語で表示されます。しかし、ネット上などで、そのエラーの対応を検索するときは、英語のエラーメッセージで検索した方が、解決方法が得られる可能性が高いので、わたしは、英語に設定しています。英語にするには、Console で次のようにします。

言語を英語に設定: `Sys.setenv(LANG = "en")`

RStudio を終了して、もう一度起動すると、日本語に戻っていると思います。ですから、作業の最初、または、エラーが出たら、変更することをお勧めします。

日本語に戻りたいときは、次のようにします。

言語を日本語に設定: `Sys.setenv(LANG = "ja")`

さまざまな Help など、すべて日本語で表示されれば日本語を使うのは有効かもしれませんが、すくなくとも、現在は、そうではないので、上に説明したことから、英語に設定することをお勧めします。

0.3.5.4 練習

1. `head(cars, 10L)` は何が出力されますか。`head(cars, n=10L)` と同じですか。
2. `?head` または、Help の検索窓に `head` と入力して、説明を見てみてください。`head(cars, n=10L)` などについて、書いてありましたか。他には、どのようなことが分かりましたか。

3. `datasets` のデータのいくつかについて、そのデータの `help` や、`head`, `str`, `summary` などを使ってみてください。これらで表示できない場合がありますか。データについては、最初に、これら、三つを試してみることをお勧めします。わかったことをメモしておくといいでしょう。`datasets` のリストをみるには、`library(help = "datasets")` でしたね。

0.3.6 RStudio について

RStudio は多くの機能を持っています。

0.3.6.1 四つの窓枠とタブ Four Panes and Tabs

- Top Left: Source Editor
- Top Right: Environment, History, etc.
- Bottom Left: Console, Terminal, Render, Background Jobs
- Bottom Right: Files, Plots, Packages, Help, Viewer, Presentation

0.3.7 R Script 実行記録

R Script を使って、コードを実行すると、その記録を残すことができます。

0.3.7.1 R Script の作成

- RStudio の上のメニュー□バーから File > New File > R Script を選択します。
- File > Save As で、名前をつけて保存します。`{file_name}.R` が作成されます。
 - 右下の、Files から、ファイルを確認してください。
- `head(cars)`, `str(cars)`, `summary(cars)`, `plot(cars)` などと改行をしながらコードを書きます。
- 実行するには、カーソルの場所で `Ctrl+Shift+Enter` (Win) または `Cmd+Shift+Enter` (Mac) とすると、カーソルのある行か、その下の行で、最初のコードが実行されます。
 - R Script エディターの上にある、Run ボタンを押しても、同様に実行されます。
 - Run ボタンの右の、Source ボタンを押すと、そのスクリプトの、最初からすべて実行されます。
- 最後には保存しておきましょう。

0.3.7.2 R Script による実行

新しく、R Script を作成し、この下の、コード（ハイライトされている部分）をコピー□ペーストして、保存し、実行してみてください。

それぞれ、どのようなことをしているでしょうか。

```
#####  
#  
# basics.R  
#  
#####  
# 'Quick R' by DataCamp may be a handy reference:  
#   https://www.statmethods.net/management/index.html  
# Cheat Sheet at RStudio: https://www.rstudio.com/resources/cheatsheets/  
# Base R Cheat Sheet: https://github.com/rstudio/cheatsheets/raw/main/base-r  
# To execute the line: Control + Enter (Window and Linux), Command + Enter (Mac)  
## try your experiments on the console  
  
## calculator  
  
3 + 7  
  
### +, -, *, /, ^ (or **), %, %/  
  
3 + 10 / 2  
  
3^2  
  
2^3  
  
2*2*2  
  
### assignment: <-, (=  
x <- 5  
  
x
```

```
#### object_name <- value, '<-' shortcut: Alt (option) + '-' (hyphen or minus)
#### Object names must start with a letter and can only contain letter, numbers, _ and .

this_is_a_long_name <- 5^3

this_is_a_long_name

char_name <- "What is your name?"

char_name

#### Use 'tab completion' and 'up arrow'

### ls(): list of all assignments

ls()
ls.str()

#### check Environment in the upper right pane

### (atomic) vectors

5:10

a <- seq(5,10)

a

b <- 5:10

identical(a,b)

seq(5,10,2) # same as seq(from = 5, to = 10, by = 2)

c1 <- seq(0,100, by = 10)

c2 <- seq(0,100, length.out = 10)

c1
```

```
c2

length(c1)

#### ? seq    ? length    ? identical

(die <- 1:6)

zero_one <- c(0,1) # same as 0:1

die + zero_one # c(1,2,3,4,5,6) + c(0,1). re-use

d1 <- rep(1:3,2) # repeat

d1

die == d1

d2 <- as.character(die == d1)

d2

d3 <- as.numeric(die == d1)

d3

### class() for class and typeof() for mode
### class of vectors: numeric, characters, logical
### types of vectors: doubles, integers, characters, logicals (complex and ra

typeof(d1); class(d1)

typeof(d2); class(d2)

typeof(d3); class(d3)

sqrt(2)
```

```
sqrt(2)^2

sqrt(2)^2 - 2

typeof(sqrt(2))

typeof(2)

typeof(2L)

5 == c(5)

length(5)

### Subsetting

(A_Z <- LETTERS)

A_F <- A_Z[1:6]

A_F

A_F[3]

A_F[c(3,5)]

large <- die > 3

large

even <- die %in% c(2,4,6)

even

A_F[large]

A_F[even]

A_F[die < 4]
```

```
### Compare df with df1 <- data.frame(number = die, alphabet = A_F)
df <- data.frame(number = die, alphabet = A_F, stringsAsFactors = FALSE)

df

df$number

df$alphabet

df[3,2]

df[4,1]

df[1]

class(df[1])

class(df[[1]])

identical(df[[1]], die)

identical(df[1],die)

#####
# The First Example
#####

plot(cars)

# Help

? cars

# cars is in the 'datasets' package

data()

# help(cars) does the same as ? cars
# You can use Help tab in the right bottom pane
```



```
help(plot)
? par

head(cars)

str(cars)

summary(cars)

x <- cars$speed
y <- cars$dist

min(x)
mean(x)
quantile(x)

plot(cars)

abline(lm(cars$dist ~ cars$speed))

summary(lm(cars$dist ~ cars$speed))

boxplot(cars)

hist(cars$speed)
hist(cars$dist)
hist(cars$dist, breaks = seq(0,120, 10))
```

0.3.7.2.1 スクリプト 1: basics.R

```
# https://coronavirus.jhu.edu/map.html
# JHU Covid-19 global time series data
# See R package coronavirus at: https://github.com/RamiKrispin/coronavirus
# Data taken from: https://github.com/RamiKrispin/coronavirus/tree/master/csv
# Last Updated
Sys.Date()

## Download and read csv (comma separated value) file
```

```

coronavirus <- read.csv("https://github.com/RamiKrispin/coronavirus/raw/master/
# write.csv(coronavirus, "data/coronavirus.csv")

## Summaries and structures of the data
head(coronavirus)
str(coronavirus)
coronavirus$date <- as.Date(coronavirus$date)
str(coronavirus)

range(coronavirus$date)
unique(coronavirus$country)
unique(coronavirus$type)

## Set Country
COUNTRY <- "Japan"
df0 <- coronavirus[coronavirus$country == COUNTRY,]
head(df0)
tail(df0)
(pop <- df0$population[1])
df <- df0[c(1,6,7,13)]
str(df)
head(df)
### alternatively,
head(df0[c("date", "type", "cases", "population")])
###

## Set types
df_confirmed <- df[df$type == "confirmed",]
df_death <- df[df$type == "death",]
df_recovery <- df[df$data_type == "recovery",]
head(df_confirmed)
head(df_death)
head(df_recovery)

## Histogram
plot(df_confirmed$date, df_confirmed$cases, type = "h")
plot(df_death$date, df_death$cases, type = "h")
# plot(df_recovered$date, df_recovered$cases, type = "h") # no data for recovered

## Scatter plot and correlation

```

```

plot(df_confirmed$cases, df_death$cases, type = "p")
cor(df_confirmed$cases, df_death$cases)

## In addition set a period
start_date <- as.Date("2022-07-01")
end_date <- Sys.Date()
df_date <- df[df$date >=start_date & df$date <= end_date,]
##

## Set types
df_date_confirmed <- df_date[df_date$type == "confirmed",]
df_date_death <- df_date[df_date$type == "death",]
df_date_recovery <- df_date[df_date$data_type == "recovery",]
head(df_date_confirmed)
head(df_date_death)
head(df_date_recovery)

## Histogram
plot(df_date_confirmed$date, df_date_confirmed$cases, type = "h")
plot(df_date_death$date, df_date_death$cases, type = "h")
# plot(df_date_recovered$date, df_date_recovered$cases, type = "h") # no data for recovery

plot(df_date_confirmed$cases, df_date_death$cases, type = "p")
cor(df_date_confirmed$cases, df_date_death$cases)

#### Extra
plot(df_confirmed$date, df_confirmed$cases, type = "h",
     main = paste("Confirmed Cases in",COUNTRY),
     xlab = "Date", ylab = "Number of Cases")

```

0.3.7.2.2 スクリプト 2: coronavirus.T

0.3.7.3 練習

上の、coronavirus.R について

1. COUNTRY <- "Japan" の Japan を他の国に変えてみましょう。
2. start_date <- as.Date("2022-07-01") の日付を、他の日付に変えてみましょう。
3. df_confirmed\$cases と df_death\$cases についてどんなことがわかりま

すか。

4. 発見や、問いがあれば、書き出してみましょう。

0.3.7.4 Tips

キーボードショートカットと言われる、さまざまな機能があります。

- 上のメニューバー: Help > Keyboard Short Cut Help 確認してみてください。
- 右下の窓枠: Files タブから、ファイルの確認ができます。

0.3.8 パッケージ - Packages

R packages are extensions to the R statistical programming language containing code, data, and documentation in a standardised collection format that can be installed by users of R using Tool > Install Packages in the top menu bar of R Studio. https://en.wikipedia.org/wiki/R_package

R パッケージは、R の拡張機能で、コード、データ、ドキュメントを標準化されたコレクション形式で含んでおり、標準的なものは、R Studio の Top Bar の Tool > Install Packages からインストールできます。

0.3.8.1 パッケージのインストール

いずれ使いますので、まずは、三つのパッケージをインストールしてみましょう。

- `tidyverse`
- `rmarkdown`
- `tinytex`

インストール方法はいくつかあります。

一つ目は、上のメニューバーの Tool から、Install Packages ... を選択します。そして、パッケージズにインストールしたい、パッケージ名を入力します。そのパッケージ名が下にも出れば、Install ボタンを押してください。入力した名前の下にパッケージ名が出ない場合は、スペルが間違っている可能性がありますから、確認して、入れ直してください。

Console に、`install.packages("tidyverse")` などと表示され、たくさんメッセージが出ます。終了すると、> のマークがでます。

二つ目は、`install.packages("tidyverse")` のような書式で書いて、Console に入れる方法です。

三つ目は、右下の窓枠の Packages のタブにある、Install というボタンを押す方法です。すると、一番目の方法に、戻り、パッケージ名を入力できるようになりま

す。

この Packages タブにある、ものが、すでに、インストールされているパッケージです。そのなかで、**base** や、**datasets** などいくつかは、チェックがついていると思いますが、それらは、ロードされていて、いつでも、使える状態になっていることを意味しています。ロードは、たとえば、`library(tidyverse)` のようにしますが、それは、いずれもう一度説明します。

インストールは一回だけ。ときどき、Tools > Check for Package Update をつかって、Update しておくといいでしょう。

0.3.8.2 備考

Package によっては、Source から Compile するかと聞いてくる場合があります。どちらでも、良いのですが、特に、問題が起こっていなければ、No でよいと思います。コンピュータにあった形でインストールすることが必要な場合は、Yes とします。

同じパッケージをもう一度、インストールしたり、または、関連するパッケージがあるような場合、R をリスタートするかと聞いてくる場合があります。特に問題が起こらなければ、No で構いません。ただ、エラーが起こって、それに関連して、特別なパッケージをインストールする必要がある場合がありますが、そのときは、Yes としてください。

0.3.9 クラウド - Posit Cloud

RStudio Cloud は、誰でもオンラインでデータサイエンスを行い、共有し、教え、学ぶことができる、軽量でクラウドベースのソリューションです。

0.3.9.1 クラウドサービス How to Start Posit Cloud

まず、サインアップして、使ってください。一ヶ月の利用時間の限度など、設定されていますが、どこからでも、インターネットにつながっていれば使えるので、わたしは、いつくかアカウントを持って、活用しています。

1. Go to <https://posit.cloud/>
2. Sign Up: top right
3. Email address or Google account
4. New Project: Project Name

0.3.10 練習問題 Posit Primers

Posit Primers <https://posit.cloud/learn/primers>

教科書 “R for Data Science” は、**tidyverse** パッケージを中心に、データサイエンスについて解説したのですが、Posit Primers は、演習問題をしながら、教科書の内容を理解できるように構成されています。

0.3.10.1 最初の演習 The Basics – r4ds: Explore, I

- Visualization Basics
- Programming Basics

ぜひこれら二つの演習問題を、トライしてください。解説を読んでいただければ、データサイエンスは身につきます。

0.3.11 参考文献 References

一番目は、すでに紹介した、教科書です。二番目は、この文書を作成している、Bookdown というパッケージのサイトですが、そこに、たくさんの本が、無償で公開されています。素晴らしい本がたくさん含まれています。

- R For Data Science, by H. Wickham: <https://r4ds.had.co.nz>
 - Introduction: <https://r4ds.had.co.nz/explore-intro.html#explore-intro>
- Bookdown: <https://bookdown.org>, Archive

下の一番目は、R 入門を、2 時限の講義でしたときのものです。二番目と三番目は、講義で使ったものを、まとめたものです。教科書のようには、できていませんが、参考になる部分もあるかと思いますので、紹介しておきます。

- Introducton to R
- Data Analysis for Researchers 2022
- Data Analysis for Researchers 2021

0.3.12 YouTube Video - getstarted

- ファイル: <https://ds-sl.github.io/intro2r/getstarted.html>

0.4 R Markdown

0.4.1 Reproducible and Literate Programming

データサイエンスは、サイエンス（科学）ということばもついています。特に、根拠に基づいた（evidence based）とか、データに基づいた（data based）ということばを使うときには、なおさら、再現可能性（reproducibility）や、コードの内

容の説明などのコミュニケーションにも注力する必要があります。このことを心がけて、データサイエンスを学んでいきましょう。

表題にある、“Reproducible and Literate Programming” は、Reproducible（再現可能）かつ、Literate な（理解できるように記述した）Program（プログラム□コード）を共有することをたいせつにしましょうということです。

0.4.1.1 目的、問いなど

プロジェクトの目的、問いなどは、途中で変わっていくこともあります。その都度に、メモをしておくといいでしょう。

0.4.1.2 データについて

どのようなデータをどのように取得してきたかを、記録し、伝えられるようにすることが、必要です。データを取得するときから、取得方法や、それを伝える方法にも常に気をつけましょう。

0.4.1.3 コードについて

どのようなコードでそのグラフ（chart）などが得られたかも、単にコードを記述するだけでなく、それぞれの部分に、説明を付与することも有効です。

0.4.1.4 グラフについて

視覚化（visualization）は、とても有効です。そこで、見て理解したこと、観察したこと（observations）などは、簡単でも構いませんから、必ず、記録しておきましょう。

0.4.1.5 まとめ: R Markdown の目的

まさに、このようなことを可能にするのが、R Markdown です。少しずつ学んでいきましょう。

0.4.2 準備: パッケージのインストール

R パッケージは、R の拡張機能で、コード、データ、ドキュメントを標準化されたコレクション形式で含んでおり、標準的なものは、R Studio の Top Bar の Tool > Install Packages からインストールできます。

- tidyverse
- rmarkdown
- tinytex

インストールを複数回しても問題はありませんが、インストールされているかどうかは、Packages タブから確認することができます。

インストールは一回だけ。ときどき、Tools > Check for Package Update をつかって、Update しておくといいでしょう。

0.4.3 R Notebook

R Markdown はデータサイエンスのためのオーサリングフレームワーク。

コード（プログラム）とその実行結果、を記録□表示し、高品質のレポートの作成を可能にします。

R Notebook は、独立してインタラクティブに実行できるチャンクを持つ R Markdown ドキュメントの一つの形式で、入力のすぐ下に出力が表示することができます。

1. File > New File > R Notebook
2. Save with a file name, say, test-notebook
3. Preview by [Preview] button
4. Run Code Chunk `plot(cars)` and then Preview again.

0.4.4 日本語のテンプレート

下のリンクを開き、右上の Code ボタンから、Download Rmd を選択すると、ダウンロードできますから、ダウンロードしたものを、プロジェクト□フォルダーに移動またはコピーしてください。ダウンロードできないときは、Ctrl を押しながら、Download Rmd をクリックすると、Save As で保存できると思います。ブラウザによって仕様が異なりますから、適切な方法を選んでください。

- <https://ds-sl.github.io/intro2r/RNotebook-J.nb.html>
- <https://ds-sl.github.io/intro2r/Rmarkdown-J.nb.html>

Windows でも、Mac でも提供されている、Google Chrome の場合には、Code ボタンから、ダウンロードされるはずです。

0.4.5 R Markdown いくつかの Output

```
title: "Testing R Markdown Formats"
author: "ID Your Name"
header-includes:
- \usepackage{ctex}
```



```
output:
  html_notebook: default
  html_document: default
  pdf_document: default
    latex_engine: xelatex
  word_document: default
  powerpoint_presentation: default
  ioslides_presentation: default
---
```

PDF でエラー? コンソールで `tinytex::install_tinytex()`

- TeX システムがインストールされている場合は不要

0.4.6 YouTube Video - rmarkdown

第 III 部

PART III INSTITUTIONAL DATA

0.5 World Bank

第 IV 部

PART IV EDA

0.6 What is EDA?

第 V 部

PART V EXAMPLES

0.7 Example 1

.1 日本語の扱いについて

.1.1 日本語・中国語・韓国語

文字化けが、起こることが多く、対応が、一定せず、難しかったのですが、どうやら、現在は、どの場合も、次の設定で、解決しているようです。下の例を確認してください。

```
# showtext を、インストールしていない場合は、一回だけ、右上の三角をクリックして実行  
install.packages('showtext')
```

.1.1.1 パッケージをロード

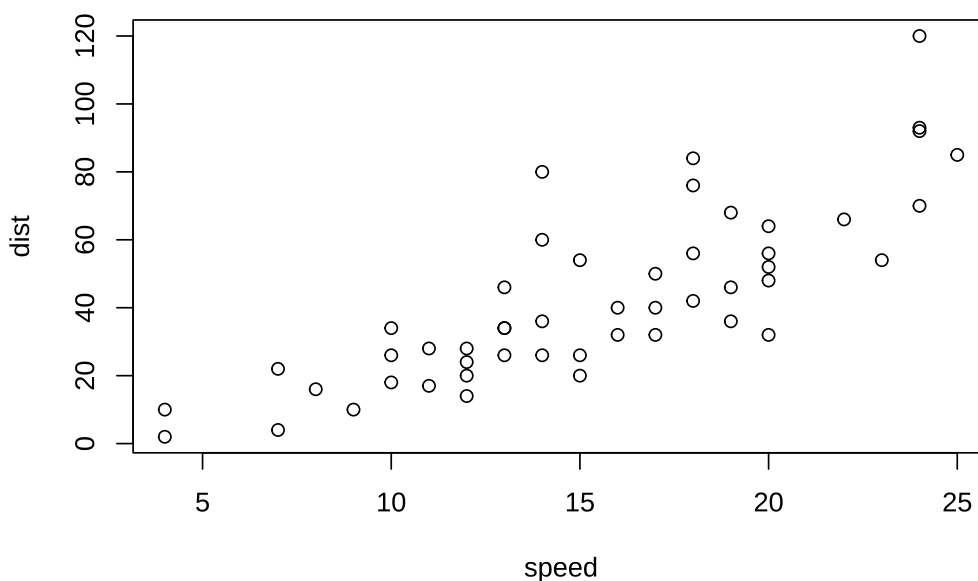
`library` によって、Package をロード（いつでも使えるように）します。

```
library(tidyverse)  
library(showtext)  
font_add_google('Noto Sans')  
showtext_auto()
```

.1.2 Base R でタイトルに日本語

```
plot(cars, main=" 散布図")
```

散布図



.1.3 列名や、データに日本語

```
df_iris <- iris
colnames(df_iris) <- c(" 萼長", " 萼幅", " 葉長", " 葉幅", "Species" )
tab <- data.frame(Species = c("setosa", "versicolor", "virginica"),
                  " 種別" = c(" ヒオウギアヤメ", " ブルーフラッグ", " バージニカ" ) )
df_iris <- df_iris %>% left_join(tab, by=c("Species" = "Species")) %>% select(
df_iris %>% slice(1:2)
#>   萼長 萼幅 葉長 葉幅      種別
#> 1  5.1  3.5  1.4  0.2 ヒオウギアヤメ
#> 2  4.9  3.0  1.4  0.2 ヒオウギアヤメ
```

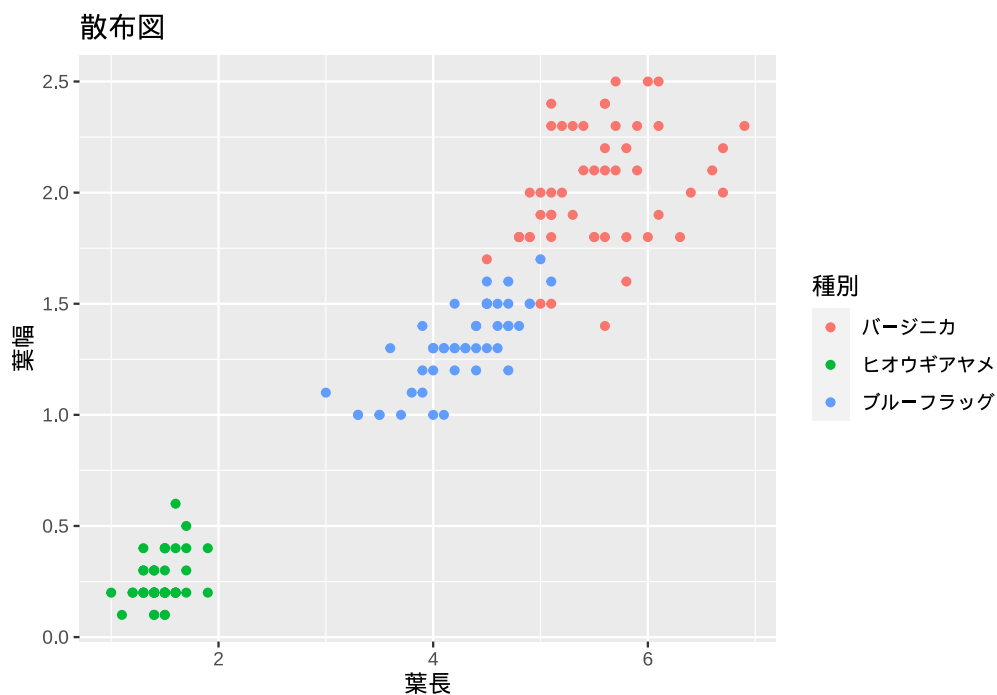
.1.4 kable で表示

```
knitr::kable(df_iris[1:6, ])
```

萼長	萼幅	葉長	葉幅	種別
5.1	3.5	1.4	0.2	ヒオウギアヤメ
4.9	3.0	1.4	0.2	ヒオウギアヤメ
4.7	3.2	1.3	0.2	ヒオウギアヤメ
4.6	3.1	1.5	0.2	ヒオウギアヤメ
5.0	3.6	1.4	0.2	ヒオウギアヤメ
5.4	3.9	1.7	0.4	ヒオウギアヤメ

.1.5 ggplot でグラフを作成

```
ggplot(df_iris, aes(x = `葉長`, y = `葉幅`, col = `種別`)) +
  geom_point() + labs(title = " 散布図", x = " 葉長", y = " 葉幅")
```



.1.6 備考:

実は、一番難しいのが、PDF の作成だと思いますが、一応、上のものも、PDF を作成することが可能です。下のリンクのファイルを、いろいろな、形式で、出力してみてください。R Notebook と、PDF に出力したもののリンクを付けておきます。

- R Notebook
 - 右上の Code ボタンから、Rmd ファイルも取得できます。
- PDF

.1.7 参考：日本語の表示について

日本語が適切に表示されない!?

簡単ではなく、未解決の部分が何かなどを含め、わたしも十分理解できているか不明であるが、理解できていると思われる範囲で、備忘録のように記す。

R を使うという場合に限っても、R Studio IDE を使う場合、RStudio Cloud を使う場合、Google colab を使う場合、他のプラットフォームで使う場合で違ってくると思われる。一応、上にあげた、三種類のプラットフォームで確かめられるものについて書いていく。上に書いた以外に、R Studio IDE を、Windows 上で使う場合と、Mac 上で使う場合（Mac のシステムは Unix 系であるが、さまざまな Linux）でも、状況が異なる。そこで、場合分けをして書いていくほうが安全であるが、それは、極力避け、どれにでも適用可能な方法を模索しながら書いていこ

うと思っている。個人的に、日常的に分断を避ける努力をすることが大切だとももっていることも背景にある。さらに、ソフトウェア開発者は、むしろ、そのような差異を理解して、どの環境でも、可能なように設計することを心がけていると思われるし、そのようなものが、R Project の正規のパッケージとして採用されていくべきだとも考えているので、多少、理想も入っているが、これを基本として書いていこうと思う。十分なチェックができていないものもあるので、不具合などは、ぜひ、お知らせ願いたい。この文章も少しずつ、改善していければと思う。

通常、日本語、中国語、韓国語などが適切に表示できない場合は、文字のエンコーディング (Encoding: どのような情報として記録されているか) と、フォントの問題、さらに、システムがこれらをどう処理しているかの問題があると思われる。しかし、R の利用者として考えると、文字化けが起きたり、適切に文字が表示されないのは、以下の三つに分けられるように思われる。

1. データファイルなどを読み込んだときに適切に表示されない
2. 図の中のタイトルなどが、適切に表示されない
3. R Markdown の出力において、適切に表示されない

1.7.1 データファイルの読み込み

- tidyverse に含まれる readr には、guess_encoding が含まれており、一般的には、たとえば、
 - read_csv("./data/file_name.csv") とすると、一番可能性の高いエンコーディングで読み込まれる。
- 使い方: guess_encoding(file, n_max = 10000, threshold = 0.2) とあり、10000 行で推測されたエンコーディング、または、確率を計算することを Default にしている。Help によると、すべての行をチェックする場合は、n_max = -1 とすることが書かれている。
- これで問題がない場合が多い。他の、readr 関数も同様である。
- なお、read_csv などにも、guess_max = min(1000, n_max) も含まれるが、これは、column type を決めるためのものである。
- read.csv() など、base R では、fileEncoding = “”, encoding = “unknown” がオプションに含まれていたので、指定して読み込むことが通常であった。

1.7.2 図の中のテキスト

- 基本的には、図の表示の前に library(showtext); font_add_google('Noto Sans'); showtext_auto() となっていれば、これ以降の図は、Google Fonts ‘Noto Sans’ が使われ、表示されるはずである。
- 二種類以上のフォントを使い分けたいときは、名前をつけて、それを family = name で指定する。
 - showtext: Using Fonts More Easily in R Graphs 参照。

.1.7.3 R Markdown の出力

- PDF 作成における問題が最後まで残っていたが、最近は、showtext Package で解決しているようである。設定については、図の中のテキストの場合と同じ。

.1.7.4 参考としたもの

.1.7.4.1 showtext: Using Fonts More Easily in R Graphs

- <https://CRAN.R-project.org/package=showtext>
 - <https://cran.r-project.org/web/packages/showtext/readme/README.html>
 - showtext: Using Fonts More Easily in R Graphs:
 - * <https://cran.r-project.org/web/packages/showtext/vignettes/introduction.html>
 - * <https://fonts.google.com>

.1.7.4.2 sysfonts: Loading Fonts into R

- <https://CRAN.R-project.org/package=sysfonts>
 - <https://cran.r-project.org/web/packages/sysfonts/sysfonts.pdf>

.1.7.4.3 foods4all: Examples of Graphs

- https://foods4all.github.io/examples/examples_of_graphs.html
 - 77.2 Japanese Environments 日本語環境（昔の記事: Last Updated: 2020-04-22)

.2 Bookdown

.2.1 About

This is a *sample* book written in **Markdown**. You can use anything that Pandoc's Markdown supports; for example, a math equation $a^2 + b^2 = c^2$.

.2.1.1 Usage

Each **bookdown** chapter is an `.Rmd` file, and each `.Rmd` file can contain one (and only one) chapter. A chapter *must* start with a first-level heading: `# A good chapter`, and can contain one (and only one) first-level heading.

Use second-level and higher headings within chapters like: `## A short section` or `### An even shorter section`.

The `index.Rmd` file is required, and is also your first book chapter. It will be the homepage when you render the book.

.2.1.2 Render book

You can render the HTML version of this example book without changing anything:

1. Find the **Build** pane in the RStudio IDE, and
2. Click on **Build Book**, then select your output format, or select “All formats” if you’d like to use multiple formats from the same book source files.

Or build the book from the R console:

```
bookdown::render_book()
```

To render this example to PDF as a `bookdown::pdf_book`, you’ll need to install XeLaTeX. You are recommended to install TinyTeX (which includes XeLaTeX): <https://yihui.org/tinytex/>.

.2.1.3 Preview book

As you work, you may start a local server to live preview this HTML book. This preview will update as you edit the book when you save individual `.Rmd` files. You can start the server in a work session by using the RStudio add-in “Preview book”, or from the R console:

```
bookdown::serve_book()
```

.2.2 Hello bookdown

All chapters start with a first-level heading followed by your chapter title, like the line above. There should be only one first-level heading (`#`) per `.Rmd` file.

.2.2.1 A section

All chapter sections start with a second-level (`##`) or higher heading followed by your section title, like the sections above and below here. You can have as many as you want within a chapter.

An unnumbered section Chapters and sections are numbered by default. To unnumber a heading, add a `{.unnumbered}` or the shorter `{-}` at the end of the heading, like in this section.

.2.3 Cross-references

Cross-references make it easier for your readers to find and link to elements in your book.

.2.3.1 Chapters and sub-chapters

There are two steps to cross-reference any heading:

1. Label the heading: `# Hello world {#nice-label}`.
 - Leave the label off if you like the automated heading generated based on your heading title: for example, `# Hello world = # Hello world {#hello-world}`.
 - To label an un-numbered heading, use: `# Hello world {-#nice-label}` or `{# Hello world .unnumbered}`.
2. Next, reference the labeled heading anywhere in the text using `\@ref(nice-label)`; for example, please see Chapter .2.3.
 - If you prefer text as the link instead of a numbered reference use: any text you want can go here.

.2.3.2 Captioned figures and tables

Figures and tables *with captions* can also be cross-referenced from elsewhere in your book using `\@ref(fig:chunk-label)` and `\@ref(tab:chunk-label)`, respectively.

See Figure 1.

```
par(mar = c(4, 4, .1, .1))
plot(pressure, type = 'b', pch = 19)
```

Don't miss Table 1.

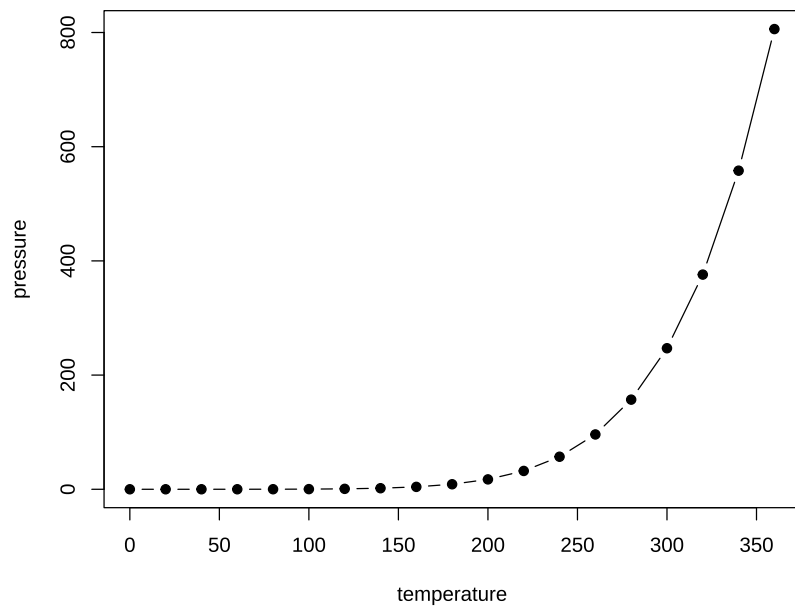


图1 Here is a nice figure!

表1 Here is a nice table!

temperature	pressure
0	0.0002
20	0.0012
40	0.0060
60	0.0300
80	0.0900
100	0.2700
120	0.7500
140	1.8500
160	4.2000
180	8.8000

```
knitr::kable(
  head(pressure, 10), caption = 'Here is a nice table!',
  booktabs = TRUE
)
```

2.4 Parts

You can add parts to organize one or more book chapters together. Parts can be inserted at the top of an .Rmd file, before the first-level chapter heading in that

same file.

Add a numbered part: `# (PART) Act one {-}` (followed by `# A chapter`)

Add an unnumbered part: `# (PART*) Act one {-}` (followed by `# A chapter`)

Add an appendix as a special kind of un-numbered part: `# (APPENDIX) Other stuff {-}` (followed by `# A chapter`). Chapters in an appendix are prepended with letters instead of numbers.

.2.5 Footnotes and citations

.2.5.1 Footnotes

Footnotes are put inside the square brackets after a caret `^[]`. Like this one ^{*1}.

.2.5.2 Citations

Reference items in your bibliography file(s) using `@key`.

For example, we are using the **bookdown** package (Xie, 2023) (check out the last code chunk in `index.Rmd` to see how this citation key was added) in this sample book, which was built on top of R Markdown and **knitr** (Xie, 2015) (this citation was added manually in an external file `book.bib`). Note that the `.bib` files need to be listed in the `index.Rmd` with the YAML `bibliography` key.

The `bs4_book` theme makes footnotes appear inline when you click on them. In this example book, we added `cs1: chicago-fullnote-bibliography.cs1` to the `index.Rmd` YAML, and include the `.cs1` file. To download a new style, we recommend: <https://www.zotero.org/styles/>

The RStudio Visual Markdown Editor can also make it easier to insert citations: <https://rstudio.github.io/visual-markdown-editing/#/citations>

.2.6 Blocks

.2.6.1 Equations

Here is an equation.

$$f(k) = \binom{n}{k} p^k (1-p)^{n-k} \quad (1)$$

You may refer to using `\@ref(eq:binom)`, like see Equation (1).

^{*1} This is a footnote.

.2.6.2 Theorems and proofs

Labeled theorems can be referenced in text using `\@ref(thm:tri)`, for example, check out this smart theorem .2.1.

Theorem .2.1. *For a right triangle, if c denotes the length of the hypotenuse and a and b denote the lengths of the **other** two sides, we have*

$$a^2 + b^2 = c^2$$

Read more here <https://bookdown.org/yihui/bookdown/markdown-extensions-by-bookdown.html>.

.2.6.3 Callout blocks

The `bs4_book` theme also includes special callout blocks, like this `.rmdnote`.

You can use **markdown** inside a block.

```
head(beaver1, n = 5)
#>   day time temp activ
#> 1  346  840 36.33     0
#> 2  346  850 36.34     0
#> 3  346  900 36.35     0
#> 4  346  910 36.42     0
#> 5  346  920 36.55     0
```

It is up to the user to define the appearance of these blocks for LaTeX output.

You may also use: `.rmdcaution`, `.rmdimportant`, `.rmdtip`, or `.rmdwarning` as the block name.

The R Markdown Cookbook provides more help on how to use custom blocks to design your own callouts: <https://bookdown.org/yihui/rmarkdown-cookbook/custom-blocks.html>

.2.7 Sharing your book

.2.7.1 Publishing

HTML books can be published online, see: <https://bookdown.org/yihui/bookdown/publishing.html>

.2.7.2 404 pages

By default, users will be directed to a 404 page if they try to access a webpage that cannot be found. If you'd like to customize your 404 page instead of using the default, you may add either a `_404.Rmd` or `_404.md` file to your project root and use code and/or Markdown syntax.

.2.7.3 Metadata for sharing

Bookdown HTML books will provide HTML metadata for social sharing on platforms like Twitter, Facebook, and LinkedIn, using information you provide in the `index.Rmd` YAML. To setup, set the `url` for your book and the path to your `cover-image` file. Your book's `title` and `description` are also used.

This `bs4_book` provides enhanced metadata for social sharing, so that each chapter shared will have a unique description, auto-generated based on the content.

Specify your book's source repository on GitHub as the `repo` in the `_output.yml` file, which allows users to view each chapter's source file or suggest an edit. Read more about the features of this output format here:

https://pkgs.rstudio.com/bookdown/reference/bs4_book.html

Or use:

```
?bookdown::bs4_book
```


参考文献

- Xie, Y. (2015). *Dynamic Documents with R and knitr*. Chapman and Hall/CRC, Boca Raton, Florida, 2nd edition. ISBN 978-1498716963.
- Xie, Y. (2023). *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.32.