MA 544 Final Project

By: Nick, Gokul, Samruth

Introduction

Purpose:

- Which type of recommendation system is best for suggesting new media to a user?

Dataset:

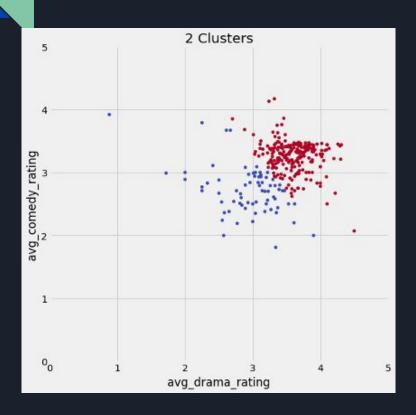
- The 100k MovieLens dataset will be used for this experiment
- https://grouplens.org/datasets/movielens/

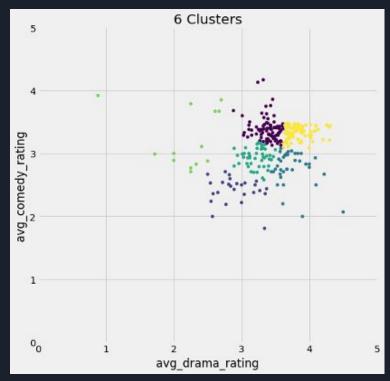
K - Means Clustering

K means Clustering is a method of vector quantization that aims to partition n observations into k clusters in which each observation belongs to the cluster with the nearest mean (cluster centers or cluster centroid), serving as a prototype of the cluster.

$$\argmin_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2 = \arg\min_{\mathbf{S}} \sum_{i=1}^k |S_i| \operatorname{Var} S_i$$

K-Means: 2D Representation





K - Means Clustering System

Recommender Function via K-Means

All this information is useful, but it needs to be in a singular function. A function is created to input a user. The the function will cluster the data, find the user's groups and make reccomendations based on the group

```
movie recommender k means cluster (user movie ratings, user id= 143, cluster number = 20):
Finds recomendations for specific user based on number of clusters
Parameters
user movie ratings : pandas
    data frame of users and movies row by col .
user id : int
    userid number for recommendatin.
cluster number : int, optional
    Number of clusters. The default is 20.
Returns
list of recommended movies : list
    list of movies recommended.
## Variables
user idx = user id -1
pick from top movies = 20
```

K- Means Clustering Results

```
# User 143, 15 Clusters
start time = time.time()
movie recommender k means cluster(user movie ratings, user id= 143, cluster number = 15)
print(f'Total Time to Run: {time.time() - start time}')
Users ratings for movies in cluster:
Kill Bill: Vol. 2 (2004)
Gone with the Wind (1939)
                              5.0
Finding Neverland (2004)
Forrest Gump (1994)
Clockwork Orange, A (1971)
Name: 142, dtype: float64
User 143 has 77 movie reviews
Users top ratings:
title
Kill Bill: Vol. 2 (2004)
Gone with the Wind (1939)
Finding Neverland (2004)
Forrest Gump (1994)
                              5.0
Clockwork Orange, A (1971)
                              5.0
Name: 143, dtype: float64
Movies Average Rating in Cluster not seen by user:
Ratatouille (2007)
                                                  5.000000
African Queen, The (1951)
                                                  4.833333
Cinema Paradiso (Nuovo cinema Paradiso) (1989)
                                                  4.812500
Psycho (1960)
                                                  4.750000
Stand by Me (1986)
                                                  4.636364
dtype: float64
The list of recommended movies for user 143 is:
['Psycho (1960)', 'Third Man, The (1949)', 'Citizen Kane (1941)', 'Hoop Dreams (1994)']
Total Time to Run: 4.063762903213501
```

SVD Clustering

SVD is a matrix factorisation technique, which reduces the number of features of a dataset by reducing the space dimension from N-dimension to K-dimension (where K<N) using:

$$A = USV^T$$

Where A is a m x n utility matrix, U is a m x r orthogonal left singular matrix, S is a r x r diagonal matrix and V is a r x n diagonal right singular matrix

SVD Cluserting Code

```
In []: 1 df title = pd.read csv('../input/movie titles.csv', encoding = "ISO-8859-1", header = None, names = ['Movie Id', 'Year', 'Na
           2 df title.set index('Movie Id', inplace = True)
           3 print (df title.head(10))
In [39]: 1 reader = Reader()
           3 # get just top 100K rows for faster run time
           4 data = Dataset.load_from_df(ratings_data[['userId', 'new_movie_id', 'rating']][:], reader)
           5 #data.split(n_folds=3)
           8 cross_validate(svd, data, measures=['RMSE', 'MAE'])
Out[39]: {'test_rmse': array([0.87007861, 0.87520837, 0.87542801, 0.87252011, 0.869551 ]),
           'test_mae': array([0.6684968 , 0.67203905, 0.6709185 , 0.66903059, 0.66902215]),
           'fit time': (3.4100866317749023,
           3.498782157897949,
           3.44977068901062,
           3.565463066101074,
           3,4956393241882324),
           'test_time': (0.12617945671081543,
           0.12566661834716797.
           0.09174680709838867.
           0.08976411819458008,
           0.13264799118041992)}
In [66]: 1 f = ['count', 'mean']
           3 df_movie_summary = ratings_data.groupby('movieId')['rating'].agg(f)
           4 df_movie_summary.index = df_movie_summary.index.map(int)
           5 movie_benchmark = round(df_movie_summary['count'].quantile(θ.7),θ)
           6 drop_movie_list = df_movie_summary[df_movie_summary['count'] < movie_benchmark].index
```

```
In [67]: 1 def movie recs(user):
                user liked movies = ratings data[(ratings data['userId'] == user) & (ratings data['rating'] == 5)]
                user liked movies = user liked movies.set index('movieId')
                user liked movies = pd.merge(movies data, user liked movies, on="movieId")
                Movies liked = user liked movies['title']
                user liked movie = movies data.copv()
                user liked movie = user liked movie.reset index()
                user liked movie = user liked movie[~user liked movie['movieId'].isin(drop movie list)]
        10
        11
                # getting full dataset
                data = Dataset.load_from_df(ratings_data[['userId', 'movieId', 'rating']], reader)
        14
                trainset = data.build full trainset()
        15
                svd.fit(trainset)
        16
        17
                user liked movie['Estimate_Rating'] = user_liked_movie['movieId'].apply(lambda x: svd.predict(user, x).est)
        18
        19
                user liked movie = user liked movie.drop('movieId', axis = 1)
        20
                user liked movie = user liked movie.sort values('Estimate Rating', ascending=False)
        22
                recommended movies = user liked movie[['title', 'Estimate Rating']]
        23
        24
                print('Movies The User Likes')
        25
                print('')
                print(Movies_liked.head(20))
        27
                print('')
        28
                print('')
                print('Recommended Movies with estimated Ratings:')
                print('')
        31
                print(recommended movies.head(10))
```

SVD Clustering Results

```
1 movie recs(user = 143)
Movies The User Likes
              Forrest Gump (1994)
           Sixteen Candles (1984)
                     Shrek (2001)
            Legally Blonde (2001)
      Walk to Remember, A (2002)
             Uptown Girls (2003)
           Chasing Liberty (2004)
            50 First Dates (2004)
            13 Going on 30 (2004)
9
       Bill Cosby, Himself (1983)
10
                        Up (2009)
11
                  Fired Up (2009)
             Hangover, The (2009)
12
13
      (500) Days of Summer (2009)
14
                   Tangled (2010)
Name: title, dtype: object
Recommended Movies with estimated Ratings:
                                                  title Estimate Rating
937
        Seventh Seal, The (Sjunde inseglet, Det) (1957)
                                                                4.245808
512
                            Beauty and the Beast (1991)
                                                                4.223449
6571
                                Eastern Promises (2007)
                                                                4.192814
7467
                                         Tangled (2010)
                                                                4.174376
                                           Shrek (2001)
3194
                                                                4.164004
                                  Cool Hand Luke (1967)
975
                                                                4.158392
711
                                       Notorious (1946)
                                                                4.132487
                      Manchurian Candidate, The (1962)
966
                                                                4.123778
4076
         Harry Potter and the Chamber of Secrets (2002)
                                                                4.100718
      Like Water for Chocolate (Como agua para choco...
                                                                4.092270
```

Observations/Conclusion

In conclusion:

- SVD system is overall faster (on average it tends to take 4.5 seconds for this dataset),

- K-means is more customizable because you can change the cluster number which allows for more precision recommendations.