# Harvard Data Science Professional Capstone Project - MovieLens Quizzes

Uday Adusumilli

## setup code

Create edx set, validation set (final hold-out test set)

step 1 Note: this process could take a couple of minutes

```r
if(!require(tidyverse)) install.packages("tidyverse", repos = "http://cran.us.r-project.org")
```

```
## Loading required package: tidyverse
```

```
## -- Attaching packages -------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5      v purrr   0.3.4
## v tibble  3.1.2      v dplyr   1.0.7
## v tidyr   1.1.3      v stringr 1.4.0
## v readr   1.4.0      v forcats 0.5.1
```

```
## -- Conflicts ----------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```r
if(!require(caret)) install.packages("caret", repos = "http://cran.us.r-project.org")
```

```
## Loading required package: caret
```

```
## Warning: package 'caret' was built under R version 4.1.1
```

```
## Loading required package: lattice
```

```
##
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:purrr':
##
##     lift
```

```
if(!require(data.table)) install.packages("data.table", repos = "http://cran.us.r-project.org")
```

```
## Loading required package: data.table
```

```
##
## Attaching package: 'data.table'
```

```
## The following objects are masked from 'package:dplyr':
##
##     between, first, last
```

```
## The following object is masked from 'package:purrr':
##
##     transpose
```

```
library(tidyverse)
library(caret)
library(data.table)
```

The following libraries were needed, and they were not mentioned in the code provided to us

```
library(dplyr)
library(tidyverse)
library(kableExtra)
```

```
## Warning: package 'kableExtra' was built under R version 4.1.1
```

```
##
## Attaching package: 'kableExtra'
```

```
## The following object is masked from 'package:dplyr':
##
##     group_rows
```

```
library(tidyr)
library(stringr)
library(forcats)
library(ggplot2)
```

```
# MovieLens 10M dataset:
# https://grouplens.org/datasets/movielens/10m/
# http://files.grouplens.org/datasets/movielens/ml-10m.zip

dl <- tempfile()
download.file("http://files.grouplens.org/datasets/movielens/ml-10m.zip", dl)

ratings <- fread(text = gsub("::", "\t", readLines(unzip(dl, "ml-10M100K/ratings.dat"))),
                 col.names = c("userId", "movieId", "rating", "timestamp"))

movies <- str_split_fixed(readLines(unzip(dl, "ml-10M100K/movies.dat")), "\\::", 3)
colnames(movies) <- c("movieId", "title", "genres")
```

the following block will depend on the r version

```
# if using R 3.6 or earlier:
# movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(levels(movieId))[movieId],
#                                            title = as.character(title),
#                                            genres = as.character(genres))
# if using R 4.0 or later:
movies <- as.data.frame(movies) %>% mutate(movieId = as.numeric(movieId),
                                           title = as.character(title),
                                           genres = as.character(genres))
```

the following block is version independent

```
movielens <- left_join(ratings, movies, by = "movieId")

# Validation set will be 10% of MovieLens data
set.seed(1, sample.kind="Rounding") # if using R 3.5 or earlier, use `set.seed(1)`
```

```
## Warning in set.seed(1, sample.kind = "Rounding"): non-uniform 'Rounding' sampler
## used
```

```
test_index <- createDataPartition(y = movielens$rating, times = 1, p = 0.1, list = FALSE)
edx <- movielens[-test_index,]
temp <- movielens[test_index,]

# Make sure userId and movieId in validation set are also in edx set
validation <- temp %>%
      semi_join(edx, by = "movieId") %>%
      semi_join(edx, by = "userId")

# Add rows removed from validation set back into edx set
removed <- anti_join(temp, validation)
```

```
## Joining, by = c("userId", "movieId", "rating", "timestamp", "title", "genres")
```

```
edx <- rbind(edx, removed)
```

```
rm(dl, ratings, movies, test_index, temp, movielens, removed)
```

## Actual Quizz Questions

### Q1

How many rows and columns are there in the edx dataset?

```
dim(edx)
```

```
## [1] 9000055       6
```

**Q2**

How many zeros were given as ratings in the edx dataset?

```
edx %>% filter(rating == 0) %>% tally()
```

```
##   n
## 1 0
```

How many threes were given as ratings in the edx dataset?

```
edx %>% filter(rating == 3) %>% tally()
```

```
##         n
## 1 2121240
```

**Q3**

How many different movies are in the edx dataset?

```
n_distinct(edx$movieId)
```

```
## [1] 10677
```

**Q4**

How many different users are in the edx dataset?

```
n_distinct(edx$userId)
```

```
## [1] 69878
```

**Q5**

How many movie ratings are in each of the following genres in the edx dataset?

```
edx %>% separate_rows(genres, sep = "\\|") %>%
  group_by(genres) %>%
  summarize(count = n()) %>%
  arrange(desc(count))
```

```
## # A tibble: 20 x 2
##    genres          count
##    <chr>           <int>
##  1 Drama         3910127
##  2 Comedy        3540930
##  3 Action        2560545
##  4 Thriller      2325899
##  5 Adventure     1908892
```

4

```
##  6 Romance             1712100
##  7 Sci-Fi              1341183
##  8 Crime               1327715
##  9 Fantasy              925637
## 10 Children             737994
## 11 Horror               691485
## 12 Mystery              568332
## 13 War                  511147
## 14 Animation            467168
## 15 Musical              433080
## 16 Western              189394
## 17 Film-Noir            118541
## 18 Documentary           93066
## 19 IMAX                   8181
## 20 (no genres listed)        7
```

**Q6**

Which movie has the greatest number of ratings?

```
edx %>% group_by(movieId, title) %>%
  summarize(count = n()) %>%
  arrange(desc(count))
```

```
## 'summarise()' has grouped output by 'movieId'. You can override using the '.groups' argument.
```

```
## # A tibble: 10,677 x 3
## # Groups:   movieId [10,677]
##    movieId title                                                          count
##      <dbl> <chr>                                                          <int>
##  1     296 Pulp Fiction (1994)                                            31362
##  2     356 Forrest Gump (1994)                                            31079
##  3     593 Silence of the Lambs, The (1991)                               30382
##  4     480 Jurassic Park (1993)                                           29360
##  5     318 Shawshank Redemption, The (1994)                               28015
##  6     110 Braveheart (1995)                                              26212
##  7     457 Fugitive, The (1993)                                           25998
##  8     589 Terminator 2: Judgment Day (1991)                              25984
##  9     260 Star Wars: Episode IV - A New Hope (a.k.a. Star Wars) (1977) 25672
## 10     150 Apollo 13 (1995)                                               24284
## # ... with 10,667 more rows
```

**Q7**

What are the five most given ratings in order from most to least?

```
edx %>% group_by(rating) %>% summarize(count = n()) %>% top_n(5) %>%
  arrange(desc(count))
```
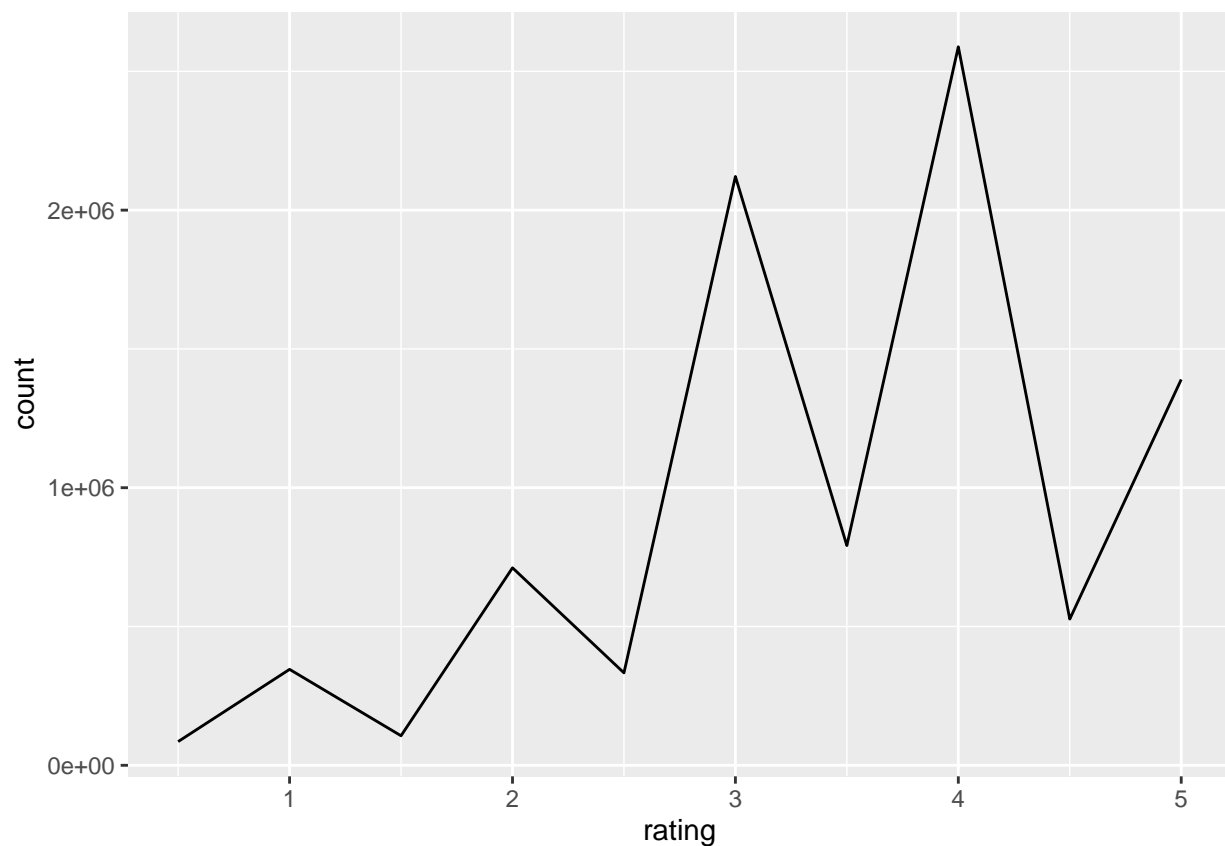
```
## Selecting by count
```

5

```
## # A tibble: 5 x 2
##    rating    count
##     <dbl>    <int>
## 1       4  2588430
## 2       3  2121240
## 3       5  1390114
## 4     3.5   791624
## 5       2   711422
```

**Q8**

True or False: In general, half star ratings are less common than whole star ratings (e.g., there are fewer
ratings of 3.5 than there are ratings of 3 or 4, etc.).

```
edx %>%
  group_by(rating) %>%
  summarize(count = n()) %>%
  ggplot(aes(x = rating, y = count)) +
  geom_line()
```

# h1

## h2

### h3

#### h4  regular line