

# IMDB.COM

WEB SCRAPING PROJECT

BY STEPHEN SHAFER

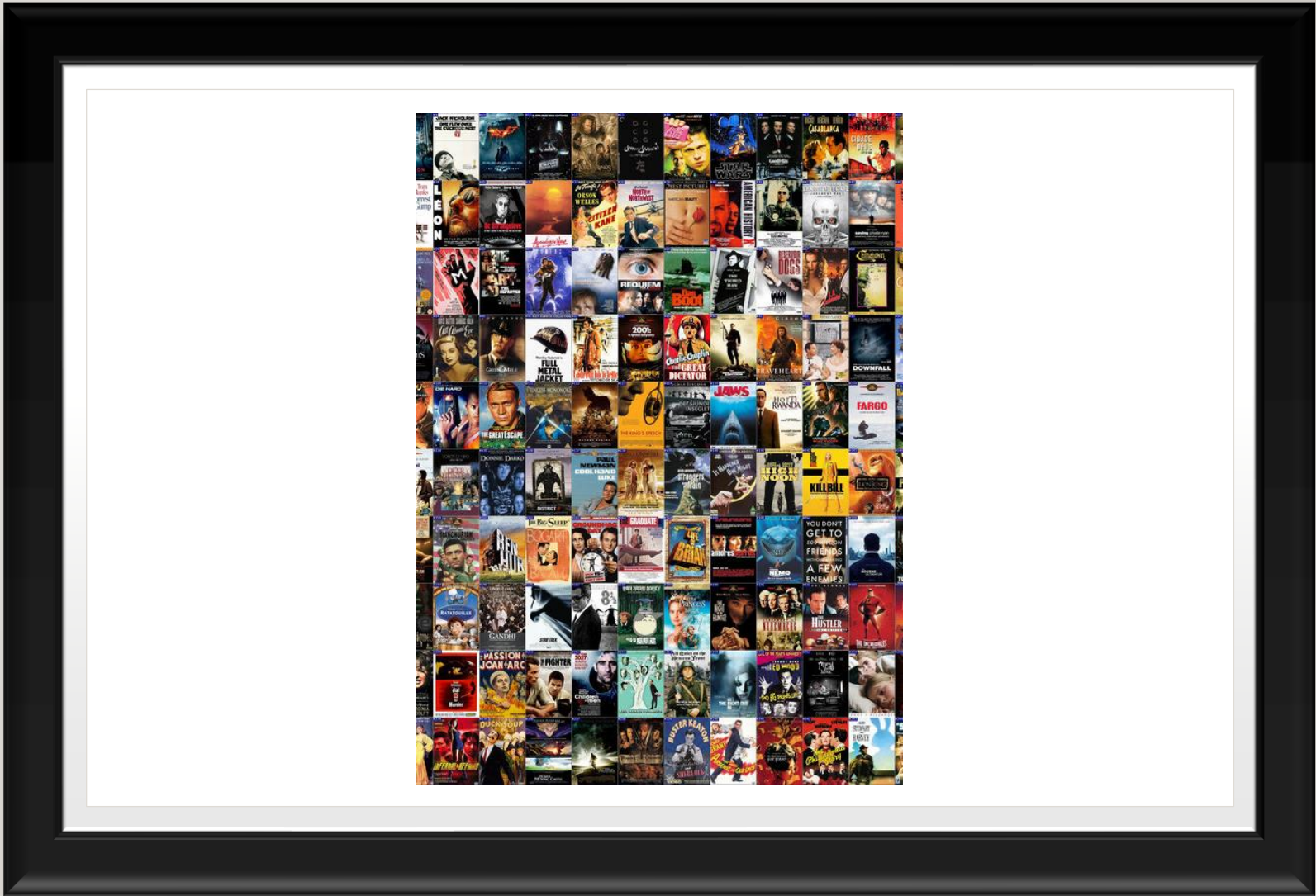
# CONTENTS

OVERVIEW

WEB SCRAPING

DATA CLEANING

ANALYSIS/VISUALIZATION



# QUESTIONS ASKED GOING IN

---

- What is the general makeup of IMDb's users?
- Do users rate movies differently based on their demographic?
- When are the highest rated movies released?
- Does voting between these demographic differ based on who the lead actor is?



# WEB SCRAPING

```
item = ImdbItem()
item['run_time'] = run_time
item['genre'] = genre
item['title'] = title
item['imdb_rating'] = imdb_rating
item['meta_rating'] = meta_rating
item['MPAA_rating'] = MPAA_rating
item['release_date'] = release_date
item['director'] = director
item['actors'] = actors
item['male_teen_rating'] = male_teen_rating
item['male_youngAdult_rating'] = male_youngAdult_rating
item['male_adult_rating'] = male_adult_rating
item['male_elder_rating'] = male_elder_rating
item['male_ratingCount'] = male_ratingCount
item['female_teen_rating'] = female_teen_rating
item['female_youngAdult_rating'] = female_youngAdult_rating
item['female_adult_rating'] = female_adult_rating
item['female_elder_rating'] = female_elder_rating
item['female_ratingCount'] = female_ratingCount
item['non_USusers'] = non_USusers
item['non_UScount'] = non_UScount
item['us_users'] = us_users
item['us_count'] = us_count

yield item
```

# CLEANING

## Dropping NAs and cleaning Data

```
In [6]: # dropping all rows with NAs in meta_rating and male_teen_rating
imdb_test = imdb_test.dropna(subset=["male_under18_rating"])
imdb_test = imdb_test.dropna(subset=["meta_rating"])
#imdb_test[imdb_test.isnull().any(axis=1)]
```

```
In [7]: # changing rating counts into int type
imdb_test.loc[:, "female_ratingCount"] = imdb_test.loc[:, "female_ratingCount"].astype(int)
imdb_test.loc[:, "male_ratingCount"] = imdb_test.loc[:, "male_ratingCount"].astype(int)
imdb_test.loc[:, "meta_rating"] = imdb_test.loc[:, "meta_rating"].astype(int)
```

```
In [8]: # converting release_date into datetime format
import re

def split_it(year):
    return re.findall('(\d+ \w+ \d+)', year)

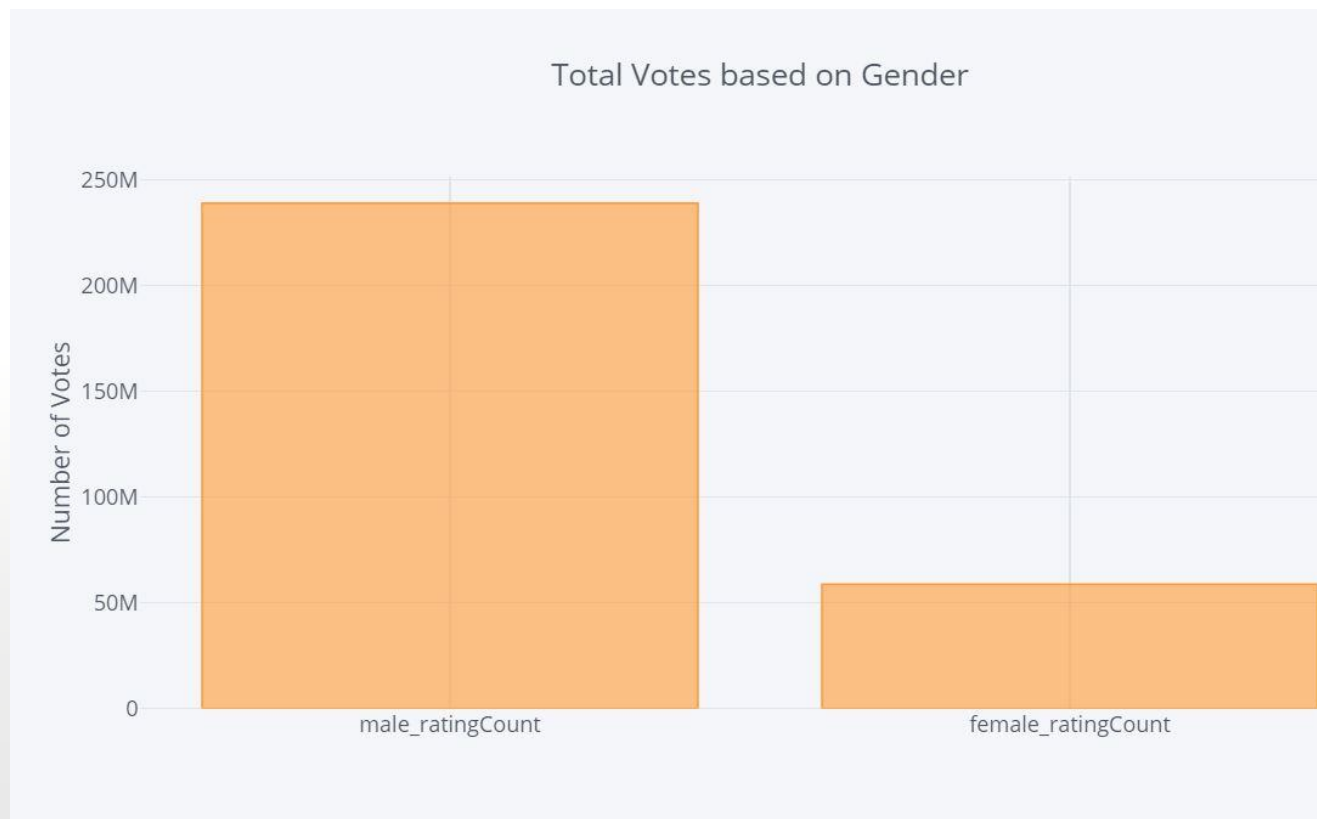
imdb_test['release_date'] = imdb_test['release_date'].apply(split_it)
imdb_test['release_date'] = imdb_test['release_date'].apply(lambda x: ','.join(map(str, x)))
imdb_test['release_date'] = pd.to_datetime(imdb_test['release_date'])
```

```
In [9]: # giving year its own column
imdb_test["release_year"] = imdb_test["release_date"].apply(lambda x : x.year)
imdb_test = imdb_test.dropna(subset=["release_year"])
imdb_test.release_year = imdb_test.release_year.astype(int)
```

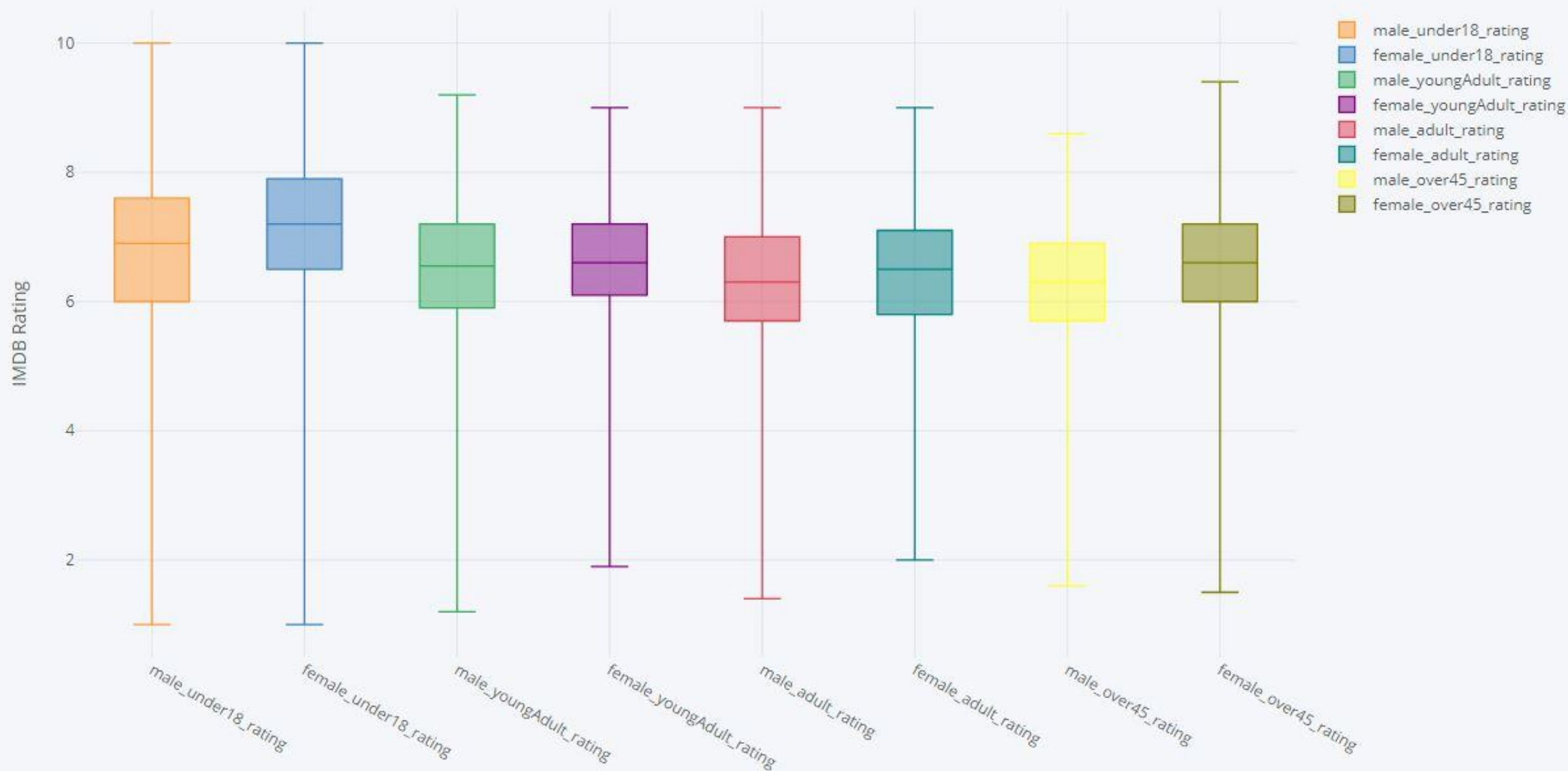
```
In [10]: # giving month its own column
imdb_test["release_month"] = imdb_test["release_date"].apply(lambda x : x.month)
imdb_test = imdb_test.dropna(subset=["release_month"])
imdb_test.release_month = imdb_test.release_month.astype(int)
```

# ANALYSIS & VISUALIZATION

---

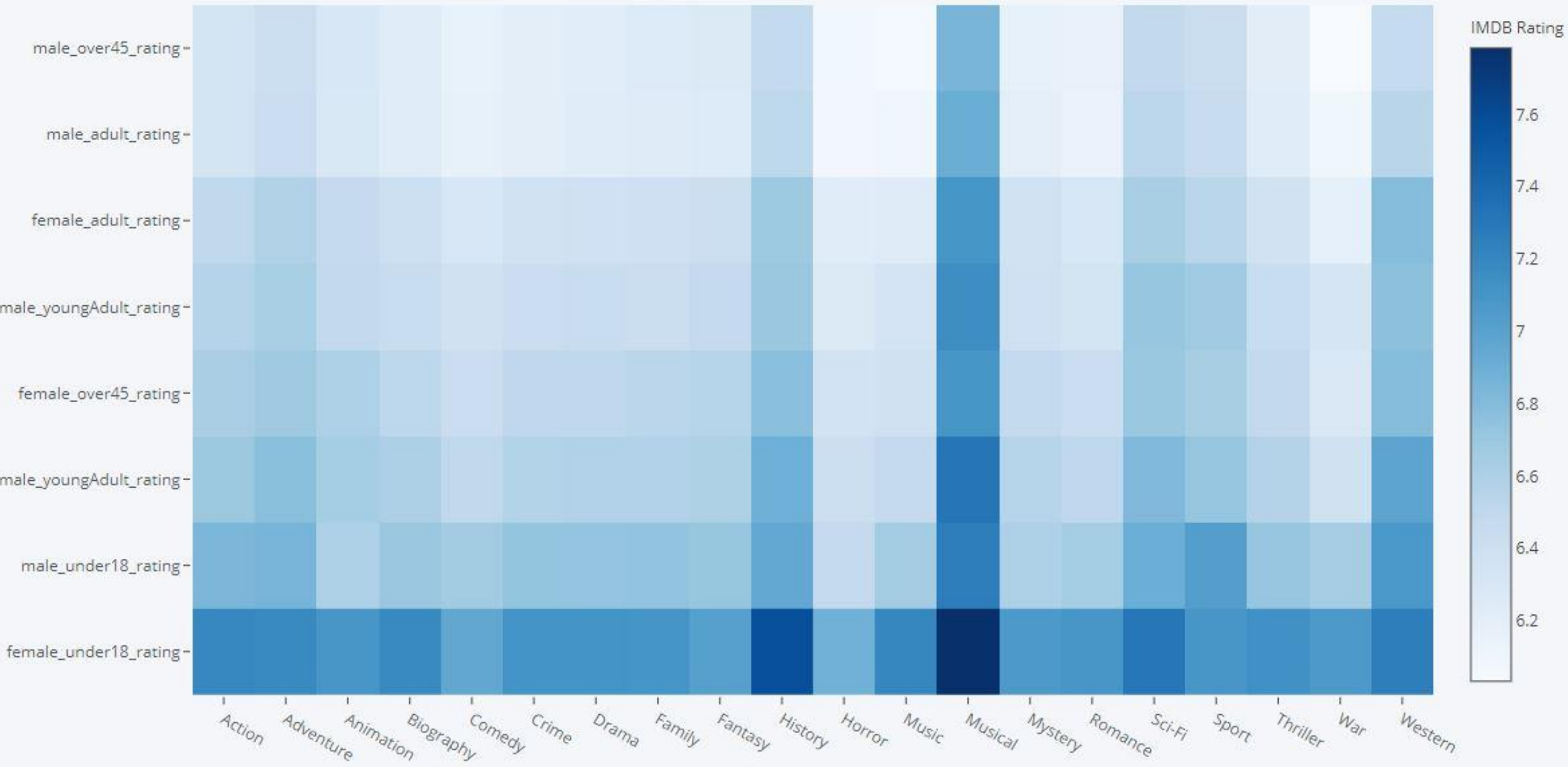


Rating Breakdown by Age/Gender

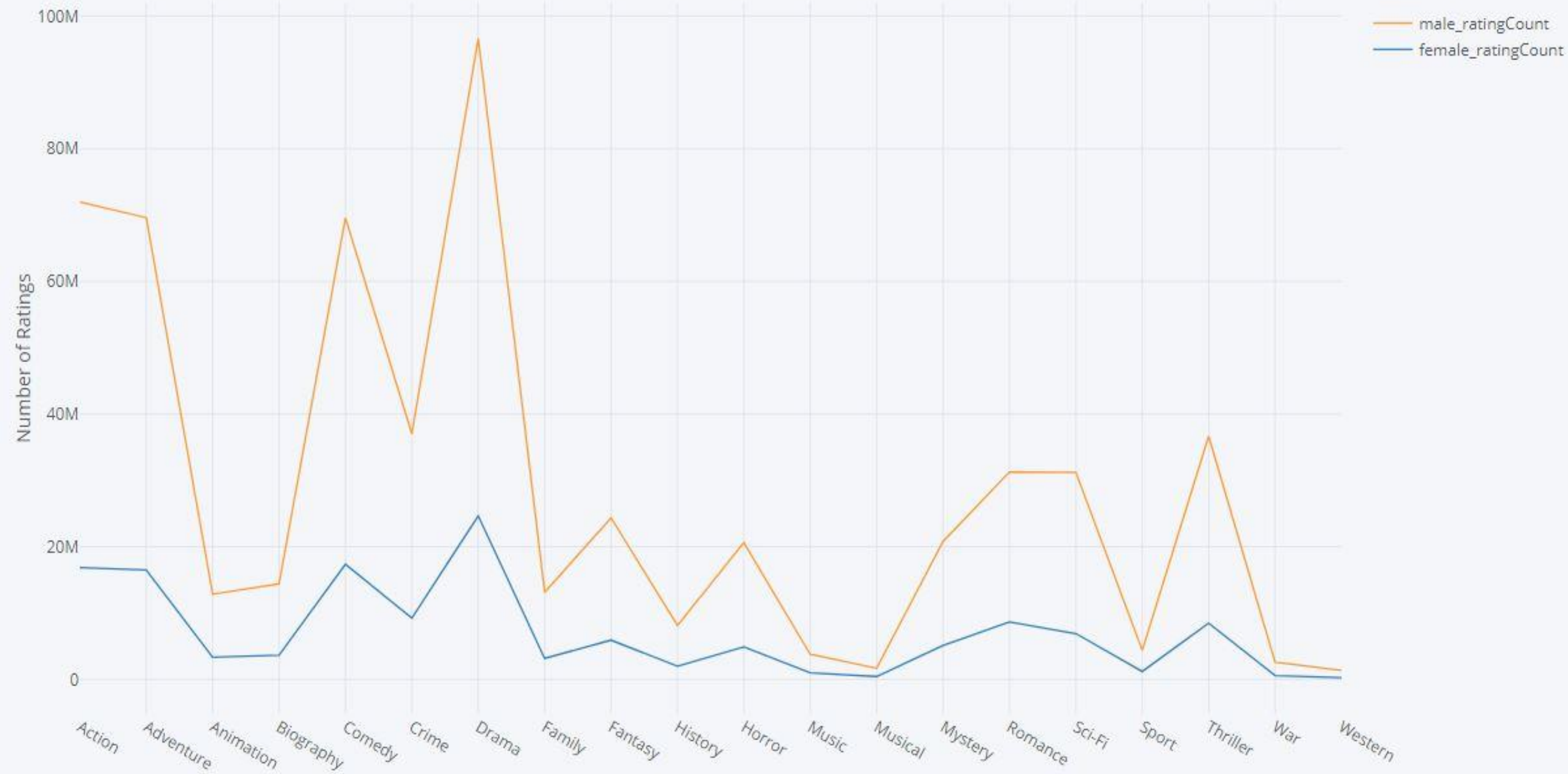


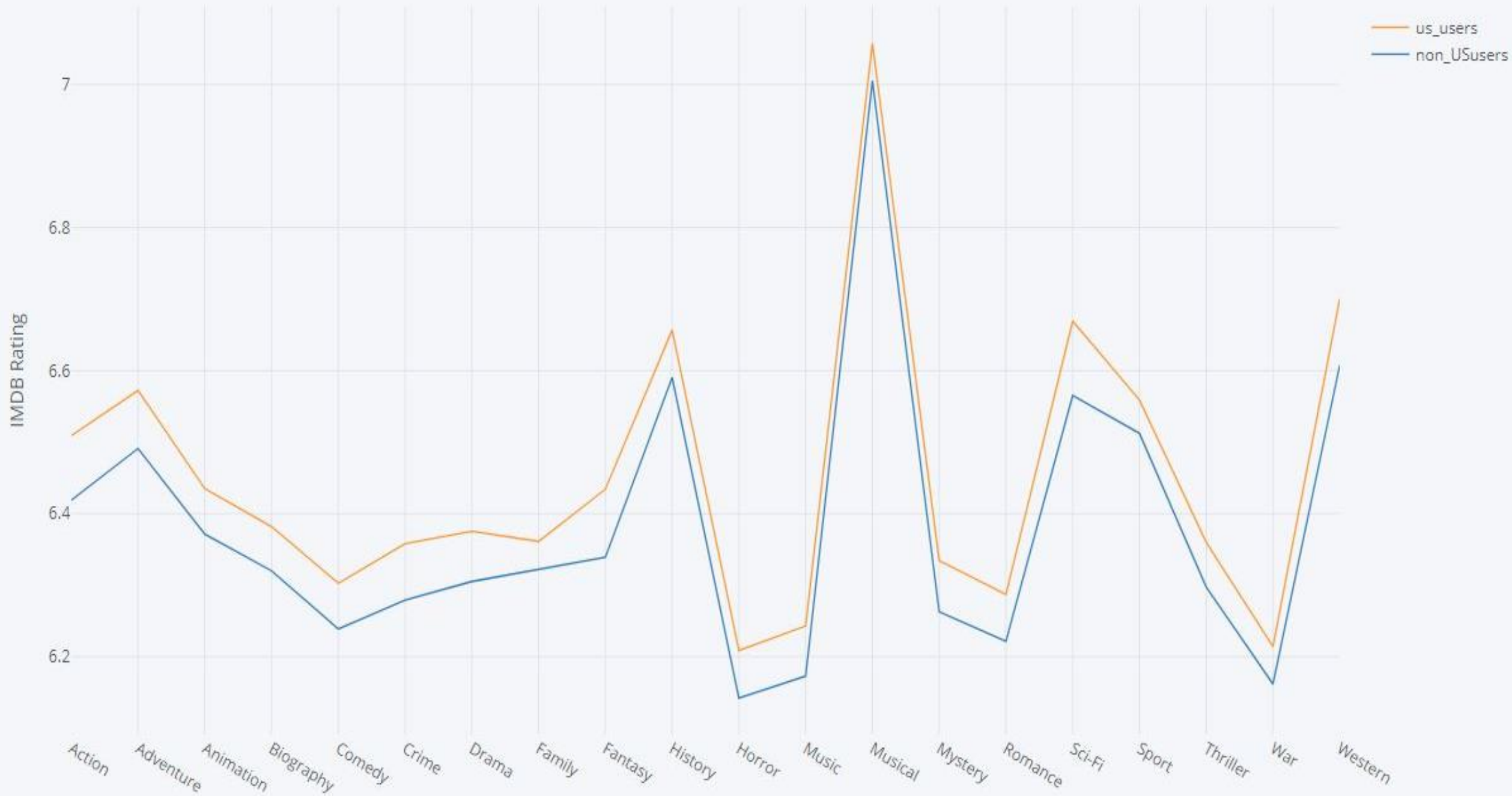


Demographic Rating by Genre

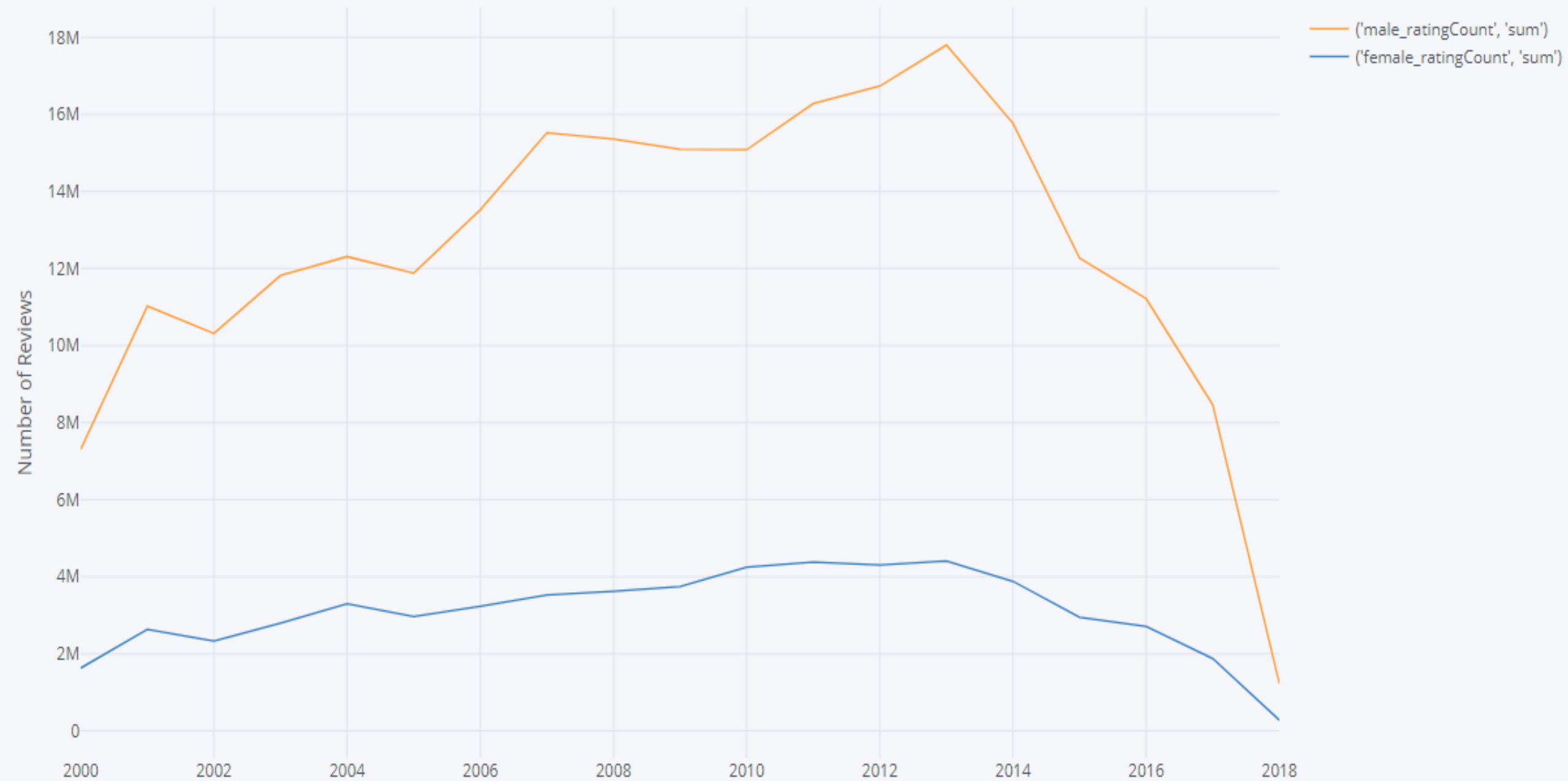


Rating Count based on Gender





Rating Count based on Year





Mean Rating by Year



Mean Rating by Month



Demographic Rating based on Actor



# CONCLUSION

