## importing libraries

```
%matplotlib inline
In [ ]:
         import pandas as pd
         import numpy as np
         import matplotlib.pyplot as plt
         import seaborn as sns
         from scipy import stats
         from ast import literal eval
         from sklearn.feature extraction.text import TfidfVectorizer , CountVectorizer
         from sklearn.metrics.pairwise import linear kernel , cosine similarity
         from nltk.stem.snowball import SnowballStemmer
         from nltk.stem.wordnet import WordNetLemmatizer
         from nltk.corpus import wordnet
         from surprise import Reader , Dataset , SVD
         from surprise.nodel selection import cross validate
         import warnings ; warnings.simplefilter('ignore')
```

# **Top N Recommendations**

```
In [ ]: md = pd.read_csv('movies_metasata.csv')
    md.head()
```

### **Preprocessing**

```
In []: md['genres'] = md['genres'].fillna('[]')
In []: md.head(100)
In []: md['genres'] = md['genres'].apply(literal_eval)
md.head()
In []: #Genres as List
In []: md['genres'] = md['genres'].apply(lambda x:[i['name']for i in x]if isinstance (x,list)else[])
```

```
In [ ]:
        md.head()
         md[md['vote count'].notnull()]
In [ ]:
         vote_count = md[md['vote_count'].notnull()]['vote_count'].astype('int')
In [ ]:
         vote count
In [ ]:
         vote average = md[md['vote average'].notnull()]['vote average'].astype('int')
         vote average
        top_movies = md.copy()
In [ ]:
         top movies = top movies.sort values('vote average', ascending = False).head(250)
         # No minimum votes requiremnt
         top_movies1
In [ ]:
         # Minimum number of votes 1000
In [ ]:
         top movies2 = top movies[top movies['vote count']>1000]
         top movies2
In [ ]:
         vote count = md[md['vote count'].notnull()]['vote count'].astype('int')
In [ ]: |
         vote_averages = md[md['vote_averages'].notnull()]['vote_averages'].astype('int')
         C = vote averages.mean()
         m = vote_count.quantile(0.95)
In [ ]:
         top_movies['year'] = pd.to_datetime(top_movies['release_date'],errors = 'coerce'),apply(lambda x: str(x).split(-)[0]if x
In [
         top movies
In [ ]:
         top movies3 = top movies[(top movies['vote count']>=m) & (top movies['vote count'].notnull()) & (top movies['vote average
In [ ]:
         top movies3['vote count'] = top movies3['vote count'].astype('int')
```

### **Top Movies**

```
qualified.head(15)
In [ ]:
         # Genre = Romance
In [ ]:
         genre_TM = top_movies.apply(lambda x: pd.Series(x['genres']),axis=1).stack().reset_index(level=1, drop= True)
In [ ]:
         genre TM.name = 'genre'
         genre_top_movies = top_movies.drop('genre' , axis=1).join(genre_TM)
         genre_top_movies
In [ ]:
In [ ]:
         def build_chart(genre, percentile=0.85):
             df=genre top movies[genre top movies['genre']==genre]
             vote counts = df[df['vote count'].notnull()]['vote count'].astype('int')
             vote averages = df[df['vote average'].notnull()]['vote average'].astype('int')
             C = vote averages.mean()
             m = vote counts.quantile(percentile)
             qualified = df[(df['vote count']>= m) & (df['vote count'].notnull()) & (df['vote average'].notnull())][['title' , 'ye
             qualified['vote count'] = qualified['vote count'].astype('int')
             qualified['vote average'] = qualified['vote average'].astype('int')
             qualified['wr'] = qualified.apply(lambda x: (x['vote_count']/(x['vote_count'])))
```

```
qualified = qualified.sort_values('wr' , ascending=False).head(250)
return qualified
```

## **Top Genre Movies**

```
In [ ]: built_chart('Animation').head(10)
In [ ]: built_chart('Family').head(10)
In [ ]: built_chart('Action').head(10)
```

#### **Content Based Recommender**

```
links small = pd.read csv('links samll csv')
         links small = links samll[links samll['tmdbId'].notnull()]['tmdbId'].astype('int')
         top movies = top movies.drop([19730,29503,35587])
In [ ]:
         top movies['id'] = top movies['id'].astype('int')
In [ ]:
         top movies4 = top movies[top movies['id'].isin(links small)]
In [ ]:
         top movies4.shape
         top movies4(head)
In [ ]:
         top movies4['tagline'] = top movies4['tagline'].fillana('')
         top movies4['description'] = top movies4['overview'] + top movies4['tagline']
         top movies4['description'] = top movies4['description'].fillana('')
        tf = TfidfVectorizer(analyzer='word' , ngram range=(1,2),min df=0, stop words='end')
In [ ]:
         tfidf matrix = ts.fit transform(top movies4['description'])
         tfidf matrix
         tfidf matrix.shape
In [ ]:
```

```
cosine_sim = linear_kernel(tfidf_matrix , tfidf_matrix)
         cosine sim
In [ ]:
         cosine_sim[0]
In [ ]:
         top_movies4 = top_movies4.reset_index()
In [ ]:
         titles = top movies4['title']
         indices = pd.Series(top movies4.index, index=top movies4['title'])
         def get recommendations(title):
In [ ]:
             idx = indices[title]
             sim scores = list(enumerate(cosine sim[idx]))
             sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)
             sim scores = sim scores[1:31]
             movie indices = [i[0] for i in sim scores]
             return titles.iloc[movie indices]
         get recommendations('GoldenEye').head(10)
         get_recommendations('The Appartment').head(10)
In [ ]:
         get recommendations('The Godfather').head(10)
In [
         get recommendations('The Dark Knight').head(10)
In [ ]:
In [ ]:
In [ ]:
```