

# Assignment02

Ted Kim

2022-09-09

## OVERVIEW

Choose six recent popular movies and ask at least five friends and family members who know to rate each movie they have seen from 1 to 5. Get the results, save them to the My SQL database, and load the information in the SQL database into the R data frame. Assuming that there is no existing database, I created a database and a table. The query generating the database and table was inserted as hardcoding in the markdown. The movie-related data was downloaded from themoviedb.org as a json file via api. Json files are registered on GitHub. It can be implemented simply by using the overwrite property of the dbWriteTable() function, but since it is not possible to create a relationship between tables, the append property will be used. Therefore, before each table is stored, duplicate data is discarded by comparing it with the existing table and json/Google sheet.

## CONNECT MySQL

Create a connection object to MySQL database.

```
con = dbConnect(RMySQL::MySQL(),
                user='root',
                password='',
                host='localhost')
```

## CREATE DATABASE

In general, databases and tables are often already created. However, if a table is not created, we must create a table and put data in it when the program runs. Create database *movies* if not exists, and use database *movies*

```
res <- dbSendQuery(con, 'create database if not exists movies;')

res <- dbSendQuery(con, 'use movies;')
```

## CREATE TABLES

Create tables *genres*, *movies*, *movies\_genres* if not exists

```
res <- dbSendQuery(con, 'create table if not exists genres (
                        id int primary key,
                        description varchar(40)
                        ) engine=innodb;')

res <- dbSendQuery(con, 'create table if not exists movies (
                        id int auto_increment primary key,
                        adult bool,
                        backdrop_path varchar(255),
                        original_language varchar(10),
```

```

        original_title varchar(255),
        overview varchar(1000),
        popularity double(8,3),
        poster_path varchar(255),
        release_date date,
        title varchar(255),
        video bool,
        vote_average double(8,3),
        vote_count int
    ) engine=innodb;')

res <- dbSendQuery(con, 'create table if not exists movies_genres (
    id int auto_increment primary key,
    movie_id int,
    genre_id int,
    foreign key (movie_id) references movies (id),
    foreign key (genre_id) references genres (id)
) engine=innodb;')

res <- dbSendQuery(con, 'create table if not exists survey_result (
    id int auto_increment primary key,
    survey_id int,
    movie_id int,
    email_address varchar(255),
    rate int,
    registered datetime,
    foreign key (movie_id) references movies (id)
) engine=innodb;')

```

Displays the tables available in the database

```
## [1] "genres"          "movies"          "movies_genres" "survey_result"
```

## LOAD GENRES DATA

Stores genre data imported from TMDB into the table on SQL Server. Since the same data may exist in the table on SQL Server, only subsets of the two data that do not exist in the table on SQL Server are appended.

**Get genres data from SQL server**

```
dfGenres <- fetch(dbSendQuery(con, 'select * from genres'), )
dbClearResult(dbListResults(con)[[1]])
```

```
[1] TRUE
```

**Get genres.json from github to load**

```
file <- 'https://raw.githubusercontent.com/blacksmilez/DATA607/main/Assignment02/json/genres.json'
jsonGenres <- fromJSON(file)
```

**Retrieve all rows on data frame *dfGenres*, *jsonGenres* to verify**

```
print(dfGenres[order(dfGenres$id),, row.names = FALSE, right = FALSE])
```

```
## id    description
##   12 Adventure
##   14 Fantasy
##   16 Animation
```

```
##      18 Drama
##      27 Horror
##      28 Action
##      35 Comedy
##      36 History
##      37 Western
##      53 Thriller
##      80 Crime
##      99 Documentary
##     878 Science Fiction
##    9648 Mystery
##   10402 Music
##   10749 Romance
##   10751 Family
##   10752 War
##   10770 TV Movie
```

```
print(jsonGenres[order(jsonGenres$id),], row.names = FALSE, right = FALSE)
```

```
## id      description
##    12 Adventure
##    14 Fantasy
##    16 Animation
##    18 Drama
##    27 Horror
##    28 Action
##    35 Comedy
##    36 History
##    37 Western
##    53 Thriller
##    80 Crime
##    99 Documentary
##   878 Science Fiction
##  9648 Mystery
## 10402 Music
## 10749 Romance
## 10751 Family
## 10752 War
## 10770 TV Movie
```

Get Subsets of the two data that do not exist in the table on SQL Server

```
## [1] id      description
## <0 rows> (or 0-length row.names)
```

Append subsets into the table on SQL Server

```
dbWriteTable(con, 'genres', dfSubsets[, c('id', 'description')], row.names=FALSE, append=TRUE)
```

```
## [1] TRUE
```

## LOAD MOVIES DATA

Stores movies data imported from TMDb into the table on SQL Server. Since the same data may exist in the table on SQL Server, only subsets of the two data that do not exist in the table on SQL Server are appended.

Get movies data from SQL server

```
dfMovies <- fetch(dbSendQuery(con, 'select * from movies'), )
dbClearResult(dbListResults(con)[[1]])
```

```
[1] TRUE
```

Get movies.json from github to load

```
file <- 'https://raw.githubusercontent.com/blacksmilez/DATA607/main/Assignment02/json/movies.json'
jsonMovies <- fromJSON(file)
```

Retrieve all rows on data frame *dfMovies*, *jsonMovies* to verify

```
print(dfMovies[order(dfMovies$id), c('id', 'title', 'release_date')],
      row.names = FALSE, right = FALSE)
```

```
## id          title          release_date
## 361743 Top Gun: Maverick      2022-05-24
## 453395 Doctor Strange in the Multiverse of Madness 2022-05-04
## 507086 Jurassic World Dominion 2022-06-01
## 539681 DC League of Super-Pets 2022-07-27
## 616037 Thor: Love and Thunder 2022-07-06
## 629176 Samaritan             2022-08-25
## 634649 Spider-Man: No Way Home 2021-12-15
## 755566 Day Shift             2022-08-10
## 766507 Prey                  2022-08-02
## 848123 Black Site            2022-05-05
## 927341 Hunting Ava Bravo      2022-04-01
## 951368 Your Boyfriend Is Mine 2022-03-19
## 997120 Sniper: Rogue Mission  2022-08-15
## 1006851 Office Invasion        2022-08-10
## 1008779 The princess          2022-08-05
```

```
print(jsonMovies[order(jsonMovies$id), c('id', 'title', 'release_date')],
      row.names = FALSE, right = FALSE)
```

```
## id          title          release_date
## 361743 Top Gun: Maverick      2022-05-24
## 453395 Doctor Strange in the Multiverse of Madness 2022-05-04
## 507086 Jurassic World Dominion 2022-06-01
## 539681 DC League of Super-Pets 2022-07-27
## 616037 Thor: Love and Thunder 2022-07-06
## 629176 Samaritan             2022-08-25
## 634649 Spider-Man: No Way Home 2021-12-15
## 755566 Day Shift             2022-08-10
## 766507 Prey                  2022-08-02
## 848123 Black Site            2022-05-05
## 927341 Hunting Ava Bravo      2022-04-01
## 951368 Your Boyfriend Is Mine 2022-03-19
## 997120 Sniper: Rogue Mission  2022-08-15
## 1006851 Office Invasion        2022-08-10
## 1008779 The princess          2022-08-05
```

Get Subsets of the two data that do not exist in the table on SQL Server

```
## [1] id          title          release_date
## <0 rows> (or 0-length row.names)
```

Append subsets into the table on SQL Server

```
dbWriteTable(con, 'movies',
             dfSubsets[, c('id', 'adult', 'backdrop_path', 'original_language',
                           'original_title', 'overview', 'popularity', 'title',
                           'poster_path', 'release_date', 'video')],
             row.names=FALSE, append=TRUE)
```

```
## [1] TRUE
```

## LOAD MOVIES-GENRES DATA

Stores movies\_genres data imported from TMDb into the table on SQL Server. Since the same data may exist in the table on SQL Server, only subsets of the two data that do not exist in the table on SQL Server are appended.

Get movies\_genres data from SQL server

```
dfMoviesGenres <- fetch(dbSendQuery(con, 'select * from movies_genres'), )
dbClearResult(dbListResults(con)[[1]])
```

```
[1] TRUE
```

Get movies\_genres.json from github to load

```
file <- 'https://raw.githubusercontent.com/blacksmilez/DATA607/main/Assignment02/json/movies_genres.json'
jsonMoviesGenres <- fromJSON(file)
```

Retrieve all rows on data frame *dfMoviesGenres*, *jsonMoviesGenres* to verify

```
print(dfMoviesGenres[order(dfMoviesGenres$movie_id, dfMoviesGenres$genre_id),],
      row.names = FALSE, right = FALSE)
```

```
##  id movie_id genre_id
##  18  361743      18
##  17  361743      28
##  39  453395      12
##  37  453395      14
##  38  453395      28
##  14  507086      12
##  15  507086      28
##  16  507086     878
##   9  539681      16
##  10  539681      28
##  13  539681      35
##  12  539681     878
##  11  539681    10751
##   5  616037      12
##   6  616037      14
##   4  616037      28
##   2  629176      18
##  45  629176      18
##   1  629176      28
##  44  629176      28
##   3  629176     878
##  46  629176     878
##  41  634649      12
##  40  634649      28
##  42  634649     878
```

```
## 29 755566      14
## 30 755566      27
## 28 755566      28
## 31 755566      35
##  8 766507      28
##  7 766507      53
## 35 848123      28
## 36 848123      53
## 20 927341      27
## 21 927341      28
## 19 927341      53
## 27 951368      28
## 25 951368      53
## 26 951368 10770
## 22 997120      28
## 23 997120      53
## 24 997120     9648
## 34 1006851     28
## 33 1006851     35
## 32 1006851     878
## 43 1008779     28
```

```
print(jsonMoviesGenres[order(jsonMoviesGenres$movie_id, jsonMoviesGenres$genre_id),],
      row.names = FALSE, right = FALSE)
```

```
## movie_id genre_id
## 361743      18
## 361743      28
## 453395      12
## 453395      14
## 453395      28
## 507086      12
## 507086      28
## 507086     878
## 539681      16
## 539681      28
## 539681      35
## 539681     878
## 539681 10751
## 616037      12
## 616037      14
## 616037      28
## 629176      18
## 629176      18
## 629176      28
## 629176      28
## 629176     878
## 629176     878
## 634649      12
## 634649      28
## 634649     878
## 755566      14
## 755566      27
## 755566      28
## 755566      35
```

```
## 766507 28
## 766507 53
## 848123 28
## 848123 53
## 927341 27
## 927341 28
## 927341 53
## 951368 28
## 951368 53
## 951368 10770
## 997120 28
## 997120 53
## 997120 9648
## 1006851 28
## 1006851 35
## 1006851 878
## 1008779 28
```

**Get Subsets of the two data that do not exist in the table on SQL Server**

```
## [1] movie_id genre_id
## <0 rows> (or 0-length row.names)
```

**Append subsets into the table on SQL Server**

```
dbWriteTable(con, 'movies_genres',
             dfSubsets[, c('movie_id', 'genre_id')],
             row.names=FALSE, append=TRUE)
```

```
## [1] TRUE
```

※ The following error occurred when running “dbWriteTable()” for the first time:

“*ERROR: Loading local data is disabled - this must be enabled on both the client and server sides*”

error occurs while copying data frames to database tables using dbWriteTable(), it is handled as follows:

```
# 1. open mysql terminal
# 2. check the local_infile
# mysql> show global variables like 'local_infile'
# +-----+-----+
# | Variable_name | Value |
# +-----+-----+
# | local_infile | OFF |
# +-----+-----+
# (this means local_infile is disable)
# 3. put set command
# mysql> set global local_infile=true;
# mysql> exit
```

## LOAD SURVEY DATA

Stores survey data imported from Google Sheet into the table on SQL Server. Since the same data may exist in the table on SQL Server, only subsets of the two data that do not exist in the table on SQL Server are appended.

**Get survey result data from SQL server**

Obtain survey data from SQL Server. There is no record set returned because there is no data at the time of initial execution.

```
dfSurveyResults <- fetch(
  dbSendQuery(con, 'select * from survey_result where survey_id = 1'),)
dbClearResult(dbListResults(con)[[1]])
```

```
[1] TRUE
```

```
print(dfSurveyResults[order(dfSurveyResults$movie_id),], row.names = FALSE, right = FALSE)
```

```
id survey_id movie_id email_address rate registered
1 1 361743 blacksmilez@gmail.com 4 2022-09-11 12:19:59 2 1 361743 negativetae@gmail.com 5 2022-09-11
12:21:38 3 1 361743 nury95@hotmail.com 5 2022-09-11 12:47:01 4 1 361743 delight_32@hotmail.com 0 2022-09-
11 13:13:58 5 1 361743 s88724@gmail.com 5 2022-09-11 13:21:36 6 1 507086 blacksmilez@gmail.com 5 2022-09-11
12:19:59 7 1 507086 negativetae@gmail.com 1 2022-09-11 12:21:38 8 1 507086 nury95@hotmail.com 5 2022-09-11
12:47:01 9 1 507086 delight_32@hotmail.com 0 2022-09-11 13:13:58 10 1 507086 s88724@gmail.com 2 2022-09-11
13:21:36 16 1 539681 blacksmilez@gmail.com 3 2022-09-11 12:19:59 17 1 539681 negativetae@gmail.com 1 2022-
09-11 12:21:38 18 1 539681 nury95@hotmail.com 5 2022-09-11 12:47:01 19 1 539681 delight_32@hotmail.com 4
2022-09-11 13:13:58 20 1 539681 s88724@gmail.com 3 2022-09-11 13:21:36 21 1 629176 blacksmilez@gmail.com
2 2022-09-11 12:19:59 22 1 629176 negativetae@gmail.com 3 2022-09-11 12:21:38 23 1 629176 nury95@
hotmail.com 5 2022-09-11 12:47:01 24 1 629176 delight_32@hotmail.com 4 2022-09-11 13:13:58 25 1 629176
s88724@gmail.com 4 2022-09-11 13:21:36 11 1 634649 blacksmilez@gmail.com 4 2022-09-11 12:19:59 12 1
634649 negativetae@gmail.com 2 2022-09-11 12:21:38 13 1 634649 nury95@hotmail.com 3 2022-09-11 12:47:01
14 1 634649 delight_32@hotmail.com 3 2022-09-11 13:13:58 15 1 634649 s88724@gmail.com 4 2022-09-11
13:21:36 26 1 755566 blacksmilez@gmail.com 0 2022-09-11 12:19:59 27 1 755566 negativetae@gmail.com 4 2022-
09-11 12:21:38 28 1 755566 nury95@hotmail.com 3 2022-09-11 12:47:01 29 1 755566 delight_32@hotmail.com
4 2022-09-11 13:13:58 30 1 755566 s88724@gmail.com 3 2022-09-11 13:21:36
```

### Get survey result from google sheet to load

The survey was conducted with Google Forms that can be easily used.

hyperlink: <https://docs.google.com/forms/d/e/1FAIpQLSeL4Ymj956wxJ9rMH-ie-XHgm6P-d25iHjvxAyNmKc7QIvIg/viewform>

It brings up the Google sheet where the data input through Google Form is stored.

```
gs4_deauth()
file <- 'https://docs.google.com/spreadsheets/d/1n8U9Ab0SKMI871oHoycPK-WXcKkC3_VmfdSx742hbTo/edit?usp=s
df <- as.data.frame(read_sheet(file))
```

### Get survey result from google sheet to load

The ID of the movie used for the survey is included in the array. A separate metric table should be created, but this time it will be omitted.

```
movies_id <- c(361743, 507086, 634649, 539681, 629176, 755566)
```

Creates an empty data frame for storing subsets.

```
dfSurveySubsets <- data.frame(matrix(ncol=5, nrow=0))
colnames(dfSurveySubsets) <- c('survey_id', 'movie_id', 'email_address', 'rate', 'registered')
```

Only new survey data that does not exist in the existing survey table is selected from the Google sheet.

```
count <- 1
for(id in movies_id) {
  dftmp <- df[, c('survey_id', paste0('movie', count), 'email_address', 'registered')]
  names(dftmp)[names(dftmp) == paste0('movie', count)] <- 'rate'
  dftmp['movie_id'] <- id
  dfSubsets <- subset(dftmp, !(movie_id %in% dfSurveyResults$movie_id
                                &&email_address %in% dfSurveyResults$email_address))
  dfSurveySubsets <- rbind(dfSurveySubsets, dfSubsets)
  count = count + 1
}
```



```
}
print(dfSurveySubsets)
```

```
[1] survey_id movie_id email_address rate registered
<0 rows> (or 0-length row.names)
```

### Append subsets into the table on SQL Server

Only non-duplicated data is stored in the SQL Server table.

```
dbWriteTable(con, 'survey_result',
             dfSurveySubsets,
             row.names=FALSE, append=TRUE)
```

```
## [1] TRUE
```

### Re-Get survey result data from SQL server

The survey data is retrieved from the SQL server again.

```
dfSurveyResults <- fetch(dbSendQuery(con, 'select * from survey_result where survey_id = 1'),)
dbClearResult(dbListResults(con)[[1]])
```

```
[1] TRUE
```

Make the column names the same for join between two data frames.

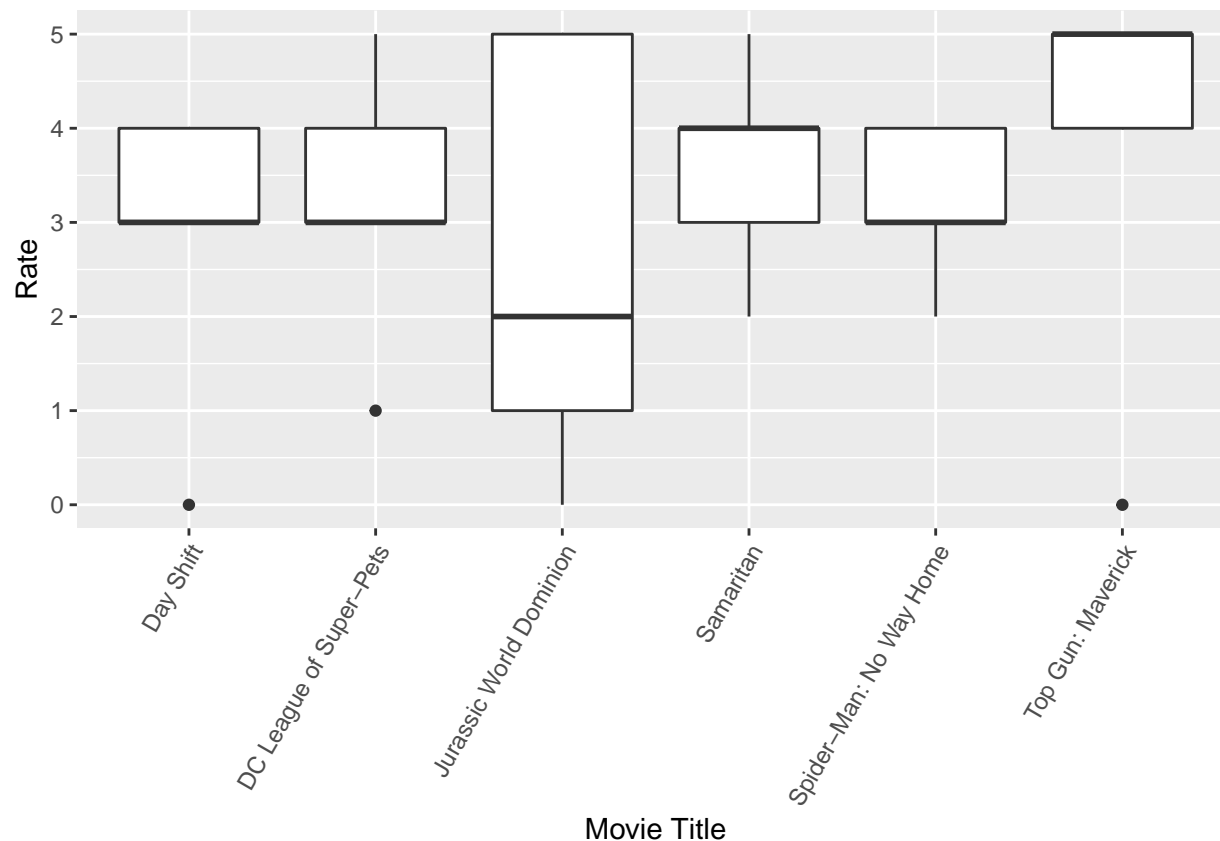
```
names(jsonMovies)[names(jsonMovies) == 'id'] <- 'movie_id'
```

Merge the two data tables and remove unnecessary columns.

```
dfResults <- merge(dfSurveyResults, jsonMovies, by = 'movie_id')
dftmp <- subset(dfResults, select = -c(registered, adult, backdrop_path, original_language,
                                     overview, original_title, poster_path, release_date,
                                     video, id, survey_id))
```

### Graphs drawn without calibration of missing data

If the graph is drawn without calibration of missing data as follows.



## Missing Data (1)

For calibration, the missing values can be filled with *mean* values.

```
colnames(dftmp)
```

```
## [1] "movie_id"      "email_address" "rate"          "popularity"
## [5] "title"
```

```
dfmean <- dftmp %>%
  group_by(movie_id) %>%
  mutate(mean = mean(rate))
```

```
dfmean$rate <- ifelse(dfmean$rate == 0, dfmean$mean, dfmean$rate)
```

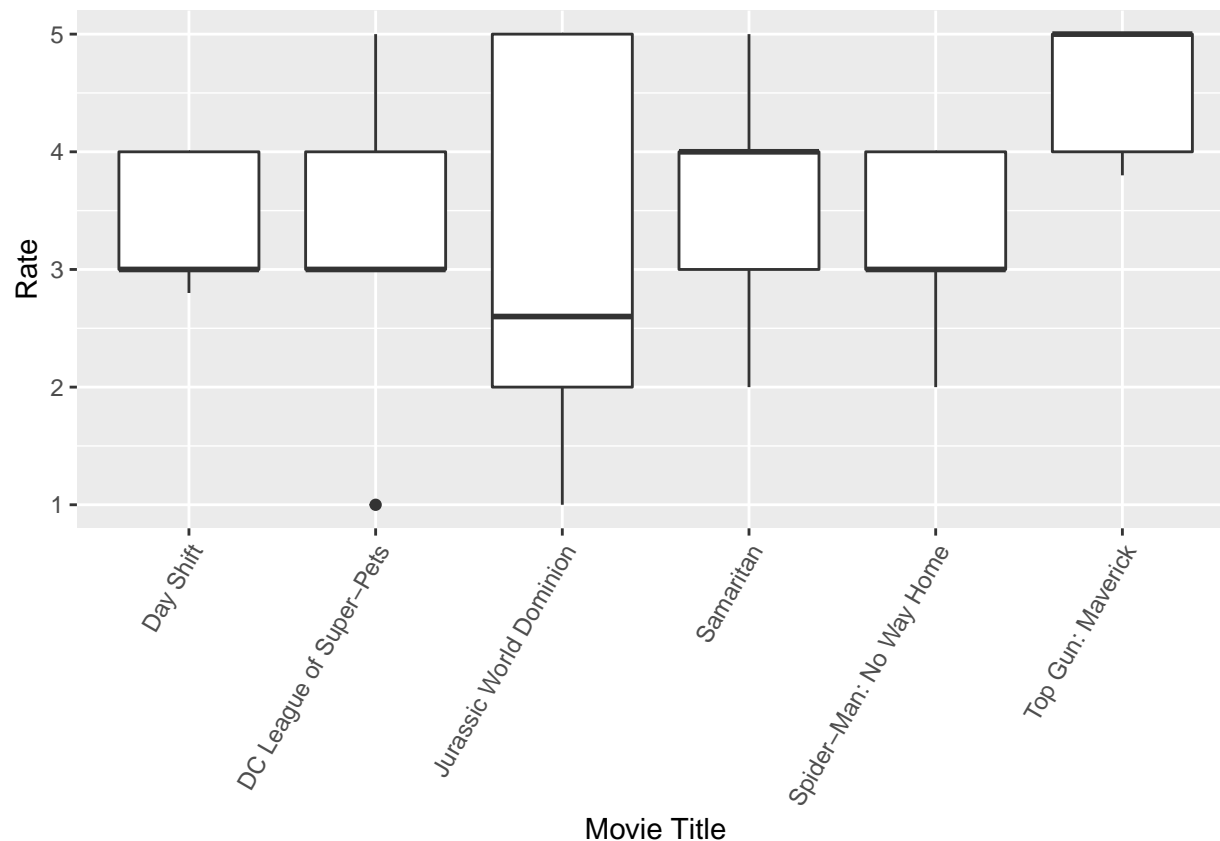
```
print(dfmean, n=100)
```

```
## # A tibble: 30 x 6
## # Groups:   movie_id [6]
##   movie_id email_address      rate popularity title      mean
##   <int> <chr>          <dbl>      <dbl> <chr>      <dbl>
## 1 361743 blacksmilez@gmail.com 4      2030. 2030. Top Gun: Maverick 3.8
## 2 361743 negativetae@gmail.com 5      2030. 2030. Top Gun: Maverick 3.8
## 3 361743 nury95@hotmail.com 5      2030. 2030. Top Gun: Maverick 3.8
## 4 361743 delight_32@hotmail.com 3.8    2030. 2030. Top Gun: Maverick 3.8
## 5 361743 s88724@gmail.com 5      2030. 2030. Top Gun: Maverick 3.8
```

##	6	507086	blacksmilez@gmail.com	5	2084.	Jurassic World Domini~	2.6
##	7	507086	negativetae@gmail.com	1	2084.	Jurassic World Domini~	2.6
##	8	507086	nury95@hotmail.com	5	2084.	Jurassic World Domini~	2.6
##	9	507086	delight_32@hotmail.com	2.6	2084.	Jurassic World Domini~	2.6
##	10	507086	s88724@gmail.com	2	2084.	Jurassic World Domini~	2.6
##	11	539681	blacksmilez@gmail.com	3	2738.	DC League of Super-Pe~	3.2
##	12	539681	negativetae@gmail.com	1	2738.	DC League of Super-Pe~	3.2
##	13	539681	nury95@hotmail.com	5	2738.	DC League of Super-Pe~	3.2
##	14	539681	delight_32@hotmail.com	4	2738.	DC League of Super-Pe~	3.2
##	15	539681	s88724@gmail.com	3	2738.	DC League of Super-Pe~	3.2
##	16	629176	blacksmilez@gmail.com	2	5115.	Samaritan	3.6
##	17	629176	negativetae@gmail.com	3	5115.	Samaritan	3.6
##	18	629176	nury95@hotmail.com	5	5115.	Samaritan	3.6
##	19	629176	delight_32@hotmail.com	4	5115.	Samaritan	3.6
##	20	629176	s88724@gmail.com	4	5115.	Samaritan	3.6
##	21	634649	blacksmilez@gmail.com	4	1142.	Spider-Man: No Way Ho~	3.2
##	22	634649	negativetae@gmail.com	2	1142.	Spider-Man: No Way Ho~	3.2
##	23	634649	nury95@hotmail.com	3	1142.	Spider-Man: No Way Ho~	3.2
##	24	634649	delight_32@hotmail.com	3	1142.	Spider-Man: No Way Ho~	3.2
##	25	634649	s88724@gmail.com	4	1142.	Spider-Man: No Way Ho~	3.2
##	26	755566	blacksmilez@gmail.com	2.8	1403.	Day Shift	2.8
##	27	755566	negativetae@gmail.com	4	1403.	Day Shift	2.8
##	28	755566	nury95@hotmail.com	3	1403.	Day Shift	2.8
##	29	755566	delight_32@hotmail.com	4	1403.	Day Shift	2.8
##	30	755566	s88724@gmail.com	3	1403.	Day Shift	2.8

Graph of missing values filled with *mean* values

```
ggplot(dfmean, aes(x=title, y=rate)) +
  geom_boxplot() +
  theme(axis.text.x = element_text(angle = 60, hjust = 1, vjust = 1.0 )) +
  labs(x='Movie Title', y='Rate')
```



## Missing Data (2)

For calibration, the missing values can be filled with *median* values.

```
dfmedian <- dftmp %>%
  group_by(movie_id) %>%
  mutate(median = median(rate))

dfmedian$rate <- ifelse(dfmedian$rate == 0, dfmedian$median, dfmedian$rate)

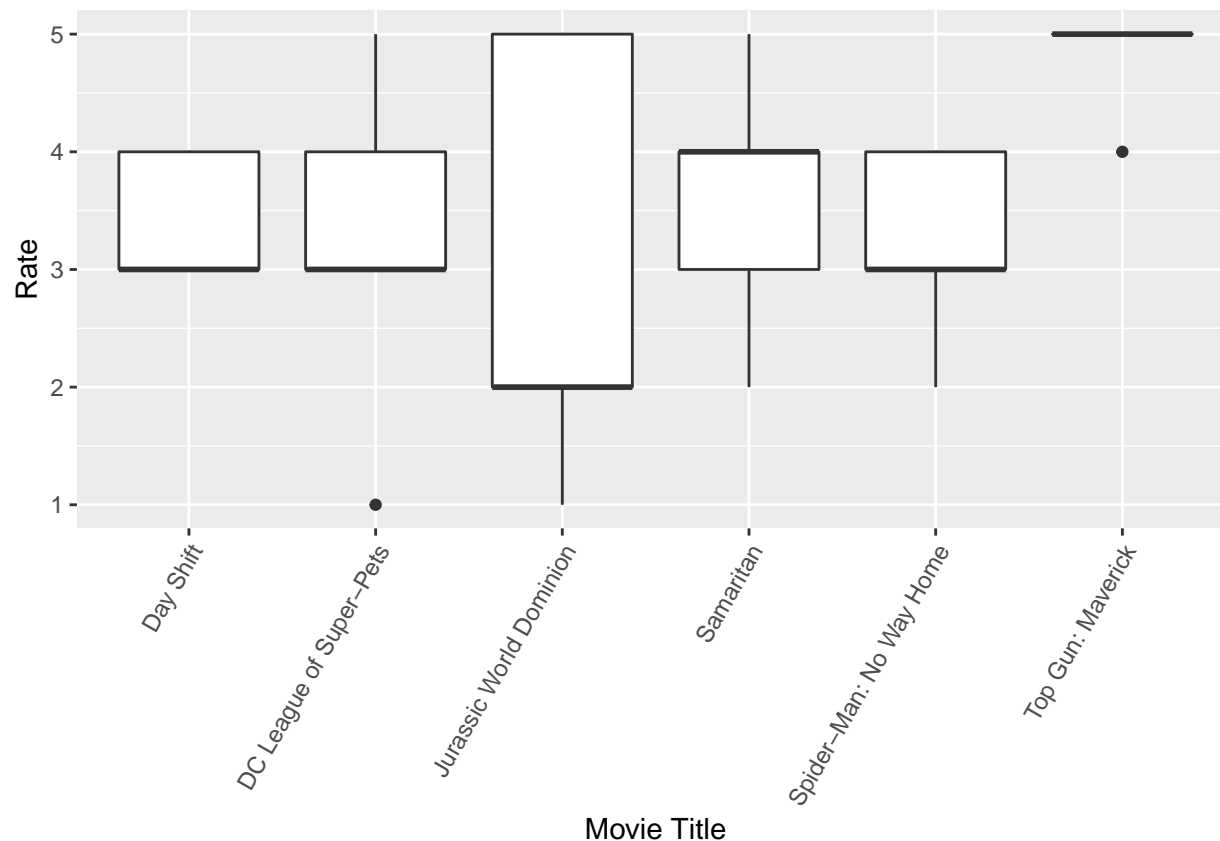
print(dfmedian, n=100)
```

```
## # A tibble: 30 x 6
## # Groups:   movie_id [6]
##   movie_id email_address      rate popularity title      median
##   <int> <chr>          <int>    <dbl> <chr>      <int>
## 1 361743 blacksmilez@gmail.com      4    2030. Top Gun: Maverick      5
## 2 361743 negativetae@gmail.com      5    2030. Top Gun: Maverick      5
## 3 361743 nury95@hotmail.com        5    2030. Top Gun: Maverick      5
## 4 361743 delight_32@hotmail.com      5    2030. Top Gun: Maverick      5
## 5 361743 s88724@gmail.com           5    2030. Top Gun: Maverick      5
## 6 507086 blacksmilez@gmail.com      5    2084. Jurassic World Domin~      2
## 7 507086 negativetae@gmail.com      1    2084. Jurassic World Domin~      2
## 8 507086 nury95@hotmail.com        5    2084. Jurassic World Domin~      2
## 9 507086 delight_32@hotmail.com      2    2084. Jurassic World Domin~      2
## 10 507086 s88724@gmail.com          2    2084. Jurassic World Domin~      2
```

## 11	539681	blacksmilez@gmail.com	3	2738. DC League of Super-P~	3
## 12	539681	negativetae@gmail.com	1	2738. DC League of Super-P~	3
## 13	539681	nury95@hotmail.com	5	2738. DC League of Super-P~	3
## 14	539681	delight_32@hotmail.com	4	2738. DC League of Super-P~	3
## 15	539681	s88724@gmail.com	3	2738. DC League of Super-P~	3
## 16	629176	blacksmilez@gmail.com	2	5115. Samaritan	4
## 17	629176	negativetae@gmail.com	3	5115. Samaritan	4
## 18	629176	nury95@hotmail.com	5	5115. Samaritan	4
## 19	629176	delight_32@hotmail.com	4	5115. Samaritan	4
## 20	629176	s88724@gmail.com	4	5115. Samaritan	4
## 21	634649	blacksmilez@gmail.com	4	1142. Spider-Man: No Way H~	3
## 22	634649	negativetae@gmail.com	2	1142. Spider-Man: No Way H~	3
## 23	634649	nury95@hotmail.com	3	1142. Spider-Man: No Way H~	3
## 24	634649	delight_32@hotmail.com	3	1142. Spider-Man: No Way H~	3
## 25	634649	s88724@gmail.com	4	1142. Spider-Man: No Way H~	3
## 26	755566	blacksmilez@gmail.com	3	1403. Day Shift	3
## 27	755566	negativetae@gmail.com	4	1403. Day Shift	3
## 28	755566	nury95@hotmail.com	3	1403. Day Shift	3
## 29	755566	delight_32@hotmail.com	4	1403. Day Shift	3
## 30	755566	s88724@gmail.com	3	1403. Day Shift	3

Graph of missing values filled with *median* values

```
ggplot(dfmedian, aes(x=title, y=rate)) +
  geom_boxplot() +
  theme(axis.text.x = element_text(angle = 60, hjust = 1, vjust = 1.0 )) +
  labs(x='Movie Title', y='Rate')
```



## Missing Data (3)

For calibration, the missing values can be filled with *max* values.

```
dfmax <- dftmp %>%
  group_by(movie_id) %>%
  mutate(max = max(rate))

dfmax$rate <- ifelse(dfmax$rate == 0, dfmax$max, dfmax$rate)

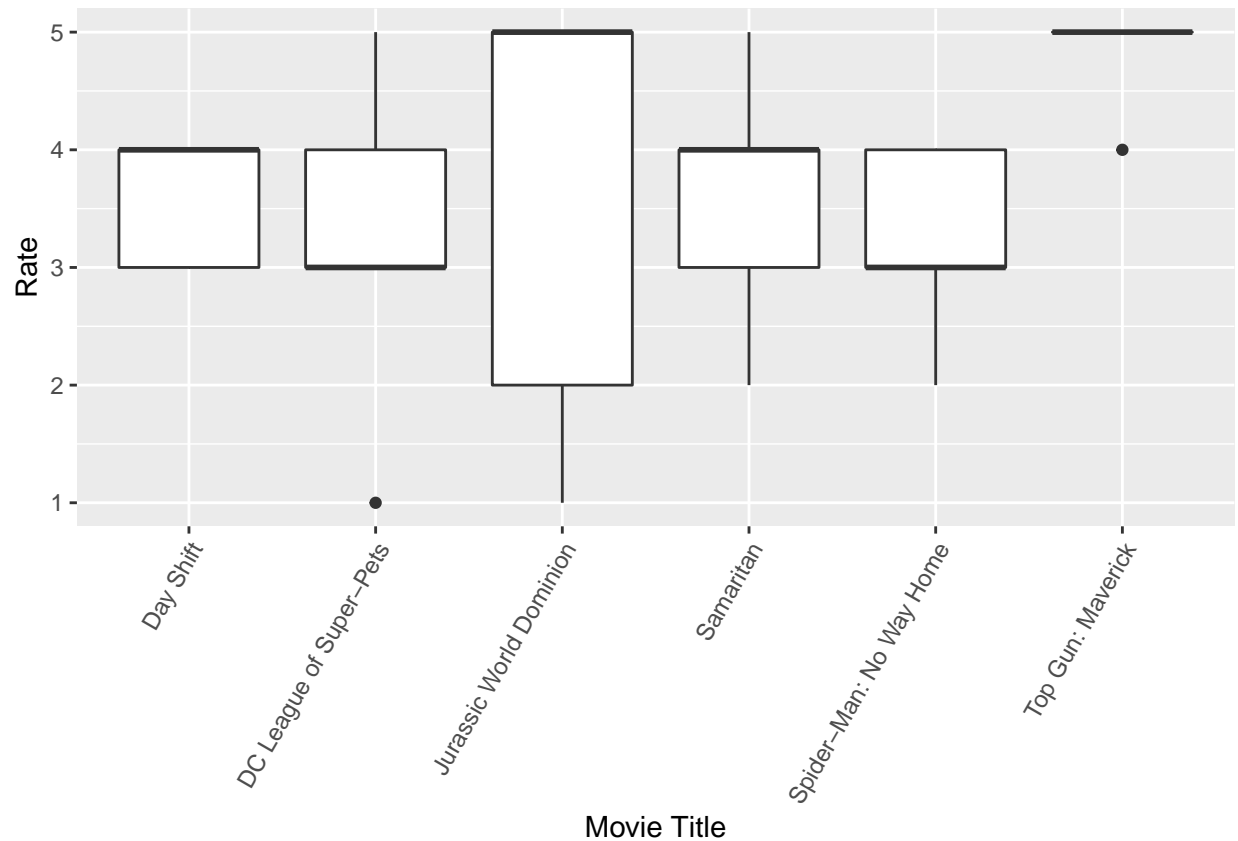
print(dfmax, n=100)
```

```
## # A tibble: 30 x 6
## # Groups:   movie_id [6]
##   movie_id email_address      rate popularity title      max
##   <int> <chr>          <int>    <dbl> <chr>    <int>
## 1 361743 blacksmilez@gmail.com      4    2030. Top Gun: Maverick      5
## 2 361743 negativetae@gmail.com      5    2030. Top Gun: Maverick      5
## 3 361743 nury95@hotmail.com        5    2030. Top Gun: Maverick      5
## 4 361743 delight_32@hotmail.com      5    2030. Top Gun: Maverick      5
## 5 361743 s88724@gmail.com           5    2030. Top Gun: Maverick      5
## 6 507086 blacksmilez@gmail.com      5    2084. Jurassic World Domini~    5
## 7 507086 negativetae@gmail.com      1    2084. Jurassic World Domini~    5
## 8 507086 nury95@hotmail.com        5    2084. Jurassic World Domini~    5
## 9 507086 delight_32@hotmail.com      5    2084. Jurassic World Domini~    5
## 10 507086 s88724@gmail.com           2    2084. Jurassic World Domini~    5
```

## 11	539681	blacksmilez@gmail.com	3	2738.	DC League of Super-Pe~	5
## 12	539681	negativetae@gmail.com	1	2738.	DC League of Super-Pe~	5
## 13	539681	nury95@hotmail.com	5	2738.	DC League of Super-Pe~	5
## 14	539681	delight_32@hotmail.com	4	2738.	DC League of Super-Pe~	5
## 15	539681	s88724@gmail.com	3	2738.	DC League of Super-Pe~	5
## 16	629176	blacksmilez@gmail.com	2	5115.	Samaritan	5
## 17	629176	negativetae@gmail.com	3	5115.	Samaritan	5
## 18	629176	nury95@hotmail.com	5	5115.	Samaritan	5
## 19	629176	delight_32@hotmail.com	4	5115.	Samaritan	5
## 20	629176	s88724@gmail.com	4	5115.	Samaritan	5
## 21	634649	blacksmilez@gmail.com	4	1142.	Spider-Man: No Way Ho~	4
## 22	634649	negativetae@gmail.com	2	1142.	Spider-Man: No Way Ho~	4
## 23	634649	nury95@hotmail.com	3	1142.	Spider-Man: No Way Ho~	4
## 24	634649	delight_32@hotmail.com	3	1142.	Spider-Man: No Way Ho~	4
## 25	634649	s88724@gmail.com	4	1142.	Spider-Man: No Way Ho~	4
## 26	755566	blacksmilez@gmail.com	4	1403.	Day Shift	4
## 27	755566	negativetae@gmail.com	4	1403.	Day Shift	4
## 28	755566	nury95@hotmail.com	3	1403.	Day Shift	4
## 29	755566	delight_32@hotmail.com	4	1403.	Day Shift	4
## 30	755566	s88724@gmail.com	3	1403.	Day Shift	4

Graph of missing values filled with *max* values

```
ggplot(dfmax, aes(x=title, y=rate)) +
  geom_boxplot() +
  theme(axis.text.x = element_text(angle = 60, hjust = 1, vjust = 1.0 )) +
  labs(x='Movie Title', y='Rate')
```



## Closure

The advantage of normalization is that it does not have unnecessary redundant data. It is possible to maintain the integrity of the data by removing the duplicate data. This is a big advantage of relational databases, but in other words, it can also be a big disadvantage. This is because emphasizing excessive normalization causes problems in system performance. Data standardization increases mutual communication by further specifying data. There are various ways to process missing data, but I checked by filling it with mean, median, and max values. Each has a slight difference, so I think we should choose and use it as needed.

Github: <https://github.com/blacksmilez/DATA607/tree/main/Assignment02>