# IMDb Movie Analysis

ECE 180 Team 7

Aiswarya Akumalla/Liyan Chen/Qiuhao Cao/Sean Kamano/Xiaoguang Wang/Xinhe Zhang/Zhisheng Huang
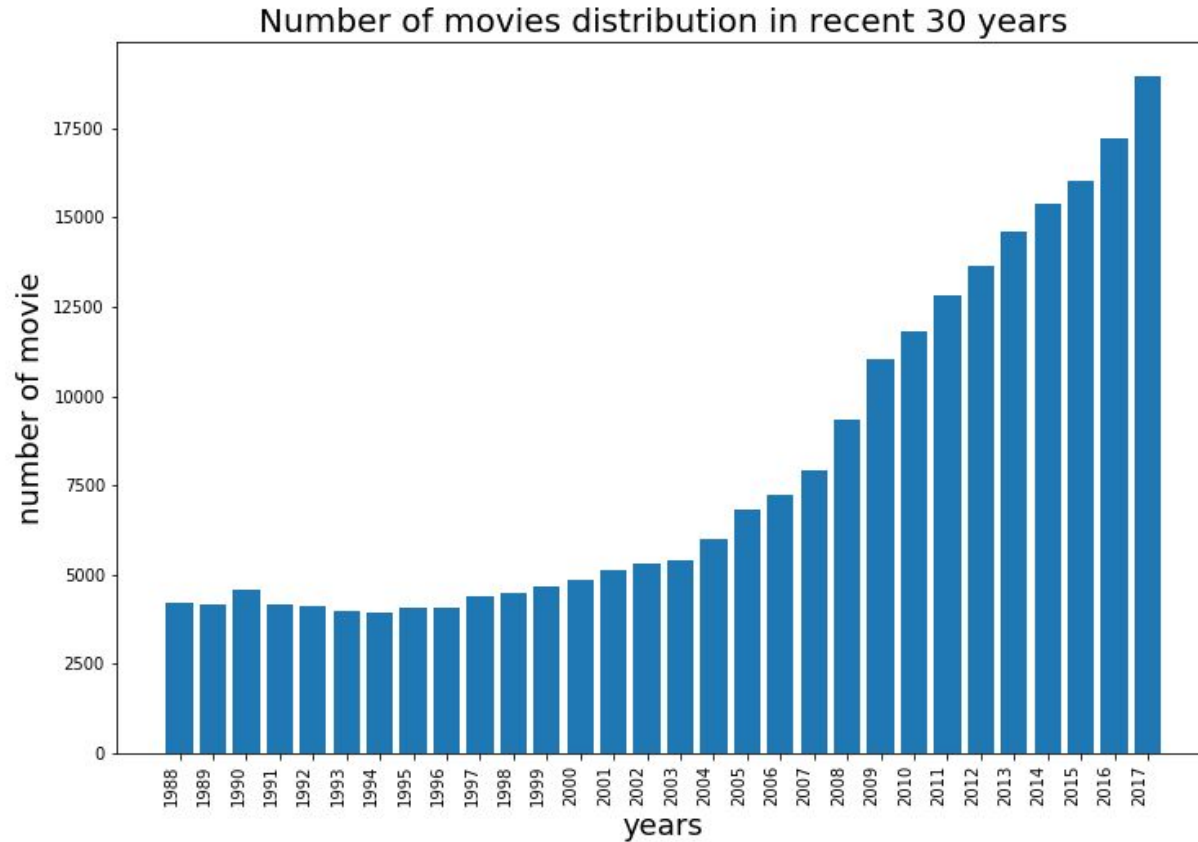
# Outline

- Motivation
- Dataset
- Data Cleaning
- Results
  - Rating
  - Time
  - Location
  - Connection Network
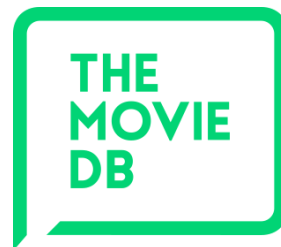- Conclusion

# Motivation

- Analyze the IMDb dataset to provide information for distributors, movie production firms and marketing firms
- Understand aspects of this industry such as:
  - The composition of all movies in the dataset on based on certain properties (genre, region, release date, etc.)
  - Professional networks and popularity of actors through visualization

# Distribution of movies over time



Number of movies distribution in recent 30 years

# Dataset

- IMDb (Internet Movie Database)
  - 478,145 movies
  - basics: genre, release year, and runtime
  - crew: directors and writers
  - principals: top-billed cast/crew
  - ratings: number of votes and weighted average rating
  - name: information about cast/crew (associated movies, birth and death year, and profession)
  - region: ISO-A2 (two-character) country codes for the country each movie was released in
- TMDb (The Movie Database)
  - person: list of popular actor along with popularity scores

# Data Cleaning

- Filtered titles by type "Movie"
- Dropped rows with missing values
- Dropped rows where number of Votes for a movie < 100,000
- Split values of a column and duplicate corresponding rows

| tconst | directors | avgRating | # of Votes |
|---|---|---|---|
| tt0000233 | nm0000690 | 8.4 | 134500 |
| tt0000234 | "nm0000400, nm0000500, nm0000600" | 7.4 | 107800 |

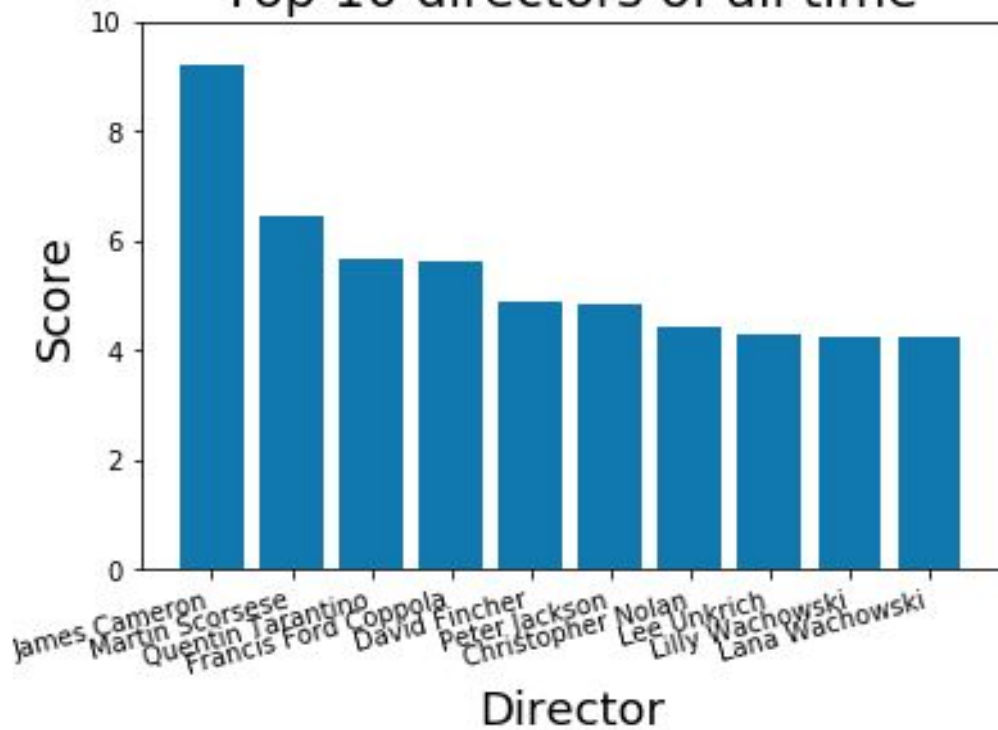| tconst | directors | avgRating | # of Votes |
|---|---|---|---|
| tt0000233 | nm0000690 | 8.4 | 134500 |
| tt0000234 | nm0000400 | 7.4 | 107800 |
| tt0000234 | nm0000500 | 7.4 | 107800 |
| tt0000234 | nm0000600 | 7.4 | 107800 |

# Results -- Top 10 directors

- Average Rating for a title = R
- Number of votes for a title = V
- # of movies made by director threshold > 4

Score for a director = R * V / no. of movies

| tconst | directors | avgRating | numVotes |
|--------|-----------|-----------|----------|
| tt0040746 | nm00000033 | 8.0 | 106921 |
| tt0044079 | nm00000033 | 8.0 | 107709 |
| tt0055600 | nm00000033 | 8.2 | 361351 |
| tt0045906 | nm00000033 | 8.5 | 124765 |

Top 10 directors of all time

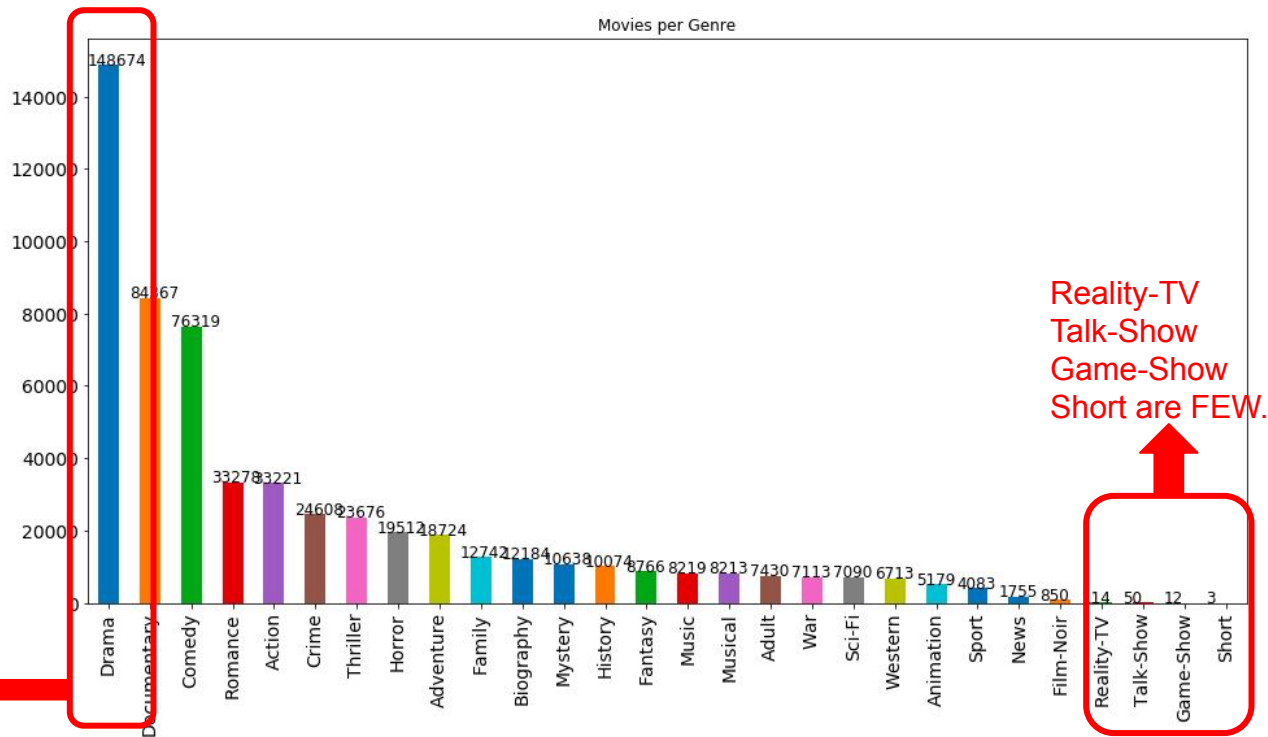| James Cameron |
| Martin Scorsese |
| Quentin Tarantino |
| Francis Ford Coppola |
| David Fincher |
| Peter Jackson |
| Christopher Nolan |
| Lee Unkrich |
| Lilly Wachowski |
| Lana Wachowski |

# Results -- Distribution of genres

- 364461 movies from 1894 to 2018
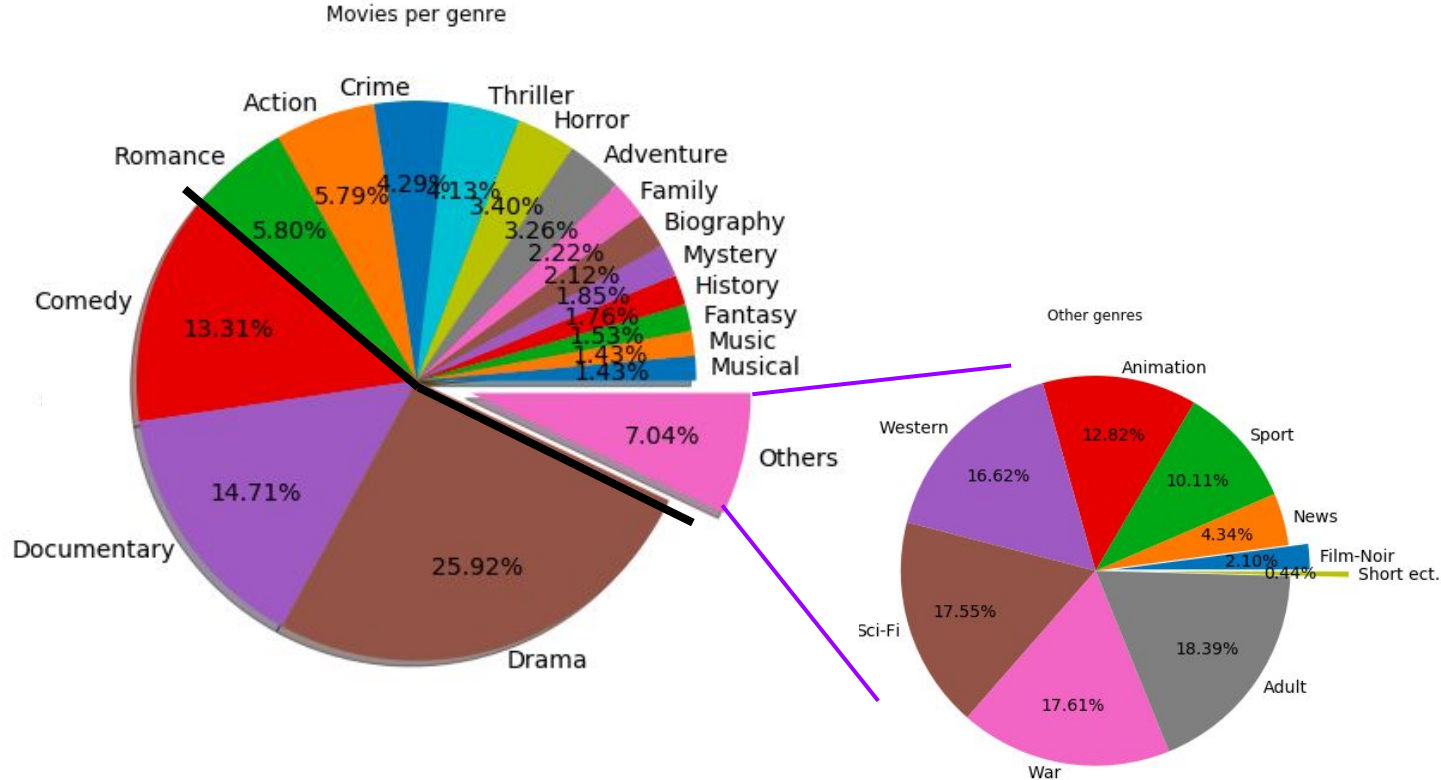- Multiple genre titles:

Ex: Movie is categorized as Drama/Comedy

Drama count increased by 1
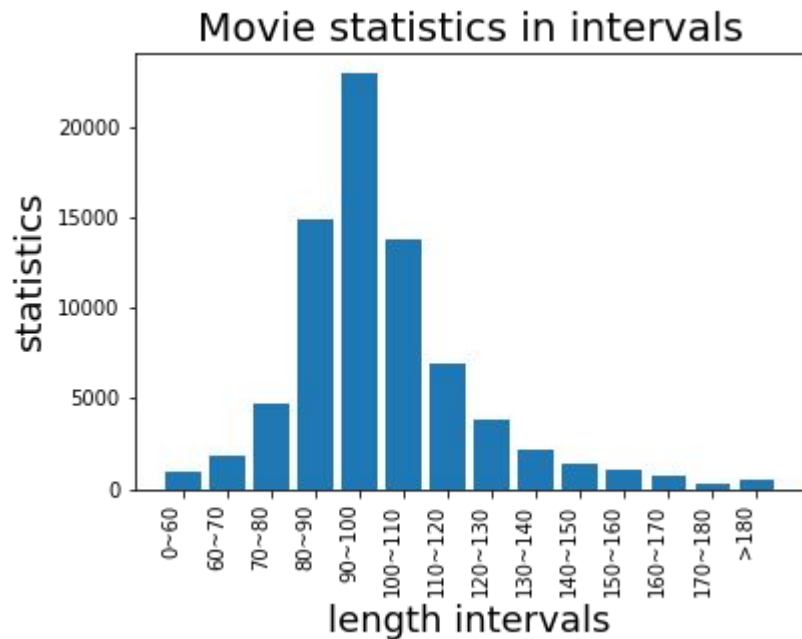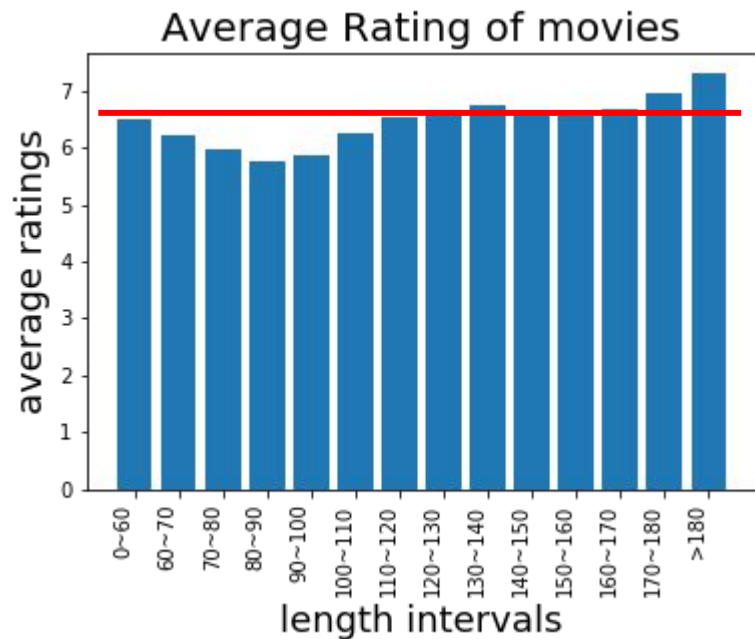Comedy count increased by 1

People love **Dramas**!

Reality-TV
Talk-Show
Game-Show
Short are FEW.



Movies per Genre

| Genre | Count |
|-------|-------|
| Drama | 148674 |
| Documentary | 84567 |
| Comedy | 76319 |
| Romance | 33278 |
| Action | 33221 |
| Crime | 24608 |
| Thriller | 23676 |
| Horror | 19512 |
| Adventure | 18724 |
| Family | 12742 |
| Biography | 12184 |
| Mystery | 10638 |
| History | 10074 |
| Fantasy | 8766 |
| Music | 8219 |
| Musical | 8213 |
| Adult | 7430 |
| War | 7113 |
| Sci-Fi | 7090 |
| Western | 6713 |
| Animation | 5179 |
| Sport | 4083 |
| News | 1755 |
| Film-Noir | 850 |
| Reality-TV | 14 |
| Talk-Show | 50 |
| Game-Show | 12 |
| Short | 3 |

# Results -- Distribution of genres(Cont.)

- Half of movies are **comedy**, **documentary** and **drama**.



Movies per genre

Crime
Action
Thriller
Horror
Romance
Adventure
5.79%  4.29% 4.13%
Family
3.40%
Biography
3.26%
Mystery
2.22%
History
2.12%
Fantasy
1.85%
Comedy
1.76%
Music
5.80%
1.53%
Musical
1.43%
13.31%
1.43%
7.04%
14.71%
Others
Documentary
25.92%
Drama

Other genres

Animation
Western
12.82%
Sport
16.62%
10.11%
News
4.34%
Film-Noir
2.10%
Sci-Fi
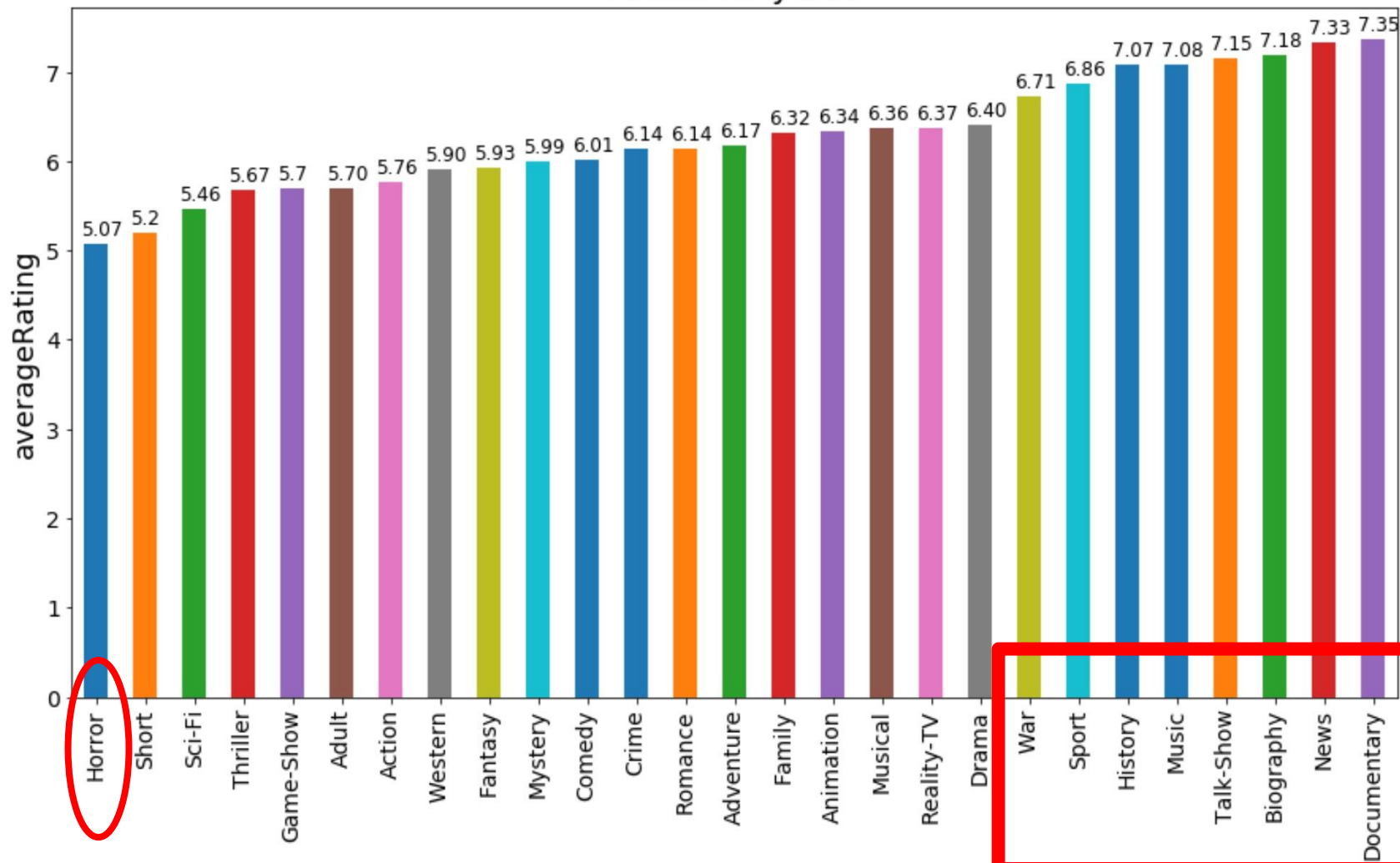0.44%  Short ect.
17.55%
18.39%
War
17.61%  Adult

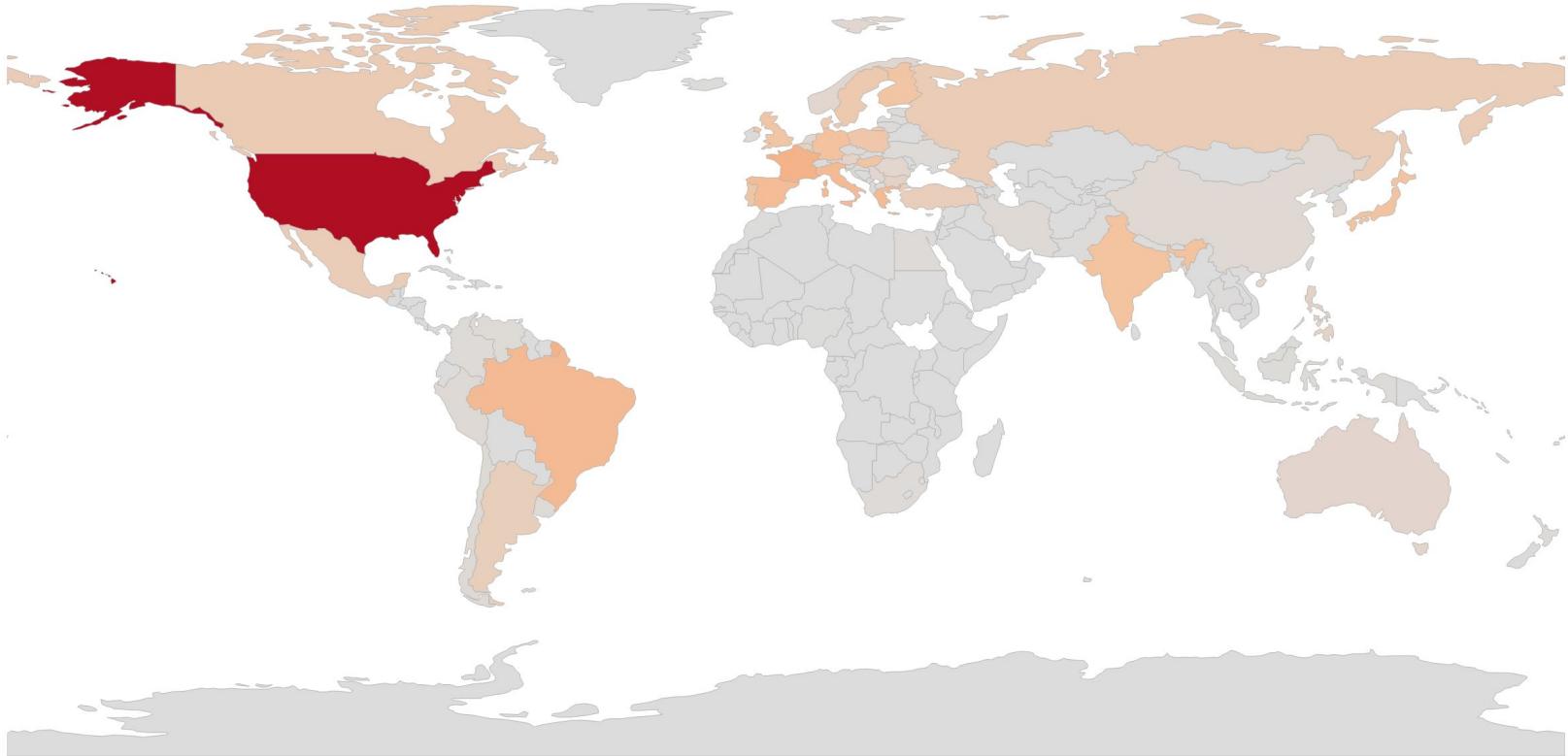# Results -- Movie length and ratings

# Results -- Genres and ratings

recent 20 years

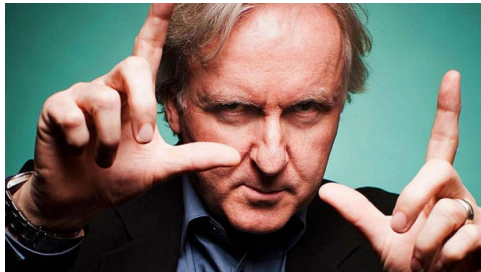# Results -- <u>Distribution of movies among countries</u>

# Results -- The connection network

- Connections between the 100 most popular actors and actresses (which movies they have starred in together)
- Shows how popular they are relative to each other based on a popularity score
  - Popularity score: 0-20, log(number of votes * average rating)

# Concluding Remarks

- Significant growth in the output of movies over the last decade
- Analysis of ratings reveals information about the movie industry such as:
  - Most well-received directors
  - Consumer preference for dramas, comedies, and documentaries
  - Largest audiences in the US and the West
  - Relationships between actors/actresses and their popularity



Drama/Documentary
U.S.

# Future Work

- Through the information gained from our project, we can figure out the most important factors which affect ratings

  Example:  ratings-based recommender system

- What matters:
  - Genre, directors, actors/actress
- What doesn't:
  - Movie length

# Thank You
# Happy movie watching!