1 Recommendation System

▼ 1.1 Summary

Using a dataset from the MovieLens research lab at the University of Minnesota that consists of about 100,000 user ratings, a recommendation system was created to suggest movies to old and new users that they are likely to enjoy. Two main types of filtering were used to achieve the results of this project; those being collaborative filtering and content based filtering. The available datasets were previewed to see which categories were necessary, which included movie ID, title, genre, user ID, and rating. In order to suggest movies to a user, we need to know which movies the user already saw and what they (as well as others) thought about it. Several models were created and libraries were used to get the optimum results. Those include the Surprise and ALS library to perform collaborative filtering. To achieve the results for content based filtering, a memory based method that is neighborhood based was performed. In the end using the model with the best root mean squared error (KNNBaseline) this project achieved a result of RMSE: 0.5536. In short, the predictions made in this project will be about 0.55 off on a scale from 0 to 5 where this range was the rating for a movie.

1.2 Import necessary libraries

```
In [1]:
         1 # Import necessary libraries
         2  # from pyspark.sql import SparkSession
         3 from surprise import Dataset
         4 from surprise import Reader
         5 from surprise import SVD
         6 from surprise.model selection import GridSearchCV
         7 from surprise.model selection import train test split
         8 from surprise.model_selection import cross_validate
         9 from surprise.prediction_algorithms import knns
        10 from surprise.similarities import cosine, msd, pearson
        11 from surprise import accuracy
        12 import pandas as pd
        13 import numpy as np
        14 from scipy.sparse import csc matrix
        15 from scipy.sparse.linalg import svds
        16 from pyspark.sql import SparkSession
        17 from pyspark.ml.evaluation import RegressionEvaluator
        18 from pyspark.ml.recommendation import ALS
        19 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder
        20 import difflib
        21 from re import search
In [2]:
           movies = pd.read csv('data/movies.csv')
           ratings = pd.read csv('data/ratings.csv')
```

▼ 1.3 Learning more about the datasets

Here is where I learned what type of categories each dataset contained.

In [3]: 1 movies.head() Out[3]: movield title genres 0 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy 1 1 2 Jumanji (1995) Adventure|Children|Fantasy 2 3 Grumpier Old Men (1995) Comedy|Romance 4 Waiting to Exhale (1995) Comedy|Drama|Romance 3 Father of the Bride Part II (1995) Comedy

```
In [4]:
               ratings.head()
Out[4]:
              userld
                      movield rating
                                      timestamp
           0
                                      964982703
                            1
                                 4.0
                            3
                                 4.0
                                      964981247
           1
                            6
                                      964982224
           2
                   1
                                 4.0
           3
                           47
                                      964983815
                                      964982931
                   1
                           50
                                 5.0
```

Viewing the shape of each dataset in order to get an idea how many rows and columns I would be dealing with once they were merged.

```
In [5]: 1 movies.shape
Out[5]: (9742, 3)
In [6]: 1 ratings.shape
Out[6]: (100836, 4)
In [7]: 1 df = pd.merge(movies, ratings, on='movieId')
2 df.head()
```

Out[7]: movield title userld rating timestamp genres Adventure|Animation|Children|Comedy|Fantasy 0 1 Toy Story (1995) 964982703 4.0 1 1 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy 5 4.0 847434962 1 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy 7 1106635946 2 4.5 3 Toy Story (1995) Adventure|Animation|Children|Comedy|Fantasy 15 1510577970 Adventure|Animation|Children|Comedy|Fantasy 1 Toy Story (1995) 17 1305696483

```
In [8]:
         1 df.info()
        <class 'pandas.core.frame.DataFrame'>
        Int64Index: 100836 entries, 0 to 100835
       Data columns (total 6 columns):
                       Non-Null Count
            Column
                                       Dtype
            _____
                       _____
        0
           movieId
                       100836 non-null int64
        1
            title
                       100836 non-null object
        2
            genres
                       100836 non-null object
        3
            userId
                       100836 non-null int64
                       100836 non-null float64
            rating
           timestamp 100836 non-null int64
       dtypes: float64(1), int64(3), object(2)
       memory usage: 5.4+ MB
```

To be safe, checking that there were no NaN values.

1.4 Creating df's that will be useful later on

```
In [10]:
         1 df.drop('timestamp', axis=1, inplace= True)
In [11]:
         1 df.info()
        <class 'pandas.core.frame.DataFrame'>
        Int64Index: 100836 entries, 0 to 100835
        Data columns (total 5 columns):
         #
             Column Non-Null Count
                                     Dtype
             _____
                     -----
                                     ____
         0
            movieId 100836 non-null int64
            title 100836 non-null object
         1
             genres 100836 non-null object
             userId 100836 non-null int64
             rating 100836 non-null float64
        dtypes: float64(1), int64(2), object(2)
        memory usage: 4.6+ MB
```

Sorting by movie ID here for mostly preference and organization of the dataset.

```
In [12]: 1 # Sorting by first movie
2 df.sort_values(by='movieId', ascending=True)
3 df.head(10)
```

Out[12]:

Out[1

	movield	title	genres	userId	rating
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1	4.0
1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	5	4.0
2	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	7	4.5
3	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	15	2.5
4	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	17	4.5
5	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	18	3.5
6	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	19	4.0
7	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	21	3.5
8	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	27	3.0
9	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	31	5.0

```
In [13]: 1 df.shape
Out[13]: (100836, 5)
```

▼ 1.4.1 Creating main dataset

```
In [14]: 1 # Separating Genres and making them into list
2 genre = df['genres'].map(lambda x: x.split('|'))
3 df['genre'] = genre
4 df.head(10)
```

14]:	n	novield	title	genres	userld	rating	genre
	0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	1	4.0	[Adventure, Animation, Children, Comedy, Fantasy]
	1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	5	4.0	[Adventure, Animation, Children, Comedy, Fantasy]
	2	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	7	4.5	[Adventure, Animation, Children, Comedy, Fantasy]
	3	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	15	2.5	[Adventure, Animation, Children, Comedy, Fantasy]
	4	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	17	4.5	[Adventure, Animation, Children, Comedy, Fantasy]
	5	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	18	3.5	[Adventure, Animation, Children, Comedy, Fantasy]
	6	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	19	4.0	[Adventure, Animation, Children, Comedy, Fantasy]
	7	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	21	3.5	[Adventure, Animation, Children, Comedy, Fantasy]
	8	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	27	3.0	[Adventure, Animation, Children, Comedy, Fantasy]
	9	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	31	5.0	[Adventure, Animation, Children, Comedy, Fantasy]

Turning the genres into a list so that it able to be manipulated easier if need be.

[Adventure, Animation, Children, Comedy, Fantasy]

[Adventure, Animation, Children, Comedy, Fantasy]

3

```
In [15]:
                 df.drop('genres', axis=1, inplace= True)
                 df.head()
Out[15]:
                movield
                                    title userld rating
                                                                                                genre
                      1 Toy Story (1995)
                                                         [Adventure, Animation, Children, Comedy, Fantasy]
             0
                                               1
             1
                      1 Toy Story (1995)
                                              5
                                                         [Adventure, Animation, Children, Comedy, Fantasy]
             2
                      1 Toy Story (1995)
                                              7
                                                        [Adventure, Animation, Children, Comedy, Fantasy]
```

1.4.2 Creating Ratings Columns

1 Toy Story (1995)

1 Toy Story (1995)

Creating another data frame here to get more information about the ratings.

15

17

```
In [16]: 1 average_ratings = pd.DataFrame(df.groupby('title')['rating'].mean())
2 average_ratings.sort_values(by='rating', ascending=False)
3
```

```
Out[16]:
                                                        rating
                                                  title
                            Gena the Crocodile (1969)
                                                          5.0
                                                          5.0
                                   True Stories (1986)
                       Cosmic Scrat-tastrophe (2015)
                                                          5.0
                                                          5.0
                             Love and Pigeons (1985)
                                                          5.0
                Red Sorghum (Hong gao liang) (1987)
                                                          0.5
                               Don't Look Now (1973)
              Journey 2: The Mysterious Island (2012)
                                                          0.5
                     Joe Dirt 2: Beautiful Loser (2015)
                                                          0.5
                                                          0.5
```

Jesus Christ Vampire Hunter (2001)

Fullmetal Alchemist 2018 (2017)

9719 rows × 1 columns

By sorting by highest rated, I could have an idea of what the best movies looked like, but by looking at this I could tell something was off. Some of these movies could have just a few people who gave it a 5.0. I needed to know the just how many people voted for these movies at this point.

Checking here for the total amount of ratings placed on all the movies.

0.5

```
In [18]:
                # Mean rating of all movie averages
               mean_rating = average_ratings['rating'].mean()
               print(mean rating)
           3.2623883953257353
In [19]:
            1 average ratings.head()
Out[19]:
                                             rating Total Ratings
                                        title
                                    '71 (2014)
                                               4.0
            'Hellboy': The Seeds of Creation (2004)
                                               4.0
                                                             1
                         'Round Midnight (1986)
                                               3.5
                                                             2
                                               5.0
                                                             1
                            'Salem's Lot (2004)
                       'Til There Was You (1997)
                                               4.0
                                                             2
In [21]:
                # Most rated movie
               most_rated = average_ratings.sort_values(by='Total Ratings', ascending=False)
            3 most rated.head(10)
Out[21]:
                                                   rating Total Ratings
                                           title
                                                4.164134
                                                                 329
                              Forrest Gump (1994)
                 Shawshank Redemption, The (1994)
                                                4.429022
                                                                 317
                                                                 307
                                                4.197068
```

```
Pulp Fiction (1994)
         Silence of the Lambs, The (1991)
                                           4.161290
                                                               279
                                           4.192446
                                                               278
                        Matrix, The (1999)
                                                               251
Star Wars: Episode IV - A New Hope (1977)
                                           4.231076
                                                               238
                      Jurassic Park (1993)
                                           3.750000
                                           4.031646
                                                               237
                        Braveheart (1995)
       Terminator 2: Judgment Day (1991)
                                           3.970982
                                                               224
                                                               220
                    Schindler's List (1993)
                                          4.225000
```

After finding the mean of each movie's rating as well as each movie's total rating, the data frame looked better. Sorting now by total ratings, these movies looked more familiar and are known as some of the best.

```
In [22]:    1  # Number of total ratings
    2  most_rated['Total Ratings'].sum()
Out[22]: 100836
```

1.4.3 Finding out more about the ratings

```
In [24]:
           1 df.rating.value_counts()
Out[24]: 4.0
                26818
         3.0
                20047
         5.0
                13211
         3.5
                13136
         4.5
                 8551
         2.0
                 7551
         2.5
                 5550
         1.0
                 2811
         1.5
                 1791
         0.5
                 1370
         Name: rating, dtype: int64
In [25]:
          1 lower rating = df['rating'].min()
           2 upper_rating = df['rating'].max()
           3 print('User rating range: {0} to {1}'.format(lower_rating, upper_rating))
```

Finding the min and max rating for all the movies was important here because this information is needed for the Surprise library calculations.

```
In [26]: 1 # Create Surprise Dataset
2 reader = Reader(rating_scale=(0.5,5))
3 data = Dataset.load_from_df(df[['userId', 'movieId', 'rating']], reader)
```

User rating range: 0.5 to 5.0

```
In [27]:
          1 svd = SVD(verbose= True, n epochs=10)
          2 cross validate(svd, data, measures= ['RMSE', 'MAE'], cv=3, verbose= True)
         Processing epoch 0
         Processing epoch 1
         Processing epoch 2
         Processing epoch 3
         Processing epoch 4
         Processing epoch 5
         Processing epoch 6
         Processing epoch 7
         Processing epoch 8
         Processing epoch 9
         Processing epoch 0
         Processing epoch 1
         Processing epoch 2
         Processing epoch 3
         Processing epoch 4
         Processing epoch 5
         Processing epoch 6
         Processing epoch 7
         Processing epoch 8
         Processing epoch 9
         Processing epoch 0
         Processing epoch 1
         Processing epoch 2
         Processing epoch 3
         Processing epoch 4
         Processing epoch 5
         Processing epoch 6
         Processing epoch 7
         Processing epoch 8
         Processing epoch 9
         Evaluating RMSE, MAE of algorithm SVD on 3 split(s).
                           Fold 1 Fold 2 Fold 3 Mean
         RMSE (testset)
                           0.8871 0.8849 0.8883 0.8868 0.0014
         MAE (testset)
                           0.6845 0.6832 0.6848 0.6842 0.0007
         Fit time
                           1.73
                                   1.74
                                           1.76
                                                   1.74
                                                           0.01
         Test time
                           0.17
                                   0.22
                                           0.23
                                                   0.21
                                                           0.03
Out[27]: {'test rmse': array([0.88713344, 0.88493684, 0.88828583]),
          'test mae': array([0.68454389, 0.68318737, 0.68478607]),
          'fit_time': (1.7331767082214355, 1.7393770217895508, 1.7623910903930664),
          'test time': (0.17146825790405273, 0.22491097450256348, 0.22838091850280762)}
```

This was my first go around of using the Surprise library features. Basically a test run for what range the RMSE I should be expecting.

1.4.4 Creating train and test sets

```
In [28]: 1 # Split into train and test set
2 trainset, testset = train_test_split(data, test_size=0.2)
In [29]: 1 print('Type trainset : ', type(trainset),'\n')
2 print('Type testset :',type(testset))
Type trainset : <class 'surprise.trainset.Trainset'>
Type testset : <class 'list'>
```

```
In [30]:
          1 print(len(testset))
          2 print(testset[0])
         20168
         (414, 30707, 4.5)
In [31]:
          1 print('Number of users: ', trainset.n_users, '\n')
           2 print('Number of items: ', trainset.n_items, '\n')
           3
         Number of users: 610
         Number of items: 8935
         Getting an idea here of whether I should set user_based to True or False. The higher number would result in more
         computing time.
          1 # Because of fewer users than items, it is more efficient to calculate
In [32]:
          2 # user-user similarity rather than item-item
          3 # However since this is set to false, this is using item-item
           4 sim_cos = {'name':'cosine', 'user_based':False}
In [33]:
          basic = knns.KNNBasic(sim_options=sim_cos)
          2 # Fit model
           3 basic.fit(data.build full trainset())
         Computing the cosine similarity matrix...
         Done computing similarity matrix.
Out[33]: <surprise.prediction_algorithms.knns.kNnBasic at 0x7f9071303040>
In [34]:
          1 # Similarity metrics of each of the movie to one another
           2 basic.sim
Out[34]: array([[1.
                            , 0.96446411, 0.97154149, ..., 0.
                                                                      , 0.
                 0.
                           ],
                [0.96446411, 1.
                                        , 0.93890126, ..., 0.
                                                                      , 0.
                 0.
                [0.97154149, 0.93890126, 1.
                                                   , ..., 0.
                                                                      , 0.
                 0.
                           ],
                 . . . ,
                [0.
                            , 0.
                                        , 0. , ..., 1.
                                                                      , 1.
```

This was the first similarity metric produced, but this didn't look like much to me so I thought I would understand it better by putting it in a data frame.

, 0.

```
In [35]: 1 basic.sim.shape
Out[35]: (9724, 9724)
```

, 0.

, 0.

0.

.01

0.

.01

], , 0.

], , 0.

]])

```
basic df = pd.DataFrame(basic.sim)
In [36]:
               basic df.head()
Out[36]:
                                                                                     7
                     0
                                                          4
                                                                   5
                                                                            6
                                                                                              8
                                                                                                             9714 9715
            0 1.000000 0.964464 0.971541 0.983870 0.955044
                                                            0.966239
                                                                     0.944239
                                                                              0.988721 0.950356
                                                                                                0.952135 ...
                                                                                                               0.0
                                                                                                                     0.0
             0.964464 1.000000 0.938901 0.960988
                                                  0.968546
                                                           0.954224 0.968611
                                                                              0.969442 0.943370 0.946736 ...
                                                                                                                     0.0
            2 0.971541 0.938901 1.000000
                                         1.000000 0.972298
                                                            0.933685 0.971202
                                                                             0.973985
                                                                                       0.942695
                                                                                                 0.921344
                                                                                                               0.0
                                                                                                                     0.0
            3 0.983870 0.960988 1.000000
                                         1.000000
                                                  0.972003
                                                            1.000000
                                                                     0.943293
                                                                              1.000000
                                                                                        0.000000
                                                                                                 0.973708 ...
                                                                                                               0.0
                                                                                                                     0.0
              0.955044 0.968546 0.972298 0.972003 1.000000 0.966924 0.974637 0.960928
                                                                                       0.963287
                                                                                                 0.952833 ...
                                                                                                               0.0
                                                                                                                     0.0
```

5 rows × 9724 columns

This helped me understand that it was comparing each movie to one another and producing a number of how similar that movie was to the other. The highest number meaning it was very similar.

1.4.5 Finding the model that will give the lowest RMSE

```
In [37]: 1 # Test model to determine how well it performed
2 predictions = basic.test(testset)
3 print(accuracy.rmse(predictions))
```

RMSE: 0.9031 0.9030818565366927

Results mean model is off by about 0.9013 points

```
In [38]: 1 # Pearson Correlation
2 sim_pearson = { 'name':'pearson', 'user_based':False}
3 basic_pearson = knns.KNNBasic(sim_options=sim_pearson)
4 basic_pearson.fit(data.build_full_trainset())
5 predictions = basic_pearson.test(testset)
6 print(accuracy.rmse(predictions))
```

Computing the pearson similarity matrix...

Done computing similarity matrix.

RMSE: 0.6960

0.6959736433654313

```
In [39]: 1 # KNN with Means. Same as basic KNN model but takes into account the mean rating
2 # of each user/item
3 sim_pearson = {'name':'pearson', 'user_based':True}
4 knn_means = knns.KNNWithMeans(sim_options=sim_pearson)
5 knn_means.fit(data.build_full_trainset())
6 predictions = knn_means.test(testset)
7 print(accuracy.rmse(predictions))
```

Computing the pearson similarity matrix... Done computing similarity matrix. RMSE: 0.6053
0.6053176017550091

```
In [40]: 1 # KNNBaseline method - This adds in bias term that is calculated by way of
2 # minimizing a cost function
3 sim_pearson = {'name':'pearson', 'user_based':False}
4 knn_baseline = knns.KNNBaseline(sim_options=sim_pearson)
5 knn_baseline.fit(data.build_full_trainset())
6 predictions = knn_baseline.test(testset)
7 print(accuracy.rmse(predictions))
```

Estimating biases using als...

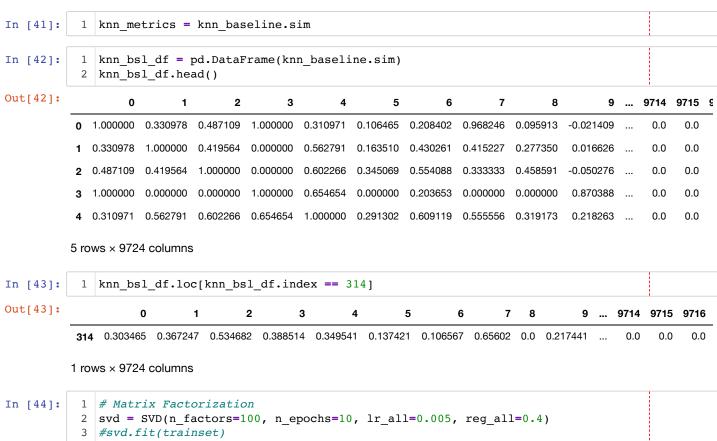
Computing the pearson similarity matrix...

Done computing similarity matrix.

RMSE: 0.5536

0.5536483298449599

KNNBaseline method had the best results with 0.5941 RMSE (user_based set to True). And RMSE: 0.5531 when user_based set to False.



```
In [44]: 1 # Matrix Factorization
2 svd = SVD(n_factors=100, n_epochs=10, lr_all=0.005, reg_all=0.4)
3 #svd.fit(trainset)
4 svd.fit(data.build_full_trainset())
5 predictions = svd.test(testset)
6 print(accuracy.rmse(predictions))
```

RMSE: 0.8577 0.857650177663376

The above methods tested were to find the right combination that would produce the lowest RMSE. Turns out that was KNNBaseline.

```
In [45]: 1 # Make Predictions
2  user_3_prediction = svd.predict('3', '25')
3  user_3_prediction

Out[45]: Prediction(uid='3', iid='25', r_ui=None, est=3.501556983616962, details={'was_impossible':
```

False})

Here I was getting familiar with how the prediction feature worked to apply later on.

```
In [46]:
              # Retrieve movie name
             def name retriever(movie id):
           3
                  for movie in df:
                     return df.loc[df['movieId']==movie id].title.values[0]
In [47]:
           1 name_retriever(318)
Out[47]: 'Shawshank Redemption, The (1994)'
         I determined a function needed to be created to turn movie ID's into their actual movie name to lessen the confusion.
In [48]:
           1 # Pandas df of unique movie ID's
             unique iids = df['movieId'].drop duplicates(keep='first')
           1 # Number of different movies
In [49]:
           2 unique_iids.shape
Out[49]: (9724,)
In [24]:
           1 # Number of different users
           2 unique users = df['userId'].drop_duplicates(keep='first')
           3 unique users.shape
Out[24]: (610,)
         A dataset of unique users and movies needed to be created so that it could be compared to later on.
In [50]:
           1 # Find unrated movies by a user function
              def find unrated(user, movie df):
                  movId = df.loc[df['userId'] == user, 'movieId']
           3
                  # Remove the iids that user n has rated from the list of all movie ids
           5
                  user_unrated = pd.concat([unique_iids,movId]).drop_duplicates(keep=False)
                  user unrated = user unrated.to frame()
           6
           7
                  return user unrated.head()
           8
In [51]:
           1 find unrated(4, df)
Out[51]:
              movield
            0
                   1
                   2
          215
                   3
          325
          377
                    4
                   5
          384
```

This function was very important and necessary to create to solve the first part of this project. I needed to know which movies any user in this dataset did not see. It wouldn't make much sense to try to predict a movie rating that a user already rated.

1.5 Recommending movies for user 4 using SVD

Here is where I put my final idea into practice. Sorting from a list of estimated ratings of movies the user has not seen to produce a top 5 movie recommendation.

```
In [52]:
            1 # Recommending movies for user 3 using SVD
            2 svd = SVD(n factors=100, n epochs=10, lr all=0.005, reg all=0.4)
            3 svd.fit(data.build full trainset())
Out[52]: <surprise.prediction algorithms.matrix factorization.SVD at 0x7f8f1f085250>
In [53]:
              unique ids = df['movieId'].unique()
              #Get list of unique ids for user 4
              newdf = df.loc[df['userId'] == 4 , 'movieId']
              # Remove rated movies
              movies to predict = np.setdiff1d(unique ids, newdf)
In [54]:
           1 # Make prediction of unrated movie
              # svd.predict(uid='3',iid='5')
In [55]:
           1
              recs = []
              for iid in movies_to_predict:
                   recs.append((iid, svd.predict(uid='4',iid=iid).est))
              top 5 rec = pd.DataFrame(recs,columns=['iid','predictions']).sort values('predictions',
              top_5_rec
Out[55]:
                 iid predictions
           256
                318
                      4.101610
                750
                      4.066828
           557
           824
               1204
                      4.059426
                858
                      4.016887
           612
                 50
                      4.016261
            42
In [56]:
            1 top_5_rec['Movie'] = top_5_rec['iid'].apply(name_retriever)
              top 5 rec
Out[56]:
                 iid predictions
                                                             Movie
                318
                      4.101610
                                        Shawshank Redemption, The (1994)
           256
           557
                750
                      4.066828 Dr. Strangelove or: How I Learned to Stop Worr...
               1204
                      4.059426
                                               Lawrence of Arabia (1962)
           824
                                                  Godfather, The (1972)
           612
                858
                      4.016887
            42
                 50
                      4.016261
                                              Usual Suspects, The (1995)
```

1.6 Recommending movies for user 4 using KNNBaseline model

Now that I knew the idea could work, I applied the same function but with KNNBaseline this time to get an even better top 5 recommendation for any user.

```
In [57]:
           1 cross validate(knn baseline, data, measures= ['RMSE', 'MAE'], cv=3, verbose= True)
          Estimating biases using als...
          Computing the pearson similarity matrix...
          Done computing similarity matrix.
          Estimating biases using als...
          Computing the pearson similarity matrix...
          Done computing similarity matrix.
          Estimating biases using als...
          Computing the pearson similarity matrix...
          Done computing similarity matrix.
          Evaluating RMSE, MAE of algorithm KNNBaseline on 3 split(s).
                             Fold 1 Fold 2 Fold 3 Mean
                                                                Std
                             0.8824
                                      0.8843
                                              0.8849
          RMSE (testset)
                                                       0.8839
                                                                0.0011
                                              0.6807 0.6791 0.0012
          MAE (testset)
                             0.6777
                                     0.6790
          Fit time
                             13.05
                                      13.15
                                               12.82
                                                       13.01
                                                                0.14
          Test time
                             8.13
                                      7.86
                                               7.89
                                                       7.96
                                                                0.12
Out[57]: {'test rmse': array([0.88239329, 0.88429518, 0.88491002]),
           'test mae': array([0.67769519, 0.67903973, 0.68065938]),
           'fit time': (13.04737401008606, 13.151367902755737, 12.81882095336914),
           'test time': (8.130614995956421, 7.863303184509277, 7.889848232269287)}
In [58]:
           1 recs = []
              for iid in movies to predict:
                  recs.append((iid, knn baseline.predict(uid='4',iid=iid).est))
              top 5 rec = pd.DataFrame(recs,columns=['iid','predictions']).sort values('predictions',
              top 5 rec
Out[58]:
                iid predictions
           256
               318
                      4.403969
          612
               858
                      4.328440
               1204
                      4.286615
           824
                750
                      4.269324
           557
                50
                      4.250099
           42
In [59]:
           1 top 5 rec['Movie'] = top 5 rec['iid'].apply(name retriever)
              top_5_rec
Out[59]:
                iid predictions
                                                           Movie
                      4.403969
                                       Shawshank Redemption, The (1994)
               318
           256
          612
               858
                      4.328440
                                                 Godfather, The (1972)
              1204
                      4.286615
                                             Lawrence of Arabia (1962)
          824
           557
                750
                      4.269324 Dr. Strangelove or: How I Learned to Stop Worr...
                50
                      4.250099
                                             Usual Suspects, The (1995)
           42
```

From this top 5 list, it is slightly different from that of the SVD method.

1.7 Using ALS

For sanity purposes, I also tried out the ALS method just in case it had a better RMSE score.

```
In [60]:
          1 print('Using ALS')
             bsl options = {'method': 'als',
                             'n epochs': 5,
          4
                             'reg u': 12,
                             'reg i': 5
          5
          6
                            }
             algo = knns.KNNBaseline(bsl options=bsl options)
            cross validate(algo, data, measures=['RMSE'], cv=3, verbose=False)
         Using ALS
         Estimating biases using als...
         Computing the msd similarity matrix...
         Done computing similarity matrix.
         Estimating biases using als...
         Computing the msd similarity matrix...
         Done computing similarity matrix.
         Estimating biases using als...
         Computing the msd similarity matrix...
         Done computing similarity matrix.
Out[60]: {'test rmse': array([0.87906848, 0.8800158, 0.88121457]),
          'fit time': (0.15366196632385254, 0.1632378101348877, 0.16462492942810059),
          'test time': (1.8678550720214844, 1.837568998336792, 1.9958178997039795)}
In [61]:
          1 # User train test split to sample a trainset and testset
          2 algo = knns.KNNBaseline(bsl options=bsl options)
          3 # Use full dataset
          4 predictions = algo.fit(data.build full trainset()).test(testset)
            accuracy.rmse(predictions)
         Estimating biases using als...
         Computing the msd similarity matrix...
         Done computing similarity matrix.
         RMSE: 0.6757
Out[61]: 0.675702953578366
```

RMSE using ALS turned out to be 0.6743 using KNNBaseline.

2 Fix to cold start problem: User input & Content Based Filtering

To solve this problem, I thought to myself what does companies like Netflix and Hulu currently do for their new users? The answer is they ask them what type of genre, show, or movies they like. In this case, I sorted earlier in this project the highest rated movies and gave the user the option to choose from the top ten. Chances are, the user has seen and possibly liked these very popular movies.

```
In [62]: 1 top_10_rated = most_rated.head(10).index.tolist()
```

```
1 | print('Here are the top 10 most rated movies.\n')
In [63]:
          2 print('Which one of these movies do you like? \n')
             print(top 10 rated)
         Here are the top 10 most rated movies.
         Which one of these movies do you like?
         ['Forrest Gump (1994)', 'Shawshank Redemption, The (1994)', 'Pulp Fiction (1994)', 'Silenc
         e of the Lambs, The (1991)', 'Matrix, The (1999)', 'Star Wars: Episode IV - A New Hope (19
         77)', 'Jurassic Park (1993)', 'Braveheart (1995)', 'Terminator 2: Judgment Day (1991)', "S
         chindler's List (1993)"]
In [64]:
          1 movie name = input('Enter your favorite movie : ')
         Enter your favorite movie : Schindler's List (1993)
In [65]:
          1 # Create list with all movies in dataset
          2 list of titles = df['title'].drop duplicates().tolist()
          3 #print(list of titles)
In [66]:
          1 # Get index of movie user likes
          2 #Get movie ID instead
          3 def get movie index from title(title):
                 return df[df['title']==movie_name].index.values[0]
In [67]:
          1 movie_index = get_movie_index_from_title(movie_name)
In [86]:
          1 old index = df[df['title'] == movie name].index.values[0]
          2 old_index
Out[86]: 14106
```

Get the index of the movie the user selected from the main dataset. The problem here is that this index number means nothing when compared to the movie index from the similarity metric.

	index	movield
0	0	1
1	215	2
2	325	3
3	377	4
4	384	5
9719	100831	193581
9720	100832	193583
9721	100833	193585
9722	100834	193587
9723	100835	193609
	1 2 3 4 9719 9720 9721 9722	0 0 1 215 2 325 3 377 4 384 9719 100831 9720 100832 9721 100833

Another problem that needed to be solved was the index number was too high. 14106 is way more than the actual unique amount of movies. Thus, each index needed another index that would be the same amount in the similarity metric array.

9724 rows × 2 columns

```
In [89]:
            1 new index =reset ind[reset ind['index']==old index].index.values[0]
               new index
Out[89]: 461
          Now I have an index I can work with that I can link back to the original movie the user choose.
In [88]:
            1 df['title'][df.index == old_index]
Out[88]: 14106
                     Schindler's List (1993)
          Name: title, dtype: object
In [72]:
            1 df.loc[df['title'] == movie name]
Out[72]:
                  movield
                                        title userId rating
                                                               genre
            14106
                      527
                          Schindler's List (1993)
                                                      5.0
                                                          [Drama, War]
           14107
                          Schindler's List (1993)
                     527
                                                 3
                                                      0.5
                                                         [Drama, War]
           14108
                          Schindler's List (1993)
                                                 5
                                                          [Drama, War]
                          Schindler's List (1993)
                                                         [Drama, War]
           14109
                     527
                                                 6
                                                      3.0
            14110
                     527
                          Schindler's List (1993)
                                                 8
                                                      5.0
                                                          [Drama, War]
            14321
                     527 Schindler's List (1993)
                                               603
                                                      3.0
                                                         [Drama, War]
            14322
                                                         [Drama, War]
                          Schindler's List (1993)
                                               606
                                                      5.0
            14323
                          Schindler's List (1993)
                                               607
                                                         [Drama, War]
                                                      5.0
            14324
                          Schindler's List (1993)
                                               608
                                                          [Drama, War]
                                                      4.0
           14325
                     527 Schindler's List (1993)
                                               610
                                                      3.5 [Drama, War]
          220 rows × 5 columns
In [90]:
               # Generating similar movies matrix
               similar movies = list(enumerate(knn metrics[new index]))
In [74]:
               # Sorting movies based on simmularity score
               # Sort movies in descending order and get second value of tuple
               sorted similar movies = sorted(similar movies, key= lambda x:x[1], reverse = True)
               #print(sorted similar movies)
           Using the new index created, this was matched with the same index in the similarity metric array.
In [75]:
               def get title from index(index):
                    orig_id = reset_ids[reset_ids.index == index]["index"].values[0]
            3
                    return df[df.index==orig id].title.values[0]
In [91]:
            1 get title from index(461)
Out[91]: "Schindler's List (1993)"
```

Another function needed to be created now to get the movie title from the new index I have been using.

```
Here are 10 movies we think you may like because you chose Schindler's List (1993). Eye for an Eye (1996)
Amazing Panda Adventure, The (1995)
Blue in the Face (1995)
Party Girl (1995)
Farinelli: il castrato (1994)
Love Affair (1994)
Priest (1994)
Forrest Gump (1994)
Blue Chips (1994)
Lassie (1994)
Puppet Masters, The (1994)
```

To put it all back together, this loop is created to go through all the similar movie (that have been sorted at this point from highest similarity score to lowest) and apply the function created above to turn that new index into a movie title.

3 Conclusion

Through a lot of trial and error, the optimum dataset and models were created to obtain the best result. Information I thought I would need like genre or minimum votes to accept in my model were ultimately removed because it was causing too much clutter in the final results.

A dataset of just user ID, movie ID, rating, and title were all that were needed. Experimenting between which similarity metrics would serve the model better also turned out to be fruitful in terms of the least amount of error. A item-item similarity metrics was ultimately selected to solve the cold start problem despite a slightly longer computing time.

In addition, multiple methods were tested to find the best RMSE. Those included the Surprise library K nearest neighbor algorithms such as KNNBasic, Pearson Correlation, KNNWithMeans, KNNBaseline, and SVD matrix factorization. Furthermore, the Alternate Least Squares (ALS) method was also attempted. Through all these models, the method that showed to lowest RMSE was the KNNBaseline method where it used an item to item similarity metric. The result of this was a RMSE of 0.5536 which was far less than the other RMSE's.

Putting it all together the first part of this project was solved using collaborative filtering. By providing a user ID already in the dataset with a movie ID the user has not seen yet, an estimate can be predicted on what that specific user may rate that movie. Once that is figured out, a simple sort to find the highest estimate a user may give a movie and be returned.

As mentioned earlier, a similarity metric was created to help solve the cold start problem. Content based filter was used in this case, or in other words if a user likes a certain movie, they may like similar movies to that one. The metric created was helpful because it was an array of how alike each movie was to one another. From there, the new user is asked which out of the top rated movies have they seen and liked. In this project case, the user selected Schindler's List and the movies most similar to that one were returned.