

# STAT 33B Workbook 13

Ming Fong (3035619833)

Dec 3, 2020

This workbook is due **Dec 3, 2020** by 11:59pm PT.

The workbook is organized into sections that correspond to the lecture videos for the week. Watch a video, then do the corresponding exercises *before* moving on to the next video.

Workbooks are graded for completeness, so as long as you make a clear effort to solve each problem, you'll get full credit. That said, make sure you understand the concepts here, because they're likely to reappear in homeworks, quizzes, and later lectures.

As you work, write your answers in this notebook. Answer questions with complete sentences, and put code in code chunks. You can make as many new code chunks as you like.

In the notebook, you can run the line of code where the cursor is by pressing **Ctrl + Enter** on Windows or **Cmd + Enter** on Mac OS X. You can run an entire code chunk by clicking on the green arrow in the upper right corner of the code chunk.

Please do not delete the exercises already in this notebook, because it may interfere with our grading tools.

You need to submit your work in two places:

- Submit this Rmd file with your edits on bCourses.
- Knit and submit the generated PDF file on Gradescope.

If you have any last-minute trouble knitting, **DON'T PANIC**. Submit your Rmd file on time and follow up in office hours or on Piazza to sort out the PDF.

## Relational Data

Watch the "Relational Data" lecture video.

No exercises for this section.

## Joins (with dplyr)

Watch the "Joins (with dplyr)" lecture video.

In this workbook, you'll use three tables from the Internet Movie Database (IMDB) to practice joining data frames and taking subsets.

The following is a description of the three tables and their columns.

- Titles (`titles10s.rds`), where each row is one movie. Columns:
  - `tconst` - alphanumeric unique identifier of the title
  - `titleType` - the type of the title (movie or tvMovie)
  - `primaryTitle` - title used by the filmmakers at the point of release
  - `originalTitle` - original title, in the original language
  - `startYear` - release year of the title

- runtimeMinutes - primary runtime of the title, in minutes
- Cast (cast10s.rds), where each row is one cast member from a movie. Columns:
  - tconst - alphanumeric unique identifier of the title
  - ordering - a number to uniquely identify rows for a given tconst
  - nconst - alphanumeric unique identifier of the name/person
  - category - the category of job that person was in
  - job - the specific job title if applicable, else NA
- People (people10s.rds), where each row is one person. Columns:
  - nconst - alphanumeric unique identifier of the name/person
  - primaryName - name by which the person is most often credited

The tables are a subset of IMDB's larger collection of data, which is available at <https://www.imdb.com/interfaces/>.

## Exercise 1

Find the names of the three people that had the most roles in movies in the data set.

*Hint 1: Think about what the rows in each table represent. Is there a table where rows represent appearances/roles in movies?*

*Hint 2: A join is only needed here to get the names of the people.*

**YOUR ANSWER GOES HERE:**

```
library(dplyr)
cast = readRDS("data/imdb/cast10s.rds")
people = readRDS("data/imdb/people10s.rds")
titles = readRDS("data/imdb/titles10s.rds")

top_3 = data.frame(sort(table(cast$nconst), decreasing = TRUE)[1:3])
joined = left_join(top_3, people, by = c("Var1" = "nconst"))
joined$primaryName

## [1] "Kevin MacLeod"      "Eric Roberts"      "William Shakespeare"
```

## Exercise 2

Compute a data frame that contains the nconst IDs, names, and roles of the primary cast from the 2018 movie "Black Panther".

*Hint 1: A common strategy for relational data is to reduce the size of a table or get specific rows by taking a subset, and then join that table with another table. Start by taking a subset of the Titles table to find the tconst ID for Black Panther.*

*Hint 2: When you have more than two tables, it is sometimes necessary to use more than one join.*

**YOUR ANSWER GOES HERE:**

```
bp_tconst = titles[titles$primaryTitle == "Black Panther", ]$tconst
bp_cast = cast[cast$tconst == bp_tconst, ]
data.frame(left_join(bp_cast, people, by = "nconst"))

##      tconst ordering   nconst category      job
## 1 tt1825683      10 nm3234869 composer    <NA>
## 2 tt1825683       1 nm1569276   actor     <NA>
## 3 tt1825683       2 nm0430107   actor     <NA>
## 4 tt1825683       3 nm2143282 actress    <NA>
## 5 tt1825683       4 nm1775091 actress    <NA>
```

```
## 6 tt1825683      5 nm3363032 director          <NA>
## 7 tt1825683      6 nm1963288  writer          written by
## 8 tt1825683      7 nm0498278  writer based on the Marvel comics by
## 9 tt1825683      8 nm0456158  writer based on the Marvel Comics by
## 10 tt1825683     9 nm0270559 producer          producer
##      primaryName
## 1   Ludwig Göransson
## 2   Chadwick Boseman
## 3   Michael B. Jordan
## 4   Lupita Nyong'o
## 5   Danai Gurira
## 6   Ryan Coogler
## 7   Joe Robert Cole
## 8   Stan Lee
## 9   Jack Kirby
## 10  Kevin Feige
```

### Exercise 3

Compute a data frame that contains the names of the primary **actors and actresses** for all Harry Potter movies in the data set.

*Hint: The `startsWith` function (or `stringr`) is helpful for identifying Harry Potter movies.*

**YOUR ANSWER GOES HERE:**

```
hp_tconst = titles[startsWith(titles$primaryTitle, "Harry Potter"), ]$tconst
hp_cast = cast[cast$tconst == hp_tconst, ]
hp_joined = left_join(hp_cast, people, by = "nconst")
hp_joined[hp_joined$category == "actor" | hp_joined$category == "actress", ]$primaryName

## [1] "Rupert Grint"      "Daniel Radcliffe" "Michael Gambon"
```

## STAT 33 Wrap-up

Watch the “STAT 33 Wrap-up” lecture video.

Please fill in the teaching evaluations for this class!