

# Itu Mukherjee - Python for Data Science Week 6 Mini-Project

## Import Modules

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from wordcloud import WordCloud, STOPWORDS, ImageColorGenerator
```

## Import Data

```
In [2]: movies = pd.read_csv('movies.csv', sep = ",")
print(movies.shape)
movies.head()
```

(27278, 3)

Out[2]:

|   | movieId | title                              | genres                                      |
|---|---------|------------------------------------|---|
| 0 | 1       | Toy Story (1995)                   | Adventure Animation Children Comedy Fantasy |
| 1 | 2       | Jumanji (1995)                     | Adventure Children Fantasy                  |
| 2 | 3       | Grumpier Old Men (1995)            | Comedy Romance                              |
| 3 | 4       | Waiting to Exhale (1995)           | Comedy Drama Romance                        |
| 4 | 5       | Father of the Bride Part II (1995) | Comedy                                      |

```
In [3]: ratings = pd.read_csv('ratings.csv', sep = ",")
print(ratings.shape)
ratings.head()
```

(20000263, 4)

Out[3]:

|   | userId | movieId | rating | timestamp  |
|---|--------|---------|--------|------------|
| 0 | 1      | 2       | 3.5    | 1112486027 |
| 1 | 1      | 29      | 3.5    | 1112484676 |
| 2 | 1      | 32      | 3.5    | 1112484819 |
| 3 | 1      | 47      | 3.5    | 1112484727 |
| 4 | 1      | 50      | 3.5    | 1112484580 |

```
In [4]: tags = pd.read_csv('tags.csv', sep = ",")
print(tags.shape)
tags.head()
```

(465564, 4)

Out[4]:

|   | userId | movieId | tag           | timestamp  |
|---|--------|---------|---------------|------------|
| 0 | 18     | 4141    | Mark Waters   | 1240597180 |
| 1 | 65     | 208     | dark hero     | 1368150078 |
| 2 | 65     | 353     | dark hero     | 1368150079 |
| 3 | 65     | 521     | noir thriller | 1368149983 |
| 4 | 65     | 592     | dark hero     | 1368150078 |

## Check for presence of NaN values in the dataframes

```
In [5]: movies.isnull().any()
```

```
Out[5]: movieId    False
title          False
genres         False
dtype: bool
```

```
In [6]: ratings.isnull().any()
```

```
Out[6]: userId      False
movieId           False
rating            False
timestamp         False
dtype: bool
```

```
In [7]: tags.isnull().any()
```

```
Out[7]: userId      False
movieId      False
tag           True
timestamp    False
dtype: bool
```

```
In [8]: tags = tags.dropna()
```

```
In [9]: tags.isnull().any()
```

```
Out[9]: userId      False
movieId      False
tag           False
timestamp    False
dtype: bool
```

## Extract year from title and store the information in a separate column

```
In [10]: movies['year'] = movies['title'].str.extract('.*\((.*)\)'.*, expand=True)
```

```
In [11]: movies.head()
```

```
Out[11]:
```

|   | movieId | title                              | genres                                      | year |
|---|---------|------------------------------------|---|------|
| 0 | 1       | Toy Story (1995)                   | Adventure Animation Children Comedy Fantasy | 1995 |
| 1 | 2       | Jumanji (1995)                     | Adventure Children Fantasy                  | 1995 |
| 2 | 3       | Grumpier Old Men (1995)            | Comedy Romance                              | 1995 |
| 3 | 4       | Waiting to Exhale (1995)           | Comedy Drama Romance                        | 1995 |
| 4 | 5       | Father of the Bride Part II (1995) | Comedy                                      | 1995 |

```
In [12]: average_rating = ratings[['movieId','rating']].groupby('movieId', as_index=False).mean()
average_rating.head()
```

Out[12]:

|   | movieId | rating   |
|---|---------|----------|
| 0 | 1       | 3.921240 |
| 1 | 2       | 3.211977 |
| 2 | 3       | 3.151040 |
| 3 | 4       | 2.861393 |
| 4 | 5       | 3.064592 |

```
In [13]: average_rating['movieId'].unique().shape
```

Out[13]: (26744,)

## Merge dataframes to have yearwise average ratings for all genres

```
In [14]: joined = movies.merge(average_rating, on='movieId', how='inner')
joined.head()
```

Out[14]:

|   | movieId | title                              | genres                                      | year | rating   |
|---|---------|------------------------------------|---|------|----------|
| 0 | 1       | Toy Story (1995)                   | Adventure Animation Children Comedy Fantasy | 1995 | 3.921240 |
| 1 | 2       | Jumanji (1995)                     | Adventure Children Fantasy                  | 1995 | 3.211977 |
| 2 | 3       | Grumpier Old Men (1995)            | Comedy Romance                              | 1995 | 3.151040 |
| 3 | 4       | Waiting to Exhale (1995)           | Comedy Drama Romance                        | 1995 | 2.861393 |
| 4 | 5       | Father of the Bride Part II (1995) | Comedy                                      | 1995 | 3.064592 |

## Merge all dataframes to have one dataframe with all information

```
In [15]: joined_with_tag = joined.merge(tags, on='movieId', how='inner')
joined_with_tag.head()
```

Out[15]:

|          | movieId | title            | genres                                      | year | rating  | userId |                                  |
|----------|---------|------------------|---|------|---------|--------|----------------------------------|
| <b>0</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 1644   | W.                               |
| <b>1</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 1741   | co<br>an                         |
| <b>2</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 1741   | Di<br>an<br>fe:                  |
| <b>3</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 1741   | Pi:<br>an                        |
| <b>4</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 1741   | T<br>Le<br>dc<br>st:<br>thi<br>m |

## Extract the movies appearing in IMDb top 250 list

```
In [16]: imdb_top_250 = joined_with_tag.loc[(joined_with_tag['tag'] == 'imdb top 250')]
print(imdb_top_250.shape)
imdb_top_250.head()
```

(1218, 8)

Out[16]:

|            | movieId | title            | genres                                      | year | rating  | userId |
|------------|---------|------------------|---|------|---------|--------|
| <b>139</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 49307  |
| <b>172</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 55313  |
| <b>292</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 103125 |
| <b>312</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 105214 |
| <b>381</b> | 1       | Toy Story (1995) | Adventure Animation Children Comedy Fantasy | 1995 | 3.92124 | 120937 |

```
In [17]: average_rating_top250 = imdb_top_250[['genres','rating']].groupby('genres', as_index=False).mean()
average_rating_top250.head()
```

Out[17]:

|          | genres   | rating   |
|----------|--|----------|
| <b>0</b> | Action Adventure                                   | 4.169116 |
| <b>1</b> | Action Adventure Animation                         | 3.719764 |
| <b>2</b> | Action Adventure Animation Children Comedy         | 3.908572 |
| <b>3</b> | Action Adventure Animation Children Comedy Fantasy | 3.750635 |
| <b>4</b> | Action Adventure Animation Drama Fantasy           | 4.096299 |

**Split the genre information and then merge them to remove the '|' and insert ',' instead**

```
In [18]: imdb_top250_genres = average_rating_top250['genres'].str.split('|', expand=True)
imdb_top250_genres.head()
```

Out[18]:

|   | 0      | 1         | 2         | 3        | 4       | 5       | 6    |
|---|--------|-----------|-----------|----------|---------|---------|------|
| 0 | Action | Adventure | None      | None     | None    | None    | None |
| 1 | Action | Adventure | Animation | None     | None    | None    | None |
| 2 | Action | Adventure | Animation | Children | Comedy  | None    | None |
| 3 | Action | Adventure | Animation | Children | Comedy  | Fantasy | None |
| 4 | Action | Adventure | Animation | Drama    | Fantasy | None    | None |

```
In [19]: imdb_top250_genres['genres'] = imdb_top250_genres[imdb_top250_genres.columns[0:]]
.apply(
    lambda x: ','.join(x.dropna().astype(str)), axis=1)
imdb_top250_genres.head()
```

Out[19]:

|   | 0      | 1         | 2         | 3        | 4       | 5       | 6    |                           |
|---|--------|-----------|-----------|----------|---------|---------|------|---------------------------|
| 0 | Action | Adventure | None      | None     | None    | None    | None | Action,Adventure          |
| 1 | Action | Adventure | Animation | None     | None    | None    | None | Action,Adventure,Animatio |
| 2 | Action | Adventure | Animation | Children | Comedy  | None    | None | Action,Adventure,Animatio |
| 3 | Action | Adventure | Animation | Children | Comedy  | Fantasy | None | Action,Adventure,Animatio |
| 4 | Action | Adventure | Animation | Drama    | Fantasy | None    | None | Action,Adventure,Animatio |

## Create a wordcloud of movie genres appearing in IMDb top 250 list

```
In [20]: text = ','.join(style for style in imdb_top250_genres.genres)
```

In [21]: text

```
Out[21]: 'Action,Adventure,Action,Adventure,Animation,Action,Adventure,Animation,Children,Comedy,Action,Adventure,Animation,Children,Comedy,Fantasy,Action,Adventure,Animation,Drama,Fantasy,Action,Adventure,Comedy,Fantasy,Action,Adventure,Comedy,Fantasy,Romance,Action,Adventure,Crime,Drama,Thriller,Action,Adventure,Crime,IMAX,Action,Adventure,Drama,Action,Adventure,Drama,Fantasy,Action,Adventure,Drama,Fantasy,Mystery,IMAX,Action,Adventure,Drama,Sci-Fi,Thriller,Action,Adventure,Drama,Thriller,Action,Adventure,Drama,War,Action,Adventure,Drama,Western,Action,Adventure,Fantasy,Horrer,Action,Adventure,Horrer,Sci-Fi,Action,Adventure,Mystery,Romance,Thriller,Action,Adventure,Romance,Action,Adventure,Sci-Fi,Action,Adventure,Sci-Fi,IMAX,Action,Adventure,Sci-Fi,Thriller,War,Action,Adventure,Thriller,IMAX,Action,Adventure,Western,Action,Comedy,Crime,Mystery,Action,Comedy,Horrer,Action,Crime,Action,Crime,Drama,Action,Crime,Drama,IMAX,Action,Crime,Drama,Mystery,Sci-Fi,Thriller,IMAX,Action,Crime,Drama,Thriller,Action,Crime,Film-Noir,Mystery,Thriller,Action,Crime,IMAX,Action,Crime,Thriller,Action,Drama,Romance,Action,Drama,Romance,War,Action,Drama,Sci-Fi,Action,Drama,Thriller,Action,Drama,Thriller,War,Action,Drama,Thriller,Western,Action,Drama,War,Action,Drama,Western,Action,Horrer,Action,Horrer,Sci-Fi,Thriller,Action,Sci-Fi,Action,Sci-Fi,IMAX,Action,Sci-Fi,Thriller,Action,Sci-Fi,Thriller,IMAX,Action,Western,Adventure,Animation,Children,Comedy,Adventure,Animation,Children,Comedy,Fantasy,Adventure,Animation,Children,Comedy,Fantasy,Romance,Adventure,Animation,Children,Romance,Sci-Fi,Adventure,Animation,Fantasy,Adventure,Animation,Fantasy,Romance,Adventure,Children,Fantasy,Musical,Adventure,Comedy,Drama,Adventure,Comedy,Fantasy,Adventure,Comedy,Romance,Adventure,Comedy,Romance,War,Adventure,Comedy,Sci-Fi,Adventure,Drama,Adventure,Drama,IMAX,Adventure,Drama,Sci-Fi,Adventure,Drama,Thriller,Adventure,Drama,War,Adventure,Fantasy,Adventure,Fantasy,IMAX,Adventure,Western,Animation,Children,Drama,Animation,Drama,War,Children,Comedy,Comedy,Comedy,Crime,Comedy,Crime,Drama,Comedy,Crime,Drama,Thriller,Comedy,Crime,Thriller,Comedy,Drama,Comedy,Drama,Romance,Comedy,Drama,Romance,War,Comedy,Drama,War,Comedy,Fantasy,Comedy,Fantasy,Romance,Comedy,Musical,Romance,Comedy,Musical,War,Comedy,Mystery,Thriller,Comedy,Romance,Comedy,War,Crime,Drama,Crime,Drama,Film-Noir,Crime,Drama,Film-Noir,Thriller,Crime,Drama,Mystery,Crime,Drama,Mystery,Thriller,Crime,Drama,Romance,Crime,Drama,Sci-Fi,Thriller,Crime,Drama,Thriller,Crime,Drama,Western,Crime,Fantasy,Horrer,Crime,Film-Noir,Crime,Film-Noir,Mystery,Crime,Film-Noir,Mystery,Thriller,Crime,Film-Noir,Thriller,Crime,Horrer,Crime,Horrer,Thriller,Crime,Mystery,Thriller,Crime,Thriller,War,Drama,Drama,Fantasy,Drama,Fantasy,Horrer,Romance,Drama,Fantasy,Romance,Drama,Fantasy,Thriller,Drama,Film-Noir,Romance,Drama,Film-Noir,Thriller,Drama,Horrer,Mystery,Drama,Horrer,Sci-Fi,Drama,Horrer,Thriller,Drama,Musical,Romance,Drama,Mystery,Drama,Mystery,Romance,Thriller,Drama,Mystery,Sci-Fi,Drama,Mystery,Sci-Fi,Thriller,Drama,Mystery,Thriller,Drama,Mystery,War,Drama,Romance,Drama,Romance,Sci-Fi,Drama,Romance,Thriller,Drama,Romance,War,Drama,Sci-Fi,Drama,Sci-Fi,Thriller,Drama,Thriller,Drama,War,Drama,Western,Fantasy,Sci-Fi,Film-Noir,Film-Noir,Mystery,Film-Noir,Mystery,Thriller,Film-Noir,Romance,Thriller,Horrer,Horrer,Mystery,Horrer,Mystery,Thriller,Horrer,Sci-Fi,Mystery,Sci-Fi,Thriller,Mystery,Thriller,Sci-Fi'
```

```
In [22]: wordcloud = WordCloud(width=800, height=500, colormap="jet",
                                background_color="white").generate(text)
```



```
In [23]: plt.figure()  
plt.imshow(wordcloud, interpolation='bilinear')  
plt.axis("off")  
wordcloud.to_file("imdbtop250_wordcloud.png")
```

Out[23]: <wordcloud.wordcloud.WordCloud at 0x27c0000f828>



**Plot temporal trends in ratings of movie genres since year 2000**

**Create a new dataframe with information only from year 2000**

```
In [24]: temp = joined.loc[(joined['year'] == '2000') | (joined['year'] == '2001') |
                        (joined['year'] == '2002') | (joined['year'] == '2003') |
                        (joined['year'] == '2004') | (joined['year'] == '2005') |
                        (joined['year'] == '2006') | (joined['year'] == '2007') |
                        (joined['year'] == '2009') | (joined['year'] == '2009') |
                        (joined['year'] == '2010') | (joined['year'] == '2011') |
                        (joined['year'] == '2012') | (joined['year'] == '2013') |
                        (joined['year'] == '2014') | (joined['year'] == '2015') |
                        (joined['year'] == '2016') | (joined['year'] == '2017') |
                        (joined['year'] == '2018')]
temp.head()
```

Out[24]:

|             | movieid | title                 | genres                    | year | rating   |
|-------------|---------|-----------------------|---------------------------|------|----------|
| <b>2683</b> | 2769    | Yards, The (2000)     | Crime Drama               | 2000 | 3.129956 |
| <b>3090</b> | 3177    | Next Friday (2000)    | Comedy                    | 2000 | 2.810680 |
| <b>3103</b> | 3190    | Supernova (2000)      | Adventure Sci-Fi Thriller | 2000 | 2.280392 |
| <b>3138</b> | 3225    | Down to You (2000)    | Comedy Romance            | 2000 | 2.644370 |
| <b>3141</b> | 3228    | Wirey Spindell (2000) | Comedy                    | 2000 | 2.104167 |

## Create separate dataframes for each of the six genres to be visualized

```
In [25]: value_list = ['Drama']
drama = temp[temp.genres.isin(value_list)]
drama.head()
```

Out[25]:

|             | movieid | title                      | genres | year | rating   |
|-------------|---------|----------------------------|--------|------|----------|
| <b>3320</b> | 3408    | Erin Brockovich (2000)     | Drama  | 2000 | 3.590389 |
| <b>3424</b> | 3514    | Joe Gould's Secret (2000)  | Drama  | 2000 | 3.290000 |
| <b>3443</b> | 3534    | 28 Days (2000)             | Drama  | 2000 | 3.092077 |
| <b>3488</b> | 3579    | I Dreamed of Africa (2000) | Drama  | 2000 | 2.488333 |
| <b>3489</b> | 3580    | Up at the Villa (2000)     | Drama  | 2000 | 2.997368 |

```
In [26]: value_list = ['Thriller']
thriller = temp[temp.genres.isin(value_list)]
thriller.head()
```

Out[26]:

|             | movieId | title                  | genres   | year | rating   |
|-------------|---------|------------------------|----------|------|----------|
| <b>3204</b> | 3291    | Trois (2000)           | Thriller | 2000 | 1.729167 |
| <b>3394</b> | 3484    | Skulls, The (2000)     | Thriller | 2000 | 2.754047 |
| <b>3706</b> | 3797    | In Crowd, The (2000)   | Thriller | 2000 | 2.272201 |
| <b>3764</b> | 3857    | Bless the Child (2000) | Thriller | 2000 | 2.614094 |
| <b>3926</b> | 4020    | Gift, The (2000)       | Thriller | 2000 | 3.390297 |

```
In [27]: value_list = ['Crime']
crime = temp[temp.genres.isin(value_list)]
crime.head()
```

Out[27]:

|              | movieId | title                                      | genres | year | rating   |
|--------------|---------|--|--------|------|----------|
| <b>6394</b>  | 6504    | Gasoline (Benzina) (2001)                  | Crime  | 2001 | 2.400000 |
| <b>12636</b> | 59418   | American Crime, An (2007)                  | Crime  | 2007 | 3.686528 |
| <b>13479</b> | 66691   | Thick as Thieves (a.k.a. Code, The) (2009) | Crime  | 2009 | 2.992593 |
| <b>15600</b> | 79536   | Hellsinki (Rööperi) (2009)                 | Crime  | 2009 | 3.144068 |
| <b>18730</b> | 93490   | Law of the Lawless (Brigada) (2002)        | Crime  | 2002 | 3.875000 |

```
In [28]: value_list = ['Comedy']
comedy = temp[temp.genres.isin(value_list)]
comedy.head()
```

Out[28]:

|             | movieId | title                   | genres | year | rating   |
|-------------|---------|-------------------------|--------|------|----------|
| <b>3090</b> | 3177    | Next Friday (2000)      | Comedy | 2000 | 2.810680 |
| <b>3141</b> | 3228    | Wirey Spindell (2000)   | Comedy | 2000 | 2.104167 |
| <b>3152</b> | 3239    | Isn't She Great? (2000) | Comedy | 2000 | 2.270732 |
| <b>3189</b> | 3276    | Gun Shy (2000)          | Comedy | 2000 | 2.925926 |
| <b>3199</b> | 3286    | Snow Day (2000)         | Comedy | 2000 | 2.222868 |

```
In [29]: value_list = ['Romance']
romance = temp[temp.genres.isin(value_list)]
romance.head()
```

Out[29]:

|             | movieId | title  | genres  | year | rating   |
|-------------|---------|--|---------|------|----------|
| <b>3930</b> | 4024    | House of Mirth, The (2000)                   | Romance | 2000 | 3.418321 |
| <b>4796</b> | 4892    | Maze (2000)                                  | Romance | 2000 | 3.190000 |
| <b>5565</b> | 5664    | Brown Sugar (2002)                           | Romance | 2002 | 3.053672 |
| <b>6093</b> | 6192    | Open Hearts (Elsker dig for evigt) (2002)    | Romance | 2002 | 3.680556 |
| <b>6446</b> | 6556    | Sea Is Watching, The (Umi wa miteita) (2002) | Romance | 2002 | 3.000000 |

```
In [30]: value_list = ['Sci-Fi']
scifi = temp[temp.genres.isin(value_list)]
scifi.head()
```

Out[30]:

|              | movieId | title   | genres | year | rating   |
|--------------|---------|---|--------|------|----------|
| <b>3267</b>  | 3354    | Mission to Mars (2000)                            | Sci-Fi | 2000 | 2.625766 |
| <b>10920</b> | 44671   | Wild Blue Yonder, The (2005)                      | Sci-Fi | 2005 | 2.803571 |
| <b>13071</b> | 62834   | Babylon 5: The Legend of the Rangers: To Live ... | Sci-Fi | 2002 | 2.952991 |
| <b>13072</b> | 62836   | Babylon 5: The Lost Tales - Voices in the Dark... | Sci-Fi | 2007 | 3.218447 |
| <b>15217</b> | 77795   | Cargo (2009)                                      | Sci-Fi | 2009 | 3.464286 |

**Plot each of the dataframes on the same matplotlib axis object, thereby creating an overlaid time series plot**

```

In [31]: with plt.style.context('seaborn-whitegrid'):
fig,ax=plt.subplots(figsize=(10,4))
ax=sns.pointplot(x='year',y='rating',data=drama,ci=None,color='blue')
ax=sns.pointplot(x='year',y='rating',data=thriller,ci=None,color='red')
ax=sns.pointplot(x='year',y='rating',data=crime,ci=None,color='magenta')
ax=sns.pointplot(x='year',y='rating',data=comedy,ci=None,color='forestgreen')
n')
ax=sns.pointplot(x='year',y='rating',data=romance,ci=None,color='orange')
ax=sns.pointplot(x='year',y='rating',data=scifi,ci=None,color='black')

plt.xticks(rotation=90)
plt.box(False)
ax.grid('both','both',linestyle='--')
ax.legend(['Drama','Thriller','Crime',
          'Comedy','Romance','Sci-Fi'], bbox_to_anchor=(1,0.6))
plt.savefig('ratings_timeline_whitegrid.png',dpi=300,bbox_inches='tight')

```

