

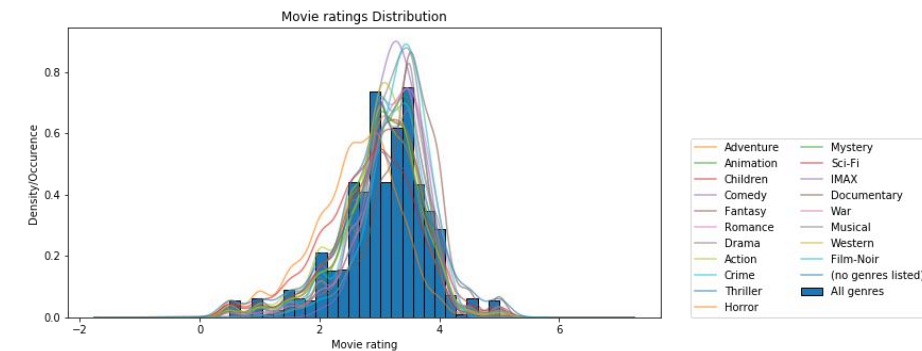
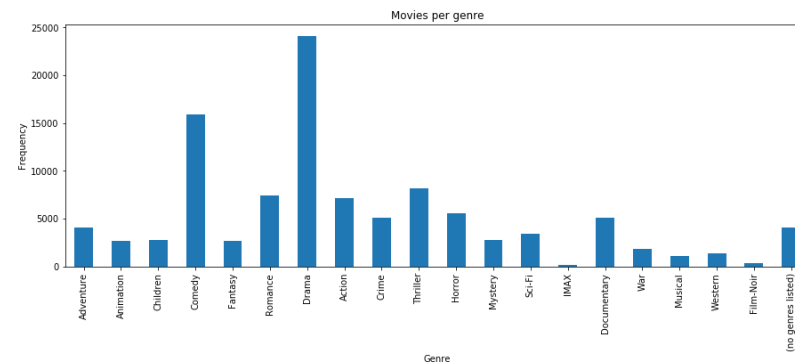
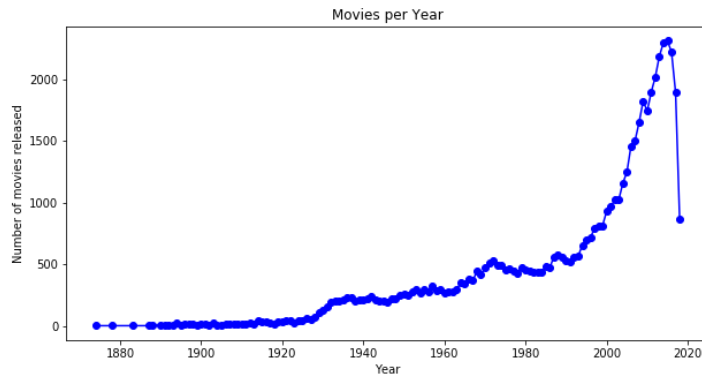
Northwestern

MSiA 423 Mid-project Review: Movie Recommender

By: Sabarish Chockalingam

Highlights

- Interesting fact: Ratings only started coming in from about 1997 (when movie lens started).
- Below are visualization that helped with EDA and will aid when modelling.



Review progress

- Imported Data and perform EDA. (2 pt)
- Stored data set in S3. (1 pt)

Demo/analysis

- Visualization of notebook showing EDA.

```
Movie Recommender Development

In [1]: M 1 import os
2 import pandas as pd
3 from matplotlib import pyplot as plt

In [92]: M 1 ratingspath = os.path.join("../..", "data", "ml-latest", "ratings.csv")
2 ratings = pd.read_csv(ratingspath)
3
4 moviespath = os.path.join("../..", "data", "ml-latest", "movies.csv")
5 movies = pd.read_csv(moviespath)
6

In [8]: M 1 ratings.head()

Out[8]:
  userid  movieid  rating  timestamp
0      1         307    3.5  1256677221
1      1         481    3.5  1256677456
2      1        1091    1.5  1256677471
3      1        1257    4.5  1256677460
4      1        1449    4.5  1256677264

In [53]: M 1 movies.head()

Out[53]:
  movieid  title  genres  year
0         1  Toy Story (1995)  Adventure|Animation|Children|Comedy|Fantasy  1995.0
1         2  Jumanji (1995)  Adventure|Children|Fantasy  1995.0
2         3  Grumpier Old Men (1995)  Comedy|Romance  1995.0
```

```
In [46]: M 1 # extracting date
2 movies['year'] = movies.title.str.extract("\\((\\d{4})\\)", expand=True)

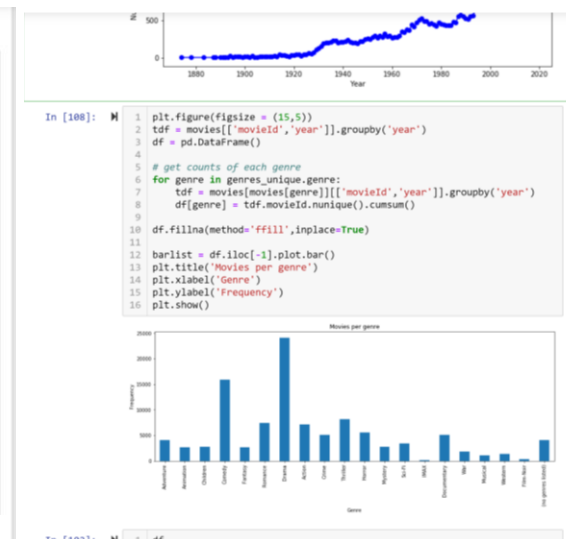
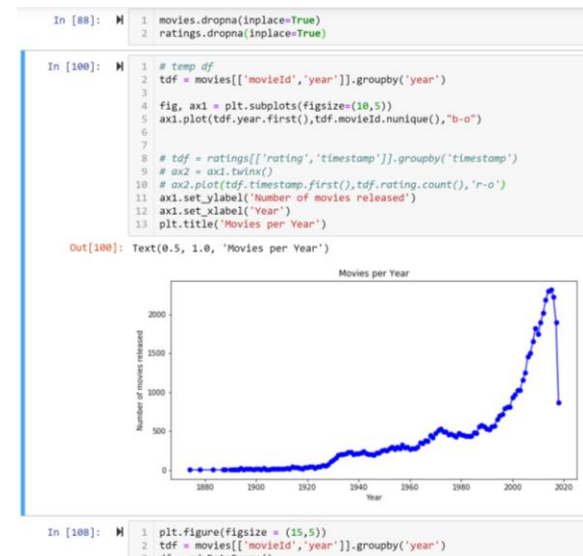
In [48]: M 1 movies.year = pd.to_datetime(movies.year, format = "%Y")
2 # removing any NAs that may result in decimals
3 movies.year = movies.year.dt.year
4 # removing year at end of title
5 #movies.title = movies.title.str[:-7]

In [60]: M 1 # all different types of genre
2 genres_unique = pd.DataFrame(movies.genres.str.split('|').tolist()).stack()
3 genres_unique = pd.DataFrame(genres_unique, columns=['genre'])
4
5 # creating boolean columns for each genre
6 movies = movies.join(movies.genres.str.get_dummies().astype(bool))
7 movies.drop('genres', inplace=True, axis=1)

In [87]: M 1 # converting timestamp to date type
2 ratings.timestamp = pd.to_datetime(ratings.timestamp, infer_datetime_format=True)
3 ratings.timestamp = ratings.timestamp.dt.year

In [88]: M 1 movies.dropna(inplace=True)
2 ratings.dropna(inplace=True)

In [100]: M 1 # temp df
2 tdf = movies[['movieId', 'year']].groupby('year')
3
4 fig, ax1 = plt.subplots(figsize=(10,5))
5 ax1.plot(tdf.year.first(), tdf.movieId.nunique(), "b-o")
6
7 # tdf = ratings[['rating', 'timestamp']].groupby('timestamp')
8 # ax2 = ax1.twinx()
9 ax2.plot(tdf.timestamp.first(), tdf.rating.count(), "r-o")
10 ax1.set_ylabel('Number of movies released')
11 ax1.set_xlabel('Year')
12 plt.title('Movies per Year')
13
```



Lessons learned

- Recommender may not be as easy to evaluate as previously expected since 5 point rating system is not very accurate. For example, it is hard to distinguish top 5 movies for users who rate most movies 5 stars.
- Planning to do the following to improve speed of recommender:
 - Ask user to pick genre.
 - Limit data to be from 1985 to present.

Recommendations

- The following stories are essential for the core app and should be completed next.
 - Create Train and Test Data set. (0 pt)
 - Build and Train Model. (2 pt)
 - Create User Input Boxes (Movie Name and their Rating). (1 pt)
 - Match Input to Movie in Dataset. (1 pt)
 - Autocomplete Title Input. (1 pt)
 - Returns top 5 recommendations from model. (1 pt)
 - Story 1: Deploy app via Flask (AWS). (4 pts)
- After completion of these stories, most of the other stories involve tuning the model and improving the app.