# Categorical Data Analysis
## Lecture 1 – Guidance · R Markdown

Graduate School of Advanced Science and Engineering
Rei Monden

2024.12.5

# Guidance

# General information for the course

▶ Course code : KA123001
▶ Credits : 2
▶ Lecture room : **EDU**K201
▶ Lecturer: Rei Monden
▶ Faculty: Graduate School of Advanced Science and Engineering
▶ Office: 総合科学・C714
▶ Email: mondenr@hiroshima-u.ac.jp

# About this course

▶ 15 lectures in total (Thursdays; 8:45 ~ 12:00)
▶ Lecture · Practical assignments
▶ Textbook : An Introduction to Categorical Data Analysis
  Third Ed. Alan Agresti
  This text book is a "bible" for categorical data analysis, but
  the translated version of this book is out of print.
▶ Textbook, Slides, R program
▶ Lecture 5 on wards, lectures will be given both in-person and
  online
▶ For English lectures, you can watch lecture videos online

# About this course

- For questions related to this course, please write it down in Teams, "Question" channel
  (I will check it regularly, but students can also help each other by answering to your colleagues' questions)
- Please DO NOT email me personally, if you have a question related to this course. Please use "Teams".

# About this course

▶ Evaluation will be done based on your contribution for the course (e.g., actively involved in Teams' discussion board)(5%), assignments(65%), final exam(30%).
▶ This course will be given mainly in Japanese.
▶ English lecture videos will be provided via **moodle** and the final exam will be provided in English.

# About the assignments

▶ In each lecture, an assignment will be given.
You have to hand in your assignment before a given deadline.
▶ Upload your assignment via Teams in a PDF format.
▶ In principle, I do not accept your report after a given deadline.
▶ I will count your "attendance" according to your submission of a report.
▶ Please write your report in English.
▶ Please set your file name as shown below, starting from your student number.
  (Example. B123456_1.pdf, B123456_Monden.pdf etc…)

# About the assignments

▶ In this lecture, you can use Posit Cloud (Rstudio Cloud) or, alternatively, an installed version of R and RStudio to prepare your report.

▶ In particular, we use the "RMarkdown" in Rstudio.

▶ Using Posit Cloud is sufficient for this course, but if you would like to install R and Rstudio to your PC, please search by yourself (Example of reference)

▶ We are going to learn very basics of how to use RMarkdown in this lecture.

# Class schedule

▶ Lec 1 (12/5) : Guidance, RMarkdown
▶ Lec 2 (12/5) : Review
▶ Lec 3 (12/12) : Analyzing contingency tables (theory)
▶ Lec 4 (12/12) : Analyzing contingency tables (practical)
▶ Lec 5 (12/19) : Generalized Linear Models (theory)
▶ Lec 6 (12/19) : Generalized Linear Models (practical)
▶ Lec 7 (12/24) : Logistic regression (theory)
▶ Lec 8 (12/24) : Logistic regression (practical)
▶ Lec 9 (1/9) : Multiple logistic regression
▶ Lec 10 (1/9) : Building logistic regression models (theory)
▶ Lec 11 (1/16) : Building logistic regression models (practical)
▶ Lec 12 (1/16) : Multicategory logit models
▶ Lec 13 (1/23) : Models for matched pairs
▶ Lec 14 (1/23) : Generalized Linear Mixed models
▶ Lec 15 (1/30) : Final Exam(8:45 ∼ 10:15)

▶ The schedule may change depending on the progress of classes.

# Goal of this course

1. Understand categorical data analytical methods, and apply them to real data.

▶ Why should we learn categorical data analytical methods? Because there is a lot of data which do not satisfy all assumptions of the general linear model (e.g., categorical dependent variables).
By learning the contents of this lecture, you will be able to analyze data that does not satisfy all the assumptions of the general linear model (highly versatile).

# Goal of this course

2. (side project) Get used to RStudio Markdown.

▶ Why RMarkdown?

Because RMarkdown can save data analysis codes, its results and interpretation in one report.

- ▶ It allows to analyze data, write a paper, or create presentation slides.
- ▶ Figures or Tables made in R can be inserted to the paper or presentation slides, directly.
- ▶ LaTeX, which is often used to write mathematical formula, can be used.
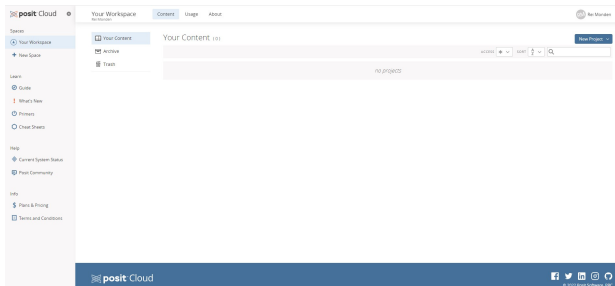
# Let's get started !

# Rstudio: Posit Cloud

1. Access to this link.
2. If you don't have an account, click on "Sign Up".
3. Fill in your email address, password and name. Then, click "Sign Up" (save your password somewhere. You will need it whenever you use PositCloud).

# Rstudio: Posit Cloud

4. You will receive an email from "Posit" to your registered account. Follow its instructions.
5. Revisit this link, but this time, click "Log In" and fill in your email address and password.
6. Select "Posit Cloud".
7. In case you reach a screen like this, your registration is completed!

# Rstudio: Notes for Posit Cloud

Notes for Posit Cloud (free version)

▶ Need access to the internet
▶ Cannot type in Japanese (In this lecture, students have to prepare their slides in English, so this should be fine)
▶ You can use <span style="color:red">up to 25 hours per month</span> (if you create more projects, your usage time gets shorter)
▶ Even when you are not typing anything, "Using hours" will be counted by opening the page (close the webpage whenever you don't need it)

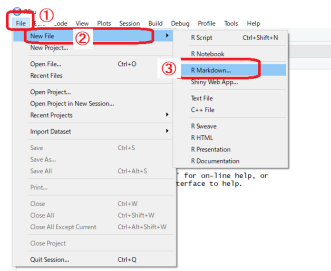# Rstudio: In case you use installed version of Rstudio

Notes

▶ To install the latest version of R and RStudio, your
  computer's OS version should satisfy the following condition
  (use 64 bit version):
  - ▶ Windows: Windows 7 or later
  - ▶ Mac: High Sierra or later
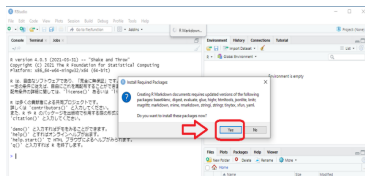  It is a good idea to update your PC to the latest version
▶ User name should be defined by Half-width characters
  - ▶ If you have double-byte character, you may be unable to install
    - ▶ BAD example: 門田　２０２３
    - ▶ Good example: Monden2023

# Let's use RMarkdown

1. Start RStudio.
2. Follow the below options.

3. The program will ask if you want to install R Markdown package. So click "Yes".

# Let's use RMarkdown

4. Start installing packages.

# Let's use RMarkdown

5. Once your installation is completed, fill in the following:

▶ ① Your student number
▶ ② Your name.

Confirm that you select HTML and then click "OK".

# Let's use RMarkdown

6. You will see something like below on your screen in RStudio. Click "Knit" inside of the red circle.

# Let's use RMarkdown

7. Declare a folder where you want to save your file.
   A file name could be "MY_STUDENTS_NUMBER_1", for instance, and click "save".



## Warnings

▶ Do not include Japanese characters or empty space in your file name.

▶ The first letter should be alphabetical (not a number or other characters).

▶ Define a file name which identifies what you will do in the code.

# Let's use RMarkdown

8. After a while, an HTML file will be generated.



9. Click "Open in Browser".

# Let's use RMarkdown

10. Open the file on a browser(Chrome, Firefox, Safari etc…),and go to "file" > "Print".

# Let's use RMarkdown

11. Select "Save as PDF" > "Save" button on the right bottom.



12. Check if a PDF was generated in a folder where you chose.

# Practice 1-1

Follow the previously mentioned steps and generate a PDF file in RMarkdown.
(In step 7, change the file name to your student number)

Why do we use RMarkdown ?

# Why do we use RMarkdown?

If you use RMarkdown, you can:

- ▶ record your analysis process.
- ▶ reuse your analysis code.
- ▶ save in various formats (HTML, PDF, DOC, etc…)
  (But if you want to save a document written in Japanese in a PDF format, you need to go through some settings and install packages.)
- ▶ insert figures or tables created in R, as well as image files (PNG, JPEG, etc…).
- ▶ easily manage reference lists when you write a paper.
- ▶ create slides for presentations.
- ▶ include and show your analysis code in R.
- ▶ use mathematical formulas, using LaTeX.
- ▶ replicate your analysis, since the analysis code and the report are saved in one file.

# Why is it important to replicate your analysis

▶ Research status in Japan

▶ The number of papers is decreasing each year



**The Countries Leading The World In Scientific Publications**

Number of science & engineering articles published in peer-reviewed journals in 2018

| | Total | Global share |
|---|---|---|
| China | 528,263 | 20.67% |
| United States | 422,808 | 16.54% |
| India | 135,788 | 5.31% |
| Germany | 104,396 | 4.08% |
| Japan | 98,793 | 3.87% |
| United Kingdom | 97,681 | 3.82% |
| Russia | 81,579 | 3.19% |
| Italy | 71,240 | 2.79% |

$\rightarrow$ Opening analysis code/results can accelerate research.

*Ref) National Science Foundation

# Why is it important to replicate your analysis

▶ The number of "influential" papers is decreasing



*Ref) 毎日新聞 (Mainichi news)

# Why is it important to replicate your analysis

▶ Moreover the number of retracted papers in Japan is one of the "top" in the world (7/30)

1. Joachim Boldt (194) See also: Editors-in-chief statement, our coverage
2. Yoshitaka Fujii (172) See also: Final report of investigating committee, our reporting, additional coverage
3. Hironobu Ueshima (124) See also: our coverage
4. Yoshihiro Sato (113) See also: our coverage
5. Ali Nazari (100) See also: our coverage
6. Jun Iwamoto (88) See also: our coverage
7. Diederik Stapel (58) See also: our coverage
8. Yuhji Saitoh (56) See also: our coverage
9. Adrian Maxim (48) See also: our coverage
10. A Salar Elahi (44) See also: our coverage
11. Chen-Yuan (Peter) Chen (43) See also: SAGE, our coverage
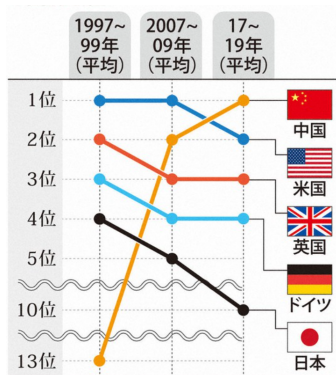12. Fazlul Sarkar (41) See also: our coverage
13. Shahaboddin Shamshirband (41) See also: our coverage
14. Hua Zhong (41) See also: journal notice
15. Shigeaki Kato (40) See also: our coverage
16. James Hunton (36) See also: our coverage
17. Hyung-In Moon (35) See also: our coverage
18. Dong Mei Wu (35) See also: National Natural Science Foundation of China finding
19. Antonio Orlandi (34) See also: our coverage
20. Dimitris Liakopoulos (33) (NB: We're counting a book he co-authored as a single retraction. The book has 13 retracted chapters with DOIs that are not included in this figure.) See also: our coverage
21. Jose L Calvo-Guirado (32) See also: our coverage
22. Jan Hendrik Schön (32) See also: our coverage
23. Amelec Viloria aka Jesus Silva (32) See also: our coverage
24. Naoki Mori (31) See also: our coverage
25. Bharat Aggarwal (30) See also: our coverage
26. Victor Grech (30) See also: our coverage
27. Soon-Gi Shin (30) See also: our coverage
28. Tao Liu (29) See also: our coverage
29. Jun Ren (29) See also: our coverage
30. Cheng-Wu Chen (28) See also: our coverage

7. Diederik Stapel (58) See also: our coverage

# R script vs RMarkdown

**Rscript**

▶ File extension .R
▶ In general, only R code and leave comments with #
▶ Purpose is to analyze data

**RMarkdown**

▶ File extension .Rmd
▶ In general, it allows for text (reports, notes).
▶ R code is written in a *code chunk*.
▶ The main purpose is to integrate analyses with reports.

# Notes when you make an RMarkdown file

▶ When you edit a file, save regularly by doing [Ctrl/Cmd + S]



▶ Each time you "Knit", your Rmd file is automatically saved.

How to write an RMarkdown file

# RMarkdown:YAML header

▶ Let's see the ".Rmd" file you made.

▶ YAML: The top part of the Rmd file enclosed by `---` marks.



```
1  ---
2  title: "M123456"
3  author: "門田　麗"
4  date: "2021/12/1"
5  output: html_document
6  ---
```

YAML

▶ The YAML is where you control information about a file or how you want to export it.

▶ This lecture is for categorical data analysis, so we only use a limited part of what RMarkdown can do.

▶ Reference) More information about YAML:

https://www.cloudbees.com/blog/yaml-tutorial-everything-you-need-get-started

# RMarkdown:YAML header

Edit the following sections of the YAML each time you prepare a report:

- ▶ `title:` "(Your student number)_(the number of lecture)"
  (Note: The title is different from the file name)
- ▶ `date:` "(The date you created your report)"
- ▶ `author:` "(Your name)"

Important:

- ▶ Colon(:) should be always placed after title, date and author.
- ▶ After colon, you need more than one half-width space (cannot use full-width space).
- ▶ Content of each element should be placed between half-width quotation marks「""」.

# RMarkdown: text

▶ Write down your text after a header of YAML (you could type in Japanese, but…).

▶ When you use Posit Cloud (RStudio Cloud), you cannot knit in Japanese (not solved yet).

▶ Prepare your report in English.

▶ Cheat sheet for RMarkdown: Click here.

▶ Besides this, more information or tutorials are available about RMarkdown. So please study about it by yourself.

# RMarkdown: syntax

▶ Syntax(left) and their output (right) in RMarkdown

| syntax | becomes |
|--------|---------|
| Plain text | Plain text |
| End a line with two spaces to start a new paragraph. | End a line with two spaces to start a new paragraph. |
| *italics* and _italics_ | *italics* and *italics* |
| **bold** and __bold__ | **bold** and **bold** |
| superscript^2^ | superscript$^2$ |
| ~~strikethrough~~ | ~~strikethrough~~ |
| [link](www.rstudio.com) | link |
| # Header 1 | **Header 1** |
| ## Header 2 | **Header 2** |
| ### Header 3 | **Header 3** |
| #### Header 4 | **Header 4** |
| ##### Header 5 | Header 5 |
| ###### Header 6 | Header 6 |

▶ # defines structure of text.
  ▶ #: The biggest section of your text      (section)
  ▶ ##: The second biggest section of your text (subsection)
▶ In Japanese, you cannot use italic font
▶ Do not put any empty space before and/or after a bold font or italic font

# RMarkdown: Syntax

▶ Syntax (left) and its output (right) in RMarkdown

```
* unordered list
* item 2
    + sub-item 1
    + sub-item 2

1. ordered list
2. item 2
    + sub-item 1
    + sub-item 2

Table Header  | Second Header
------------- | -------------
Table Cell    | Cell 2
Cell 3        | Cell 4
```

- unordered list
- item 2
    ◦ sub-item 1
    ◦ sub-item 2

1. ordered list
2. item 2
    ◦ sub-item 1
    ◦ sub-item 2

| Table Header | Second Header |
| --- | --- |
| Table Cell | Cell 2 |
| Cell 3 | Cell 4 |

▶ For numbered bullet lists, you can manually increase the number as 1., 2., etc…, but if you type 1. repeatedly, the "correct" number will be displayed in the output.

▶ You should have a half-width space right after a symbol.

# RMarkdown: Insert R code in the middle of the text

▶ Write R code between back-tick marks 「`」 (Shift+@) and write "r" right after the first back-tick mark.

| In RMarkdown... | Output... |
| --- | --- |
| the result of 1+2 is `r 1+2`. | the result of 1+2 is 3. |
| Pi is about `r round(pi, digits = 2)`. | Pi is about 3.14. |

# RMarkdown: Use LaTeX formula in a sentence (FYI)

▶ Formula in a middle of text should be placed between 「\$」
  (In Japanese PC, \ is shown as ¥)

| In RMarkdown... | Output... |
| --- | --- |
| The sum of `$x_i$` is `$\sum_{i=1}^{n}x_i$`. | The sum of $x_i$ is $\sum_{i=1}^{n} x_i$. |

▶ Displayed formula: Placed between 「\$\$」

| In RMarkdown... | Output... |
| --- | --- |
| The sum of `$x_i$` is `$$\sum_{i=1}^{n}x_i$$`. | The sum of $x_i$ is $$\sum_{i=1}^{n} x_i.$$ |

# Practice 1-2

We are keep working on the .Rmd file which you have created in Practice 1-1.

## Delete everything after "RMarkdown".

1. Write down your name in *italic* and **bold** font.
2. Insert a webpage link (any website is fine).
3. Set "Categorical Data Analysis" as a section and "RMarkdown" as a sub-section.
4. Make a sentence and insert some R code in the middle (R code can be whatever you want).
5. Make a list with bullet points.

# RMarkdown: code chunk

▶ A part between 3 back-ticks together 「```」 is called a "code chunk".
▶ You need to declare which language you use in the beginning of a chunk.
  ▶ In case you use R, write {r}
  ▶ You can also use other language besides R (currently, only Python is available)

これは通常の文章ですが、コードチャンクは
```{r}
x <- 1:10
mean(x)
```
のように書きます.

  ▶ Click on the right-top button on a code chunk to select what to display in your output.

# RMarkdown: code chunk

▶ Short cut for code chunk: Press 3 keys at the same time.
  ▶ Windows: Ctrl + Alt + I (alphabet "I")
  ▶ Mac : command + option + I (alphabet "I")
▶ Place your cursor on the row of the R script which you want to run and,
  ▶ Windows: Ctrl + Enter
  ▶ Mac: command * shift + return
▶ Or click the "Run" button (triangle sign) on the right-top of the chunk

# Learn more about RMarkdown!

For English
https://pkgs.rstudio.com/rmarkdown/articles/rmarkdown.html
https://www.dataquest.io/blog/r-markdown-guide-cheatsheet/
https://bookdown.org/yihui/rmarkdown-cookbook/

For Japanese

▶ 上にある３つ目の参考文献に関しては片桐智志氏によって和訳され
たものがあります
https://gedevan-aleksizde.github.io/rmarkdown-cookbook/

## Practice 1-3

Keep editing the .Rmd file which you have created in Practice 1-2.

Do the following in a code chunk.

1. Calculate the mean and variance of the "Sepal.Length" variable in the iris data (this data is stored in R by default) and cite the results in a text.
   Hint) To import data in R: data("iris")

2. Draw a scatter plot for variables "Sepal.Length" ($x$-axis) and "Petal.Length" ($y$-axis) from the iris data set, and cite their correlation in the middle of a sentence.
   Hint) Draw a scatter plot of variables x and y in R: plot(x, y)

3. Save the generated file in PDF format and upload it in Teams>Assignment>Assignment 1.
   Note that your file name should start from your student number!!