

dplyr Exercises

Elena Dubova

6/23/2020

Data from Harvard Database: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/O35FW8>

Entrance exam scores of students applying to a university in Brazil (Federal University of Rio Grande do Sul), along with the students' GPAs during the first three semesters at university. In this dataset, each row contains anonymized information about an applicant's scores on nine exams taken as part of the application process to the university, as well as their corresponding GPA during the first three semesters at university. The dataset has 43,303 rows, each corresponding to one student. The columns correspond to: 1) Gender. 0 denotes female and 1 denotes male. 2) Score on physics exam 3) Score on biology exam 4) Score on history exam 5) Score on second language exam 6) Score on geography exam 7) Score on literature exam 8) Score on Portuguese essay exam 9) Score on math exam 10) Score on chemistry exam 11) Mean GPA during first three semesters at university, on a 4.0 scale.

First, we import the data into R and brush it up.

```
data <- read.csv('GPA_Brazil.csv', header = FALSE)
head(data)
```

```
##   V1     V2     V3     V4     V5     V6     V7     V8     V9     V10    V11
## 1  0 622.60 491.56 439.93 707.64 663.65 557.09 711.37 731.31 509.80 1.33333
## 2  1 538.00 490.58 406.59 529.05 532.28 447.23 527.58 379.14 488.64 2.98333
## 3  1 455.18 440.00 570.86 417.54 453.53 425.87 475.63 476.11 407.15 1.97333
## 4  0 756.91 679.62 531.28 583.63 534.42 521.40 592.41 783.76 588.26 2.53333
## 5  1 584.54 649.84 637.43 609.06 670.46 515.38 572.52 581.25 529.04 1.58667
## 6  1 325.99 466.74 597.06 554.43 535.77 717.03 477.60 503.82 422.92 1.66667
```

```
colnames(data) <- c('Gender', 'Physics', 'Biology', 'History', 'Second_Language', 'Geography', 'Literature', 'Portuguese_Essay', 'Math', 'Chemistry', 'GPA_3S')
```

```
data$Student_ID <- paste0('Student_', 1:43303)
data <- data[c(12,1:11)]
```

```
head(data)
```

```
##   Student_ID Gender Physics Biology History Second_Language Geography
## 1 Student_1      0 622.60 491.56 439.93          707.64      663.65
## 2 Student_2      1 538.00 490.58 406.59          529.05      532.28
## 3 Student_3      1 455.18 440.00 570.86          417.54      453.53
## 4 Student_4      0 756.91 679.62 531.28          583.63      534.42
## 5 Student_5      1 584.54 649.84 637.43          609.06      670.46
## 6 Student_6      1 325.99 466.74 597.06          554.43      535.77
##   Literature Portuguese_Essay Math Chemistry GPA_3S
## 1      557.09          711.37 731.31      509.80 1.33333
## 2      447.23          527.58 379.14      488.64 2.98333
## 3      425.87          475.63 476.11      407.15 1.97333
## 4      521.40          592.41 783.76      588.26 2.53333
```

```
## 5      515.38          572.52 581.25    529.04 1.58667
## 6      717.03          477.60 503.82    422.92 1.66667
```

```
summary(data)
```

```
##      Student_ID      Gender      Physics      Biology
## Length:43303      Min.      :0.0000      Min.      :299.3      Min.      :263.0
## Class :character   1st Qu.:0.0000      1st Qu.:482.8      1st Qu.:492.4
## Mode  :character   Median :1.0000      Median :565.6      Median :566.4
##                                     Mean  :0.5158      Mean   :576.1      Mean   :568.7
##                                     3rd Qu.:1.0000      3rd Qu.:662.8      3rd Qu.:634.8
##                                     Max.   :1.0000      Max.   :952.1      Max.   :966.6
##      History      Second_Language      Geography      Literature
## Min.      :265.0      Min.      :222.7      Min.      :224.9      Min.      :239.1
## 1st Qu.:516.1      1st Qu.:517.7      1st Qu.:510.2      1st Qu.:516.8
## Median :578.9      Median :580.3      Median :575.5      Median :587.1
## Mean   :580.8      Mean   :574.0      Mean   :574.5      Mean   :583.3
## 3rd Qu.:650.2      3rd Qu.:640.6      3rd Qu.:637.3      3rd Qu.:648.7
## Max.   :925.8      Max.   :858.4      Max.   :941.8      Max.   :904.8
## Portuguese_Essay      Math      Chemistry      GPA_3S
## Min.      :151.6      Min.      : 298.0      Min.      : 300.5      Min.      :0.000
## 1st Qu.:491.9      1st Qu.: 489.4      1st Qu.: 484.5      1st Qu.:2.280
## Median :553.6      Median : 571.9      Median : 565.5      Median :2.920
## Mean   :551.0      Mean   : 579.2      Mean   : 571.7      Mean   :2.786
## 3rd Qu.:613.1      3rd Qu.: 665.2      3rd Qu.: 655.4      3rd Qu.:3.430
## Max.   :825.5      Max.   :1072.1      Max.   :1001.9      Max.   :4.000
```

```
sum(is.na(data))
```

```
## [1] 0
```

Data is clean, so we can proceed with extracting information from it.

Let us start with dplyr.

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
##
## The following objects are masked from 'package:stats':
##
##      filter, lag
##
## The following objects are masked from 'package:base':
##
##      intersect, setdiff, setequal, union
```

Special note: in the same fashion as with scalars, you can assign the output of the pipe to a variable or just print it out.

```
data %>%
  select(Student_ID, GPA_3S) %>%
  sample_n(5)
```

```
##      Student_ID GPA_3S
## 1 Student_21102 2.05667
## 2 Student_11481 2.39667
## 3 Student_20205 2.14000
```

```
## 4 Student_7281 2.08333
## 5 Student_1848 3.54333

my_sample <-data %>%
  select(Student_ID, GPA_3S) %>%
  sample_n(5)
```

```
my_sample
```

```
##      Student_ID GPA_3S
## 1 Student_33691 1.93333
## 2 Student_23273 3.48667
## 3 Student_8254 3.35667
## 4 Student_35314 4.00000
## 5 Student_23218 3.49000
```

```
my_sample <-data %>%
  select(Student_ID, GPA_3S) %>%
  sample_n(5)
```

```
my_sample
```

```
##      Student_ID GPA_3S
## 1 Student_40997 2.94333
## 2 Student_5008 1.61000
## 3 Student_31515 3.35000
## 4 Student_6296 3.39000
## 5 Student_17876 3.45333
```

1. Select.

1a. Select physics, chemistry and math columns. Print out first 4 rows of resulting dataframe. What is the dimensions of the dataframe?

```
#Your code goes here
```

1b. Select all columns but Mean_GPA_3S. Print out first 4 rows of resulting dataframe. Can you think about a second way to write this code (if you don't know, google it!)? Which way you prefer?

```
#Your code goes here
```

2. Filter.

2a. Filter the dataframe so that it contains only male students. How many male students are in the dataframe? Female students?

```
#Your code goes here
```

2b. Filter the dataframe so that it has only students with math exam scores above 1000. What are their IDs?

```
#Your code goes here
```

2c. Filter the dataframe so that it has only students with both math and Portuguese essay exams scores above 800. What are their IDs?

```
#Your code goes here
```

3. Mutate.

3a. Create a new column and record sum of chemistry and biology scores. Name it 'life_science_score'.

#Your code goes here

3b. Now update the values of 'life_science_score' column by dividing each value by 2. What summary statistic did you get?

Remember, you can reuse the code above by assigning it to a variable or write the code from scratch.

#Your code goes here

3. Mutate. 3c. Create 'life_science_score_1' but this time calculate mean using rowMeans() function wrapped around cbind() function. To see documentation and examples type ?rowMeans

Check if you get the same result. How would you approach this?

#Your code goes here

3d. Create 'high_gpa' column that says 'You are good!' if student's GPA is above 3.0 and 'Work harder!' otherwise. Save result as new_data. Print the last 6 rows and check if your code worked.

You can use ifelse() function. Type ?ifelse to see the documentation.

#Your code goes here

4. Group by a variable and summarize.

4a. Using the dataframe you created in the previous example, group by the new column - high_gpa. How many students have to work harder?

#Your code goes here

4b. Using the dataframe you created in the previous example, group by the new column - high_gpa. Calculate mean GPA for every group? Does the result make sense?

#Your code goes here

4c. What is the median GPA for students with high GPA?

#Your code goes here

4c. What is the median GPA for students with high GPA?

#Your code goes here

4d. Who are the students with GPA exactly equal median GPA? Lets collect them in a new dataframe named 'median_club'

#Your code goes here

4e. Group 'median_club' by gender and check the mean exam score for every discipline. Talk the results through.

#Your code goes here

Let's see the distribution of GPA in our Brazil University.

#Your code goes here