# Analyzing Movie Franchises

Vedang Mehta

Advisor: Prof. Steven Buyske

# Objective

Getting data about movie franchises by scraping IMDb website and using OMDb API, and trying to get interesting insights.

# Outline

- Get list of movie sequels from IMDb.
- Gather data about these movies using OMDb API.
- Figure out which movies belong to the same franchise.
- Tidying up the dataset.
- Find interesting insights.

# Getting the Dataset



- Found a list of sequels on IMDb. This list had 1000+ movies.
- Movies were listed in order of franchise, but there was no indication of separate franchises.
- Had to find a way to identify what franchise a movie belongs to.

# Tidying Up

- Attributes like actors were comma separated strings.
- Box-office collection was in currency ($100,000,000) format. Converted this column to numeric values with regular expression.
- Awards column was in "Won 2 Oscars. Another 11 wins & 20 nominations." format. Wrote a regular expression to get the number of Oscars.
- OMDb had a lot of missing values for box-office collection, which I had to scrape from IMDb.
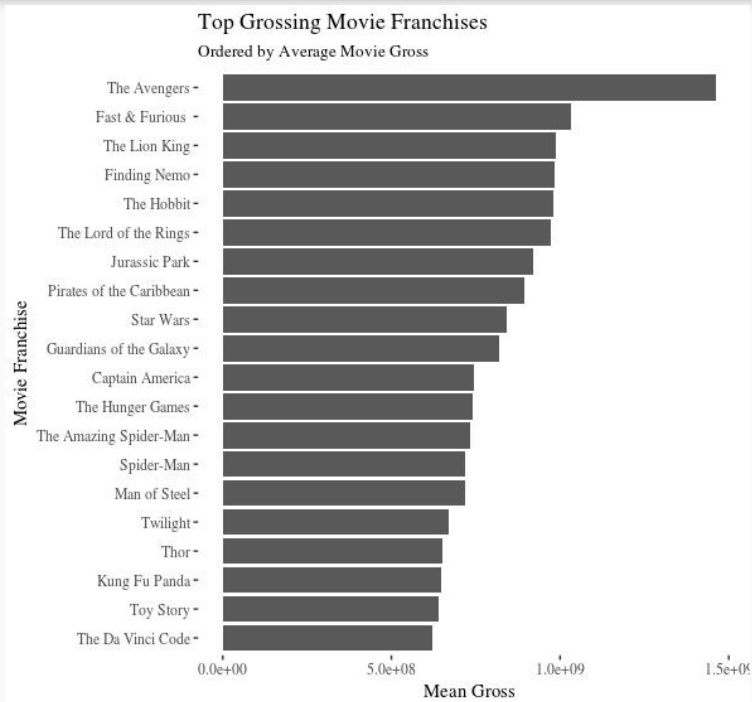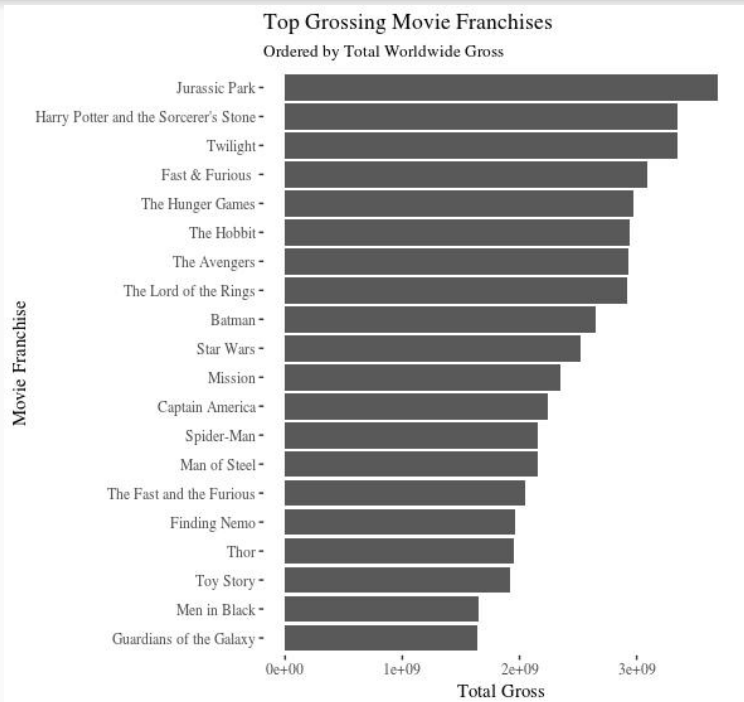
# Getting More Data with OMDb API

- Once I had the list of movie titles from IMDb, I could collect more information using OMDb API.
- OMDb API provided information such as movie title, IMDb rating, box-office collection, actors, director, producer, runtime etc.
- Data was collected in JSON format from OMDb and I stored it to a CSV file so that it becomes easier to process in the future.
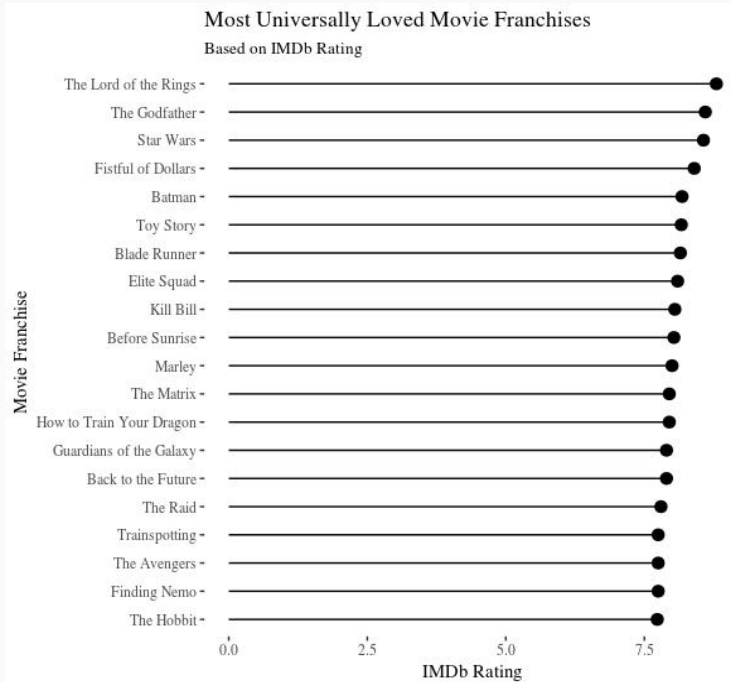
# Identifying Movie Franchise

```r
is_sequel <- function(index){
  if(ls$Year[index] < ls$Year[index - 1]){
    return(FALSE)
  }
  title_sim <- check_for_similarity(get_words(ls$Title[index]), get_words(ls$Title[index - 1]))
  actor_sim <- check_for_similarity(get_words(ls$Actors[index]), get_words(ls$Actors[index - 1]))
  director_sim <- (ls$Director[index] == ls$Director[index - 1])
  production_sim <- (ls$Production[index] == ls$Production[index - 1])
  genre_sim <- check_for_similarity(get_words(ls$Genres[index]), get_words(ls$Genres[index - 1]))
  rated_sim <- (ls$Rated[index] == ls$Rated[index - 1])
  if(actor_sim == TRUE & title_sim == TRUE){
    return(TRUE)
  }
  return(ifelse(title_sim + actor_sim + director_sim + production_sim + genre_sim + rated_sim >= 3, TRUE, FALSE))
}
```

- No need to perform any complex analysis.
- Movies were already ordered by their franchise.
- Comparing basic attributes of two consecutive movies gave satisfactory results.

# Top Grossing Franchises



Top Grossing Movie Franchises
Ordered by Total Worldwide Gross
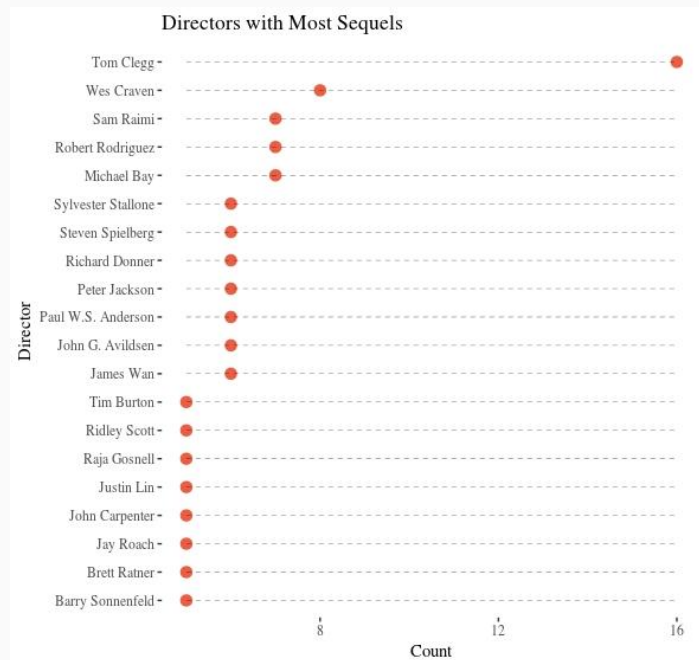
Top Grossing Movie Franchises
Ordered by Average Movie Gross

# Most Universally Loved Movie Franchises



## Number of Oscars

| Movie Franchise | Oscars |
|---|---|
| The Lord of the Rings | |
| The Godfather | |
| Star Wars | |
| The Silence of the Lambs | |
| The French Connection | |
| Terms of Endearment | |
| The Terminator | |
| The Matrix | |
| Raiders of the Lost Ark | |
| The Bourne Identity | |
| Rocky | |
| Jurassic Park | |
| Toy Story | |
| The Little Mermaid | |
| The Lion King | |
| The Exorcist | |
| Speed | |
| Beauty and the Beast | |
| Batman | |
| Alien | |

## Most Universally Loved Movie Franchises
### Based on IMDb Rating

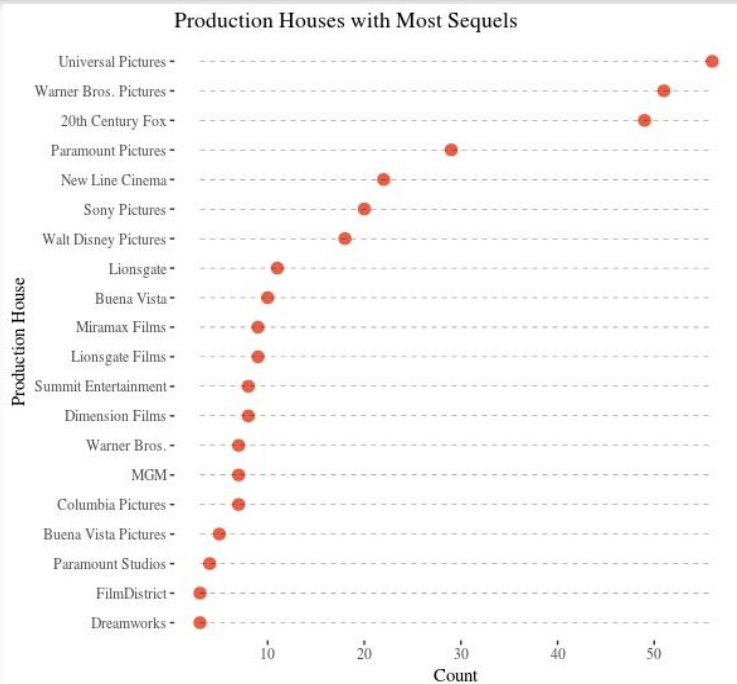| Movie Franchise | IMDb Rating |
|---|---|
| The Lord of the Rings | |
| The Godfather | |
| Star Wars | |
| Fistful of Dollars | |
| Batman | |
| Toy Story | |
| Blade Runner | |
| Elite Squad | |
| Kill Bill | |
| Before Sunrise | |
| Marley | |
| The Matrix | |
| How to Train Your Dragon | |
| Guardians of the Galaxy | |
| Back to the Future | |
| The Raid | |
| Trainspotting | |
| The Avengers | |
| Finding Nemo | |
| The Hobbit | |

# Directors with Most Sequels

Code for Dot Plot

```
df %>% group_by(Director) %>% summarise(count = n()) %>% arrange(-count) %>%
  head(20) %>% ggplot(aes(x=reorder(Director, count), y=count)) +
  geom_point(col="tomato2", size=3) +    # Draw points
  geom_segment(aes(x=Director,
                   xend=Director,
                   y=min(count),
                   yend=max(count)),
               linetype="dashed",
               size=0.1) +    # Draw dashed lines
  labs(title="Directors with Most Sequels",
       x = "Director",
       y = "Count"
       ) +
  coord_flip() + theme_tufte()
```
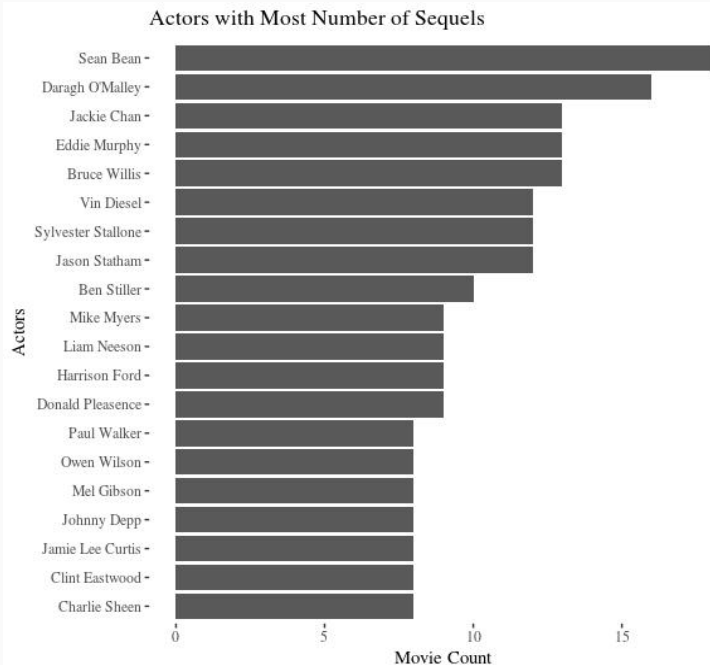


Directors with Most Sequels

# Production Houses with Most Sequels



Production Houses with Most Sequels

Top three production houses produced significantly higher number of sequels than the rest.

# Actors with Most Sequels



Actors with Most Number of Sequels

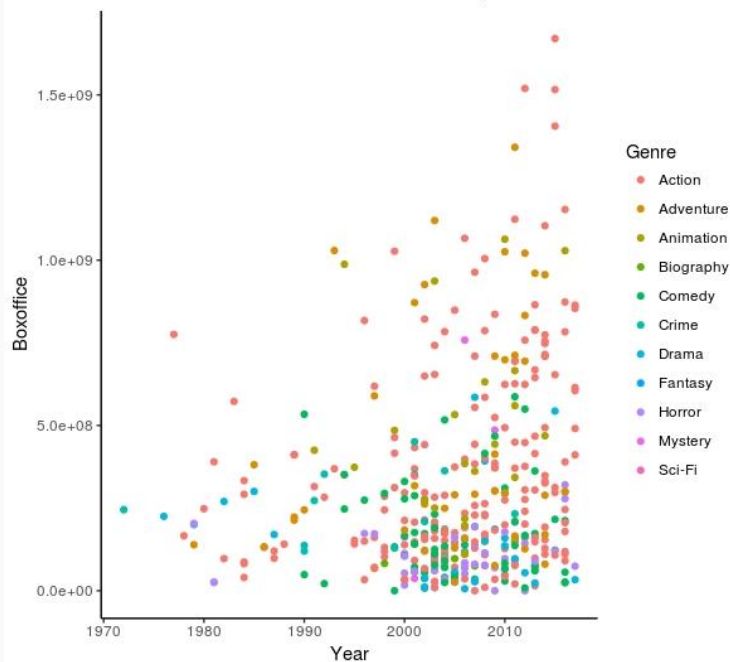As said before, actors were separated by commas. (example, "Seth Rogen, James Franco, Jonah Hill").

Wrote the following code to convert Actors column to a single list of actors -

```
strsplit(df$Actors, ", ") %>% rbind() %>%
apply(MARGIN = 1, unlist)
```
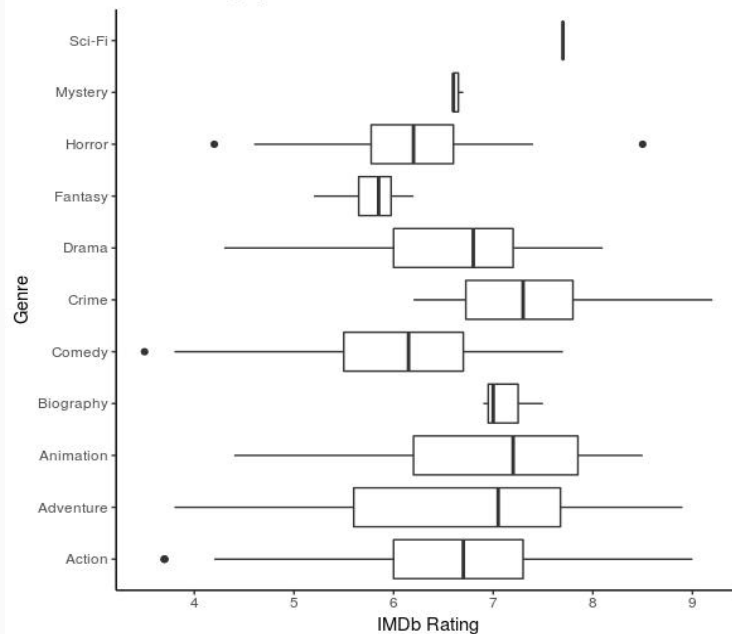
Performed count on this list and plotted a bar plot.

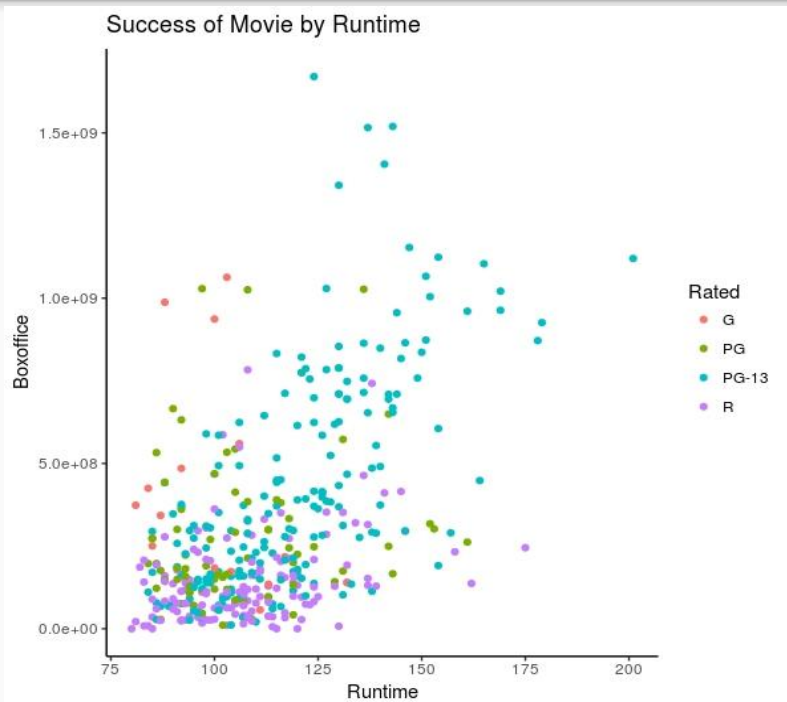# Does Genre Affect the Success of the Movie?



How Boxoffice Collection is Affected by Genre and Time



IMDb Rating by Genre

# Effect of Runtime on Box-office Success
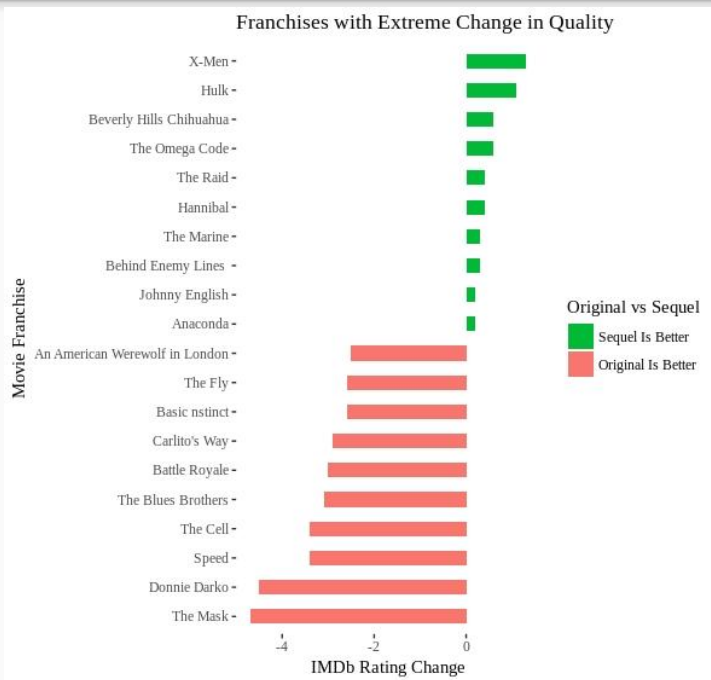


Success of Movie by Runtime

- The correlation between runtime and box-office collection is not very strong.
- Even though the correlation isn't very strong, longer movies tend to earn a bit more.
- Box-office performance correlates more strongly to how the movie is rated.

# Future Scope

- Visualize what franchises are getting better and what franchises are getting worse.*
- Show comparison with inflation adjusted box-office collection.
- How often do sequels perform better than the original movies?*
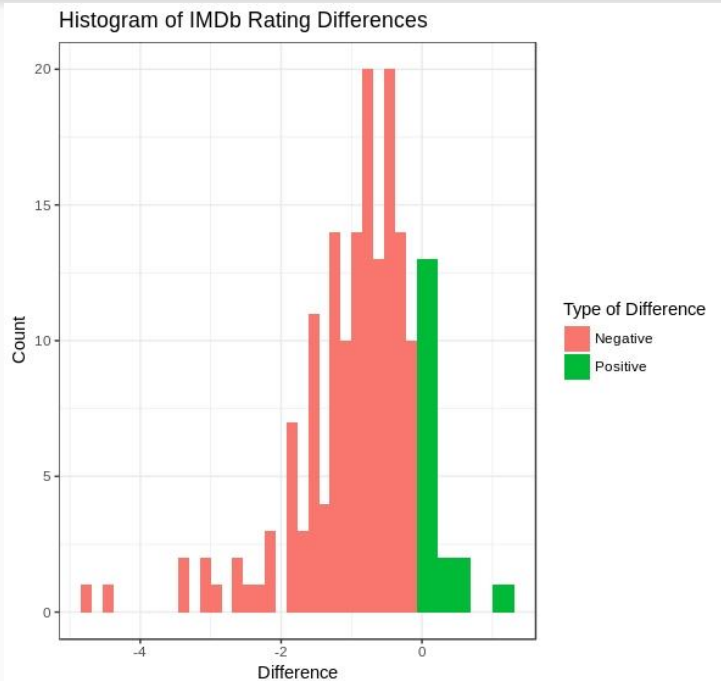
*Added after the in-class presentation. See the next two slides.

# Franchises that Got Significantly Better or Significantly Worse



Franchises with Extreme Change in Quality

- Only considering franchises with two movies.
- There were 188 such franchises in the dataset.
- Displaying 10 movies with highest positive difference and 10 movies with highest negative difference in IMDb ratings.
- Used diverging bars to visualize positive/negative change in rating.

# How Often Do Sequels Perform Better than the Original?



Histogram of IMDb Rating Differences

- Created a histogram of IMDb rating differences.
- Only 34/188 franchises in the dataset with the sequel having higher rating than the original.
- Magnitude of rating difference is also higher on the negative side than the positive.

# References

1. OMDb API - http://www.omdbapi.com/
2. IMDb to collect the list of movie sequels - https://www.imdb.com/list/ls003495084/
3. Gallery of ggplot2 visualizations - http://r-statistics.co/Top50-Ggplot2-Visualizations-MasterList-R-Code.html

# Thank You!