

dplyr & purrr を用いたデータハンドリング

17 July, 2017

Contents

1	Load Libraries	1
2	Read Data	1
2.1	Make tidy data	2
3	Transform data	3
4	Visualize	3
5	purrr	4
6	Map family	5

This exercise is taken from dplyr & purrr を用いたデータハンドリング. Original data is on e-Stat: 精神科病院の推計患者数, 年齢階級 × 性・疾病分類（精神及び行動の障害） × 入院－外来別

1 Load Libraries

```
library(tidyverse)
```

2 Read Data

Try reading data without any option. We can see untidy data.

```
data <- readr::read_csv("j0056.csv",  
                        locale = readr::locale(encoding="cp932")  
                        )  
knitr::kable(data[1:5, 1:3])
```

平成 2 6 年	患者調査	平成 2 6 年 1 0 月
上巻第 5 6 表	精神科病院の推計患者数, 年齢階級 × 性・疾病分類（精神及び行動の障害） × 入院－外来別	NA
NA	NA	総数
NA	NA	NA
入院	NA	NA
総数	総数	218.4

Let's skip first two lines to make it better.

```
data <- readr::read_csv("j0056.csv",  
                        locale = readr::locale(encoding="cp932"),  
                        skip=2)  
colnames(data)
```

2.1 Make tidy data

sex	disease_name	0 歳	1 ～ 4	5 ～ 9
男	血管性及び詳細不明の認知症	-	-	-
男	アルコール使用＜飲酒＞による精神及び行動の障害	-	-	-
男	その他の精神作用物質使用による精神及び行動の障害	-	-	-
男	統合失調症、統合失調症型障害及び妄想性障害	-	-	0.0
男	気分〔感情〕障害（躁うつ病を含む）	-	-	-

sex	disease_name	age	count
男	血管性及び詳細不明の認知症	0歳	-
男	アルコール使用＜飲酒＞による精神及び行動の障害	0歳	-
男	その他の精神作用物質使用による精神及び行動の障害	0歳	-
男	統合失調症，統合失調症型障害及び妄想性障害	0歳	-
男	気分〔感情〕障害（躁うつ病を含む）	0歳	-
男	神経症性障害，ストレス関連障害及び身体表現性障害	0歳	-
男	知的障害＜精神遅滞＞	0歳	-
男	その他の精神及び行動の障害	0歳	-
女	（再掲）精神及び行動の障害	0歳	-
女	血管性及び詳細不明の認知症	0歳	-

2

```
data %>%
  rowwise() %>%
  mutate(count = ifelse(count=="-", NA, count), # create NA
         disease_name = stringr::str_trim(disease_name, side="both") # trim spaces
  ) %>%
  mutate(count = as.numeric(count)) -> data
knitr::kable(data[1:10, 1:4])
```

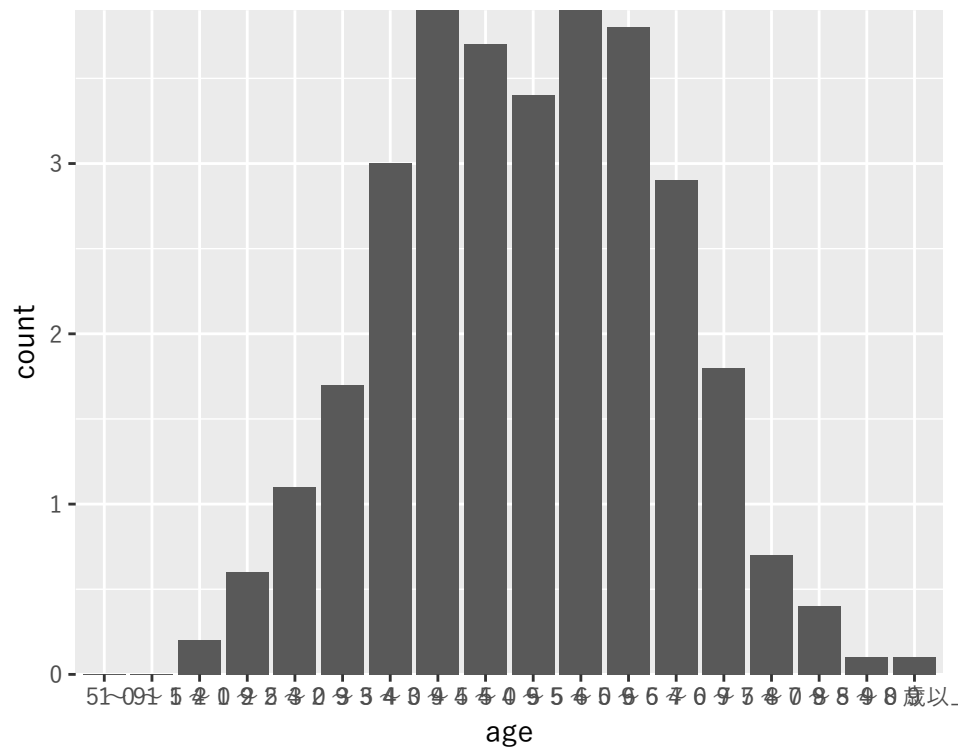
sex	disease_name	age	count
男	血管性及び詳細不明の認知症	0 歳	NA
男	アルコール使用<飲酒>による精神及び行動の障害	0 歳	NA
男	その他の精神作用物質使用による精神及び行動の障害	0 歳	NA
男	統合失調症, 統合失調症型障害及び妄想性障害	0 歳	NA
男	気分〔感情〕障害（躁うつ病を含む）	0 歳	NA
男	神経症性障害, ストレス関連障害及び身体表現性障害	0 歳	NA
男	知的障害<精神遅滞>	0 歳	NA
男	その他の精神及び行動の障害	0 歳	NA
女	（再掲）精神及び行動の障害	0 歳	NA
女	血管性及び詳細不明の認知症	0 歳	NA

3 Transform data

```
schizo <- data %>%
  filter(disease_name == "統合失調症, 統合失調症型障害及び妄想性障害",
         !is.na(count)) # do not take rows with NA
```

4 Visualize

```
ggplot(data=schizo) +
  geom_bar(aes(x=age, y=count), stat="identity") +
  scale_y_continuous(expand=c(0,0)) +
  theme_gray(base_family = "YuGo-Medium")
```



5 purrr

```
mapped <- data %>%
  split(list(. $disease_name, . $sex)) %>% # make list by name and sex
  map(group_by, age) %>% # applying a function to the all lists
  map(summarize, mean=mean(count, na.rm=T))
```

You can apply self-defined functions:

```
myhist <- function(data){
  ggplot(data) +
    geom_bar(aes(x=age, y=mean), stat="identity") +
    scale_y_continuous(expand=c(0,0)) +
    theme_gray(base_family = "YuGo-Medium")
}

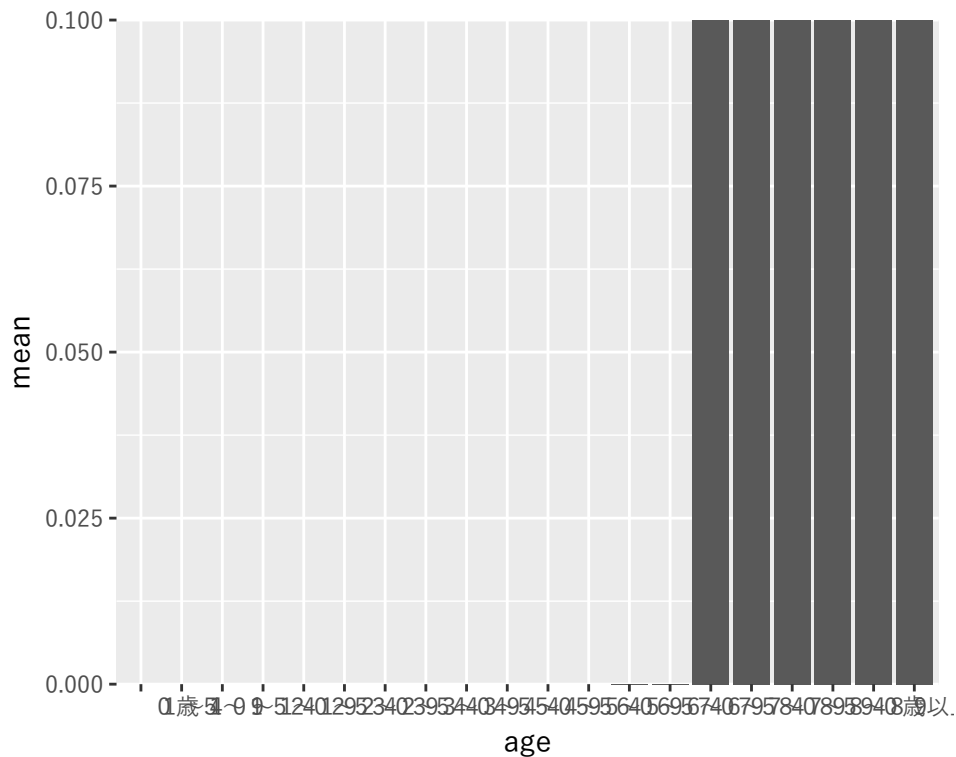
mappedhist <- data %>%
  split(list(. $disease_name, . $sex)) %>% # make list by name and sex
  map(group_by, age) %>% # applying a function to the all lists
  map(summarize, mean=mean(count, na.rm=T)) %>%
  map(myhist)
names(mappedhist)
```

```
## [1] "(再掲)精神及び行動の障害.女"
## [2] "アルコール使用<飲酒>による精神及び行動の障害.女"
## [3] "その他の精神作用物質使用による精神及び行動の障害.女"
## [4] "その他の精神及び行動の障害.女"
## [5] "気分[感情]障害(躁うつ病を含む).女"
## [6] "知的障害<精神遅滞>.女"
```

```
## [7] "神経症性障害, ストレス関連障害及び身体表現性障害.女"
## [8] "統合失調症, 統合失調症型障害及び妄想性障害.女"
## [9] "血管性及び詳細不明の認知症.女"
## [10] "(再掲)精神及び行動の障害.男"
## [11] "アルコール使用<飲酒>による精神及び行動の障害.男"
## [12] "その他の精神作用物質使用による精神及び行動の障害.男"
## [13] "その他の精神及び行動の障害.男"
## [14] "気分〔感情〕障害(躁うつ病を含む).男"
## [15] "知的障害<精神遅滞>.男"
## [16] "神経症性障害, ストレス関連障害及び身体表現性障害.男"
## [17] "統合失調症, 統合失調症型障害及び妄想性障害.男"
## [18] "血管性及び詳細不明の認知症.男"
```

You can take out any figure:

```
mappedhist$血管性及び詳細不明の認知症.男
```



You can make a list of analysis:

```
data %>% split(.$disease_name) %>%
  map(~lm(count ~ sex, data=..)) %>%
  map(summary) %>%
  map("coefficients")
```

6 Map family

- Change return values
- map_lg()
- map_chr()
- map_int()

- `map_db1()`
- `map` multiple lists
- `map2()`
- `map3()`
- `map_n()`