

1. Business Problem

1.1 Problem Description

Netflix is all about connecting people to the movies they love. To help customers find those movies, they developed world-class movie recommendation system: CinematchSM. Its job is to predict whether someone will enjoy a movie based on how much they liked or disliked other movies. Netflix use those predictions to make personal movie recommendations based on each customer's unique tastes. And while **Cinematch** is doing pretty well, it can always be made better.

Now there are a lot of interesting alternative approaches to how Cinematch works that netflix haven't tried. Some are described in the literature, some aren't. We're curious whether any of these can beat Cinematch by making better predictions. Because, frankly, if there is a much better approach it could make a big difference to our customers and our business.

Credits: <https://www.netflixprize.com/rules.html>

1.2 Problem Statement

Netflix provided a lot of anonymous rating data, and a prediction accuracy bar that is 10% better than what Cinematch can do on the same training data set. (Accuracy is a measurement of how closely predicted ratings of movies match subsequent actual ratings.)

1.3 Sources

- <https://www.netflixprize.com/rules.html>
- <https://www.kaggle.com/netflix-inc/netflix-prize-data>
- Netflix blog: <https://medium.com/netflix-techblog/netflix-recommendations-beyond-the-5-stars-part-1-55838468f429> (very nice blog)
- surprise library: <http://surpriselib.com/> (we use many models from this library)
- surprise library doc: http://surprise.readthedocs.io/en/stable/getting_started.html (we use many models from this library)
- installing surprise: <https://github.com/NicolasHug/ Surprise#installation>
- Research paper: <http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf> (most of our work was inspired by this paper)
- SVD Decomposition : <https://www.youtube.com/watch?v=P5mlg91as1c>

1.4 Real world/Business Objectives and constraints

Objectives:

1. Predict the rating that a user would give to a movie that he has not yet rated.
2. Minimize the difference between predicted and actual rating (RMSE and MAPE)

Constraints:

1. Some form of interpretability.

2. Machine Learning Problem

2.1 Data

2.1.1 Data Overview

Get the data from : <https://www.kaggle.com/netflix-inc/netflix-prize-data/data>

Data files :

Data files :

- combined_data_1.txt
- combined_data_2.txt
- combined_data_3.txt
- combined_data_4.txt
- movie_titles.csv

The first line of each file [combined_data_1.txt, combined_data_2.txt, combined_data_3.txt, combined_data_4.txt] contains the movie id followed by a colon. Each subsequent line in the file corresponds to a rating from a customer and its date in the following format:

CustomerID,Rating,Date

MovieIDs range from 1 to 17770 sequentially.

CustomerIDs range from 1 to 2649429, with gaps. There are 480189 users.

Ratings are on a five star (integral) scale from 1 to 5.

Dates have the format YYYY-MM-DD.

2.1.2 Example Data point

1:

```
1488844,3,2005-09-06
822109,5,2005-05-13
885013,4,2005-10-19
30878,4,2005-12-26
823519,3,2004-05-03
893988,3,2005-11-17
124105,4,2004-08-05
1248029,3,2004-04-22
1842128,4,2004-05-09
2238063,3,2005-05-11
1503895,4,2005-05-19
2207774,5,2005-06-06
2590061,3,2004-08-12
2442,3,2004-04-14
543865,4,2004-05-28
1209119,4,2004-03-23
804919,4,2004-06-10
1086807,3,2004-12-28
1711859,4,2005-05-08
372233,5,2005-11-23
1080361,3,2005-03-28
1245640,3,2005-12-19
558634,4,2004-12-14
2165002,4,2004-04-06
1181550,3,2004-02-01
1227322,4,2004-02-06
427928,4,2004-02-26
814701,5,2005-09-29
808731,4,2005-10-31
662870,5,2005-08-24
337541,5,2005-03-23
786312,3,2004-11-16
1133214,4,2004-03-07
1537427,4,2004-03-29
1209954,5,2005-05-09
3335563,3,2005-06-16
```

```
2381599,3,2005-09-12
525356,2,2004-07-11
1910569,4,2004-04-12
2263586,4,2004-08-20
2421815,2,2004-02-26
1009622,1,2005-01-19
1481961,2,2005-05-24
401047,4,2005-06-03
2179073,3,2004-08-29
1434636,3,2004-05-01
93986,5,2005-10-06
1308744,5,2005-10-29
2647871,4,2005-12-30
1905581,5,2005-08-16
2508819,3,2004-05-18
1578279,1,2005-05-19
1159695,4,2005-02-15
2588432,3,2005-03-31
2423091,3,2005-09-12
470232,4,2004-04-08
2148699,2,2004-06-05
1342007,3,2004-07-16
466135,4,2004-07-13
2472440,3,2005-08-13
1283744,3,2004-04-17
1927580,4,2004-11-08
716874,5,2005-05-06
4326,4,2005-10-29
```

2.2 Mapping the real world problem to a Machine Learning Problem

2.2.1 Type of Machine Learning Problem

For a given movie and user we need to predict the rating would be given by him/her to the movie.

The given problem is a Recommendation problem

It can also be seen as a Regression problem

2.2.2 Performance metric

- Mean Absolute Percentage Error: https://en.wikipedia.org/wiki/Mean_absolute_percentage_error
- Root Mean Square Error: https://en.wikipedia.org/wiki/Root-mean-square_deviation

2.2.3 Machine Learning Objective and Constraints

1. Minimize RMSE.
2. Try to provide some interpretability.

In [0]:

```
from google.colab import drive
drive.mount('/content/drive')
```

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).

In [0]:

```
# this is just to know how much time will it take to run this entire ipython notebook
from datetime import datetime
```

```

from datetime import datetime
globalstart = datetime.now()
import pandas as pd
import numpy as np
import matplotlib
matplotlib.use('nbagg')

import matplotlib.pyplot as plt
plt.rcParams.update({'figure.max_open_warning': 0})

import seaborn as sns
sns.set_style('whitegrid')
import os
from scipy import sparse
from scipy.sparse import csr_matrix

from sklearn.decomposition import TruncatedSVD
from sklearn.metrics.pairwise import cosine_similarity
import random

```

3. Exploratory Data Analysis

3.1 Preprocessing

3.1.1 Converting / Merging whole data to required format: u_i, m_j, r_ij

In [0]:

```

start = datetime.now()
if not os.path.isfile('drive/My Drive/data_folder/data.csv'):
    # Create a file 'data.csv' before reading it
    # Read all the files in netflix and store them in one big file('data.csv')
    # We re reading from each of the four files and appendig each rating to a global file
    'train.csv'
    data = open('drive/My Drive/data_folder/data.csv', mode='w')

    row = list()
    files=['drive/My Drive/data_folder/combined_data_1.txt','drive/My
Drive/data_folder/combined_data_2.txt',
          'drive/My Drive/data_folder/combined_data_3.txt', 'drive/My
Drive/data_folder/combined_data_4.txt']
    for file in files:
        print("Reading ratings from {}".format(file))
        with open(file) as f:
            for line in f:
                del row[:] # you don't have to do this.
                line = line.strip()
                if line.endswith(':'):
                    # All below are ratings for this movie, until another movie appears.
                    movie_id = line.replace(':', '')
                else:
                    row = [x for x in line.split(',') ]
                    row.insert(0, movie_id)
                    data.write(','.join(row))
                    data.write('\n')
            print("Done.\n")
    data.close()
print('Time taken :', datetime.now() - start)

```

Time taken : 0:00:00.001663

In [0]:

```

print("creating the dataframe from data.csv file..")
df = pd.read_csv('drive/My Drive/data_folder/data.csv', sep=',',nrows=300000,
                 names=['movie', 'user','rating','date'])
df.date = pd.to_datetime(df.date)

```

```
print('Done.\n')

# we are arranging the ratings according to time.
print('Sorting the dataframe by date..')
df.sort_values(by='date', inplace=True)
print('Done..')
```

creating the dataframe from data.csv file..
Done.

Sorting the dataframe by date..
Done..

In [0]:

```
df.shape
```

Out[0]:

```
(300000, 4)
```

In [0]:

```
df.head()
```

Out[0]:

	movie	user	rating	date
254722	55	1972971	1	1999-12-09
241446	46	510180	3	1999-12-20
290861	77	830363	3	1999-12-21
233922	45	355883	2	1999-12-25
239481	46	1223553	2	1999-12-30

In [0]:

```
df.describe()['rating']
```

Out[0]:

```
count      300000.000000
mean         3.616673
std          1.070734
min          1.000000
25%          3.000000
50%          4.000000
75%          4.000000
max          5.000000
Name: rating, dtype: float64
```

3.1.2 Checking for NaN values

In [0]:

```
# just to make sure that all Nan containing rows are deleted..
print("No of Nan values in our dataframe : ", sum(df.isnull().any()))
```

No of Nan values in our dataframe : 0

3.1.3 Removing Duplicates

In [0]:

```
dup_bool = df.duplicated(['movie','user','rating'])
dups = sum(dup_bool) # by considering all columns..( including timestamp)
print("There are {} duplicate rating entries in the data..".format(dups))
```

There are 0 duplicate rating entries in the data..

3.1.4 Basic Statistics (#Ratings, #Users, and #Movies)

In [0]:

```
print("Total data ")
print("-"*50)
print("\nTotal no of ratings :",df.shape[0])
print("Total No of Users      :", len(np.unique(df.user)))
print("Total No of movies     :", len(np.unique(df.movie)))
```

Total data

```
-----
Total no of ratings : 300000
Total No of Users   : 180721
Total No of movies  : 77
```

3.2 Splitting data into Train and Test(80:20)

In [0]:

```
if not os.path.isfile('drive/My Drive/train.csv'):
    # create the dataframe and store it in the disk for offline purposes..
    df.iloc[:int(df.shape[0]*0.80)].to_csv("drive/My Drive/train.csv", index=False)

if not os.path.isfile('drive/My Drive/test.csv'):
    # create the dataframe and store it in the disk for offline purposes..
    df.iloc[int(df.shape[0]*0.80):].to_csv("drive/My Drive/test.csv", index=False)

train_df = pd.read_csv("drive/My Drive/train.csv", parse_dates=['date'])
test_df = pd.read_csv("drive/My Drive/test.csv")
```

3.2.1 Basic Statistics in Train data (#Ratings, #Users, and #Movies)

In [0]:

```
# movies = train_df.movie.value_counts()
# users = train_df.user.value_counts()
print("Training data ")
print("-"*50)
print("\nTotal no of ratings :",train_df.shape[0])
print("Total No of Users      :", len(np.unique(train_df.user)))
print("Total No of movies     :", len(np.unique(train_df.movie)))
```

Training data

```
-----
Total no of ratings : 240000
Total No of Users   : 148982
Total No of movies  : 77
```

3.2.2 Basic Statistics in Test data (#Ratings, #Users, and #Movies)

In [0]:

```
print("Test data ")
print("-"*50)
```

```
print("\nTotal no of ratings :",test_df.shape[0])
print("Total No of Users      :", len(np.unique(test_df.user)))
print("Total No of movies     :", len(np.unique(test_df.movie)))
```

Test data

```
-----
Total no of ratings : 60000
Total No of Users   : 47088
Total No of movies  : 77
```

3.3 Exploratory Data Analysis on Train data

In [0]:

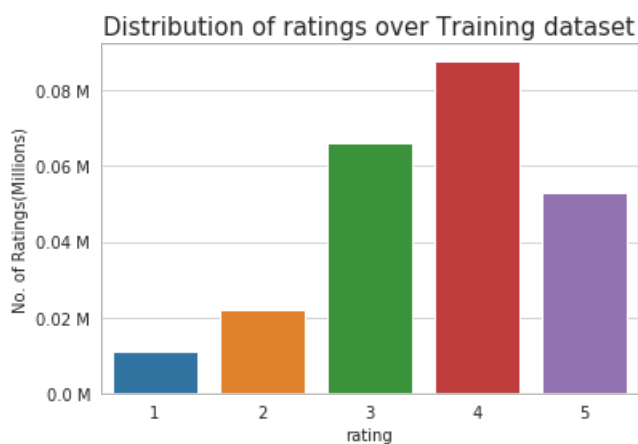
```
# method to make y-axis more readable
def human(num, units = 'M'):
    units = units.lower()
    num = float(num)
    if units == 'k':
        return str(num/10**3) + " K"
    elif units == 'm':
        return str(num/10**6) + " M"
    elif units == 'b':
        return str(num/10**9) + " B"
```

3.3.1 Distribution of ratings

In [0]:

```
%matplotlib inline
fig, ax = plt.subplots()
plt.title('Distribution of ratings over Training dataset', fontsize=15)
sns.countplot(train_df.rating)
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
ax.set_ylabel('No. of Ratings(Millions)')

plt.show()
```



Add new column (week day) to the data set for analysis.

In [0]:

```
# It is used to skip the warning ''SettingWithCopyWarning''..
pd.options.mode.chained_assignment = None # default='warn'

train_df['day_of_week'] = train_df.date.dt.weekday_name

train_df.tail()
```

Out[0]:

	movie	user	rating	date	day_of_week
239995	8	441872	3	2005-07-27	Wednesday
239996	33	178803	5	2005-07-27	Wednesday
239997	8	2551091	2	2005-07-27	Wednesday
239998	28	2436349	2	2005-07-27	Wednesday
239999	68	2088694	5	2005-07-27	Wednesday

3.3.2 Number of Ratings per a month

In [0]:

```
ax = train_df.resample('m', on='date')['rating'].count().plot()
ax.set_title('No of ratings per month (Training data)')
plt.xlabel('Month')
plt.ylabel('No of ratings(per month)')
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
plt.show()
```



3.3.3 Analysis on the Ratings given by user

In [0]:

```
no_of Rated movies per user = train_df.groupby(by='user')['rating'].count().sort_values(ascending=False)
no_of Rated movies per user.head()
```

Out[0]:

```
user
305344    73
2439493    72
387418     69
1461435    41
1932594    38
Name: rating, dtype: int64
```

In [0]:

```
import warnings
warnings.filterwarnings("ignore")

fig = plt.figure(figsize=plt.figaspect(.5))
```



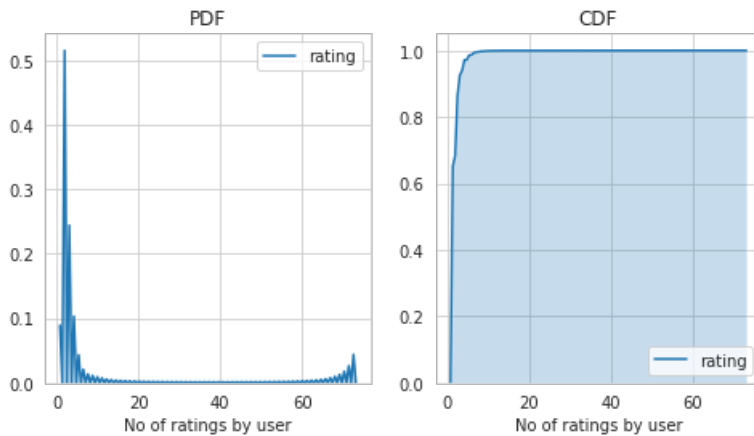
```

ax1 = plt.subplot(121)
sns.kdeplot(no_of Rated movies per user, shade=True, ax=ax1)
plt.xlabel('No of ratings by user')
plt.title("PDF")

ax2 = plt.subplot(122)
sns.kdeplot(no_of Rated movies per user, shade=True, cumulative=True, ax=ax2)
plt.xlabel('No of ratings by user')
plt.title('CDF')

plt.show()

```



In [0]:

```
no_of Rated movies per user.describe()
```

Out[0]:

```

count      148982.000000
mean         1.610933
std          1.219711
min          1.000000
25%          1.000000
50%          1.000000
75%          2.000000
max          73.000000
Name: rating, dtype: float64

```

There, is something interesting going on with the quantiles..

In [0]:

```
quantiles = no_of Rated movies per user.quantile(np.arange(0,1.01,0.01), interpolation='higher')
```

In [0]:

```

plt.title("Quantiles and their Values")
quantiles.plot()
# quantiles with 0.05 difference
plt.scatter(x=quantiles.index[::5], y=quantiles.values[::5], c='orange', label="quantiles with 0.05 intervals")
# quantiles with 0.25 difference
plt.scatter(x=quantiles.index[::25], y=quantiles.values[::25], c='m', label = "quantiles with 0.25 intervals")
plt.ylabel('No of ratings by user')
plt.xlabel('Value at the quantile')
plt.legend(loc='best')

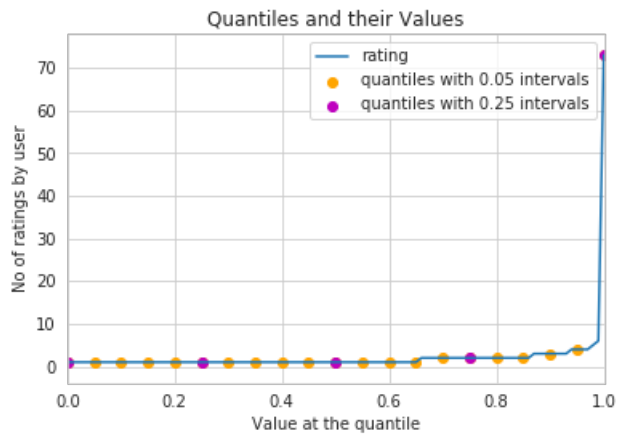
# annotate the 25th, 50th, 75th and 100th percentile values....
for x,y in zip(quantiles.index[::25], quantiles[::25]):
    plt.annotate(s="({} , {})".format(x,y), xy=(x,y), xytext=(x-0.05, y+500)
                , fontweight='bold')

```

```
plt.show()
```

(1.0 , 73)

(0.0 , 1) (0.25 , 1) (0.5 , 1) (0.75 , 2)



In [0]:

```
quantiles[:,5]
```

Out[0]:

```
0.00    1
0.05    1
0.10    1
0.15    1
0.20    1
0.25    1
0.30    1
0.35    1
0.40    1
0.45    1
0.50    1
0.55    1
0.60    1
0.65    1
0.70    2
0.75    2
0.80    2
0.85    2
0.90    3
0.95    4
1.00   73
Name: rating, dtype: int64
```

how many ratings at the last 5% of all ratings??

In [0]:

```
print('\n No of ratings at last 5 percentile : {}'.format(sum(no_of Rated movies per user >= 749)
))
```

```
No of ratings at last 5 percentile : 0
```

3.3.4 Analysis of ratings of a movie given by a user

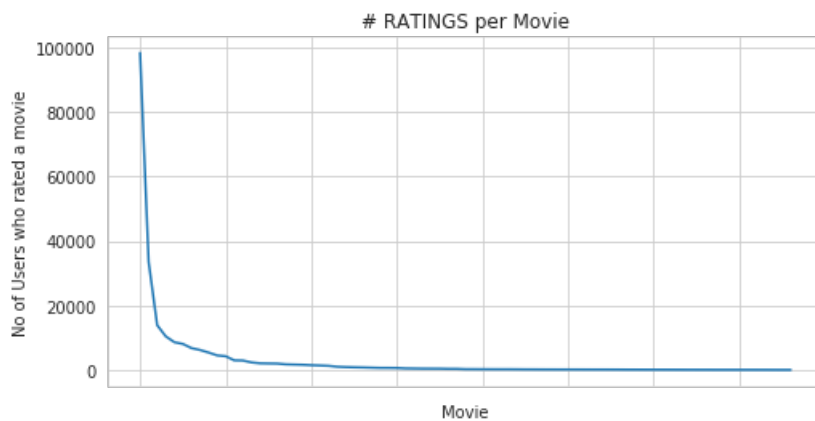
In [0]:

```
no_of_ratings_per_movie = train_df.groupby(by='movie')
['rating'].count().sort_values(ascending=False)

fig = plt.figure(figsize=plt.figaspect(.5))
ax = plt.gca()
plt.plot(no_of_ratings_per_movie.values)
plt.title('# RATINGS per Movie')
```

```
plt.xlabel('Movie')
plt.ylabel('No of Users who rated a movie')
ax.set_xticklabels([])

plt.show()
```

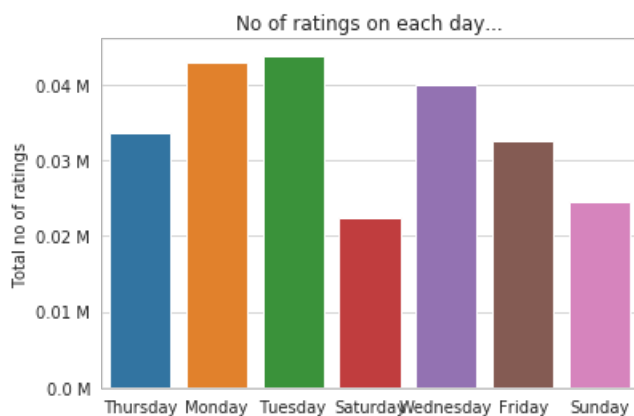


- It is very skewed.. just like number of ratings given per user.
 - There are some movies (which are very popular) which are rated by huge number of users.
 - But most of the movies (like 90%) got some hundreds of ratings.

3.3.5 Number of ratings on each day of the week

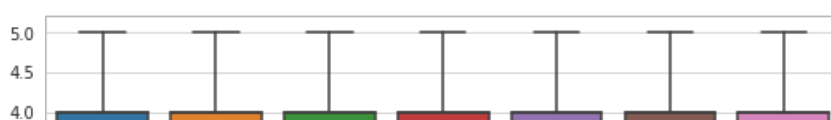
In [0]:

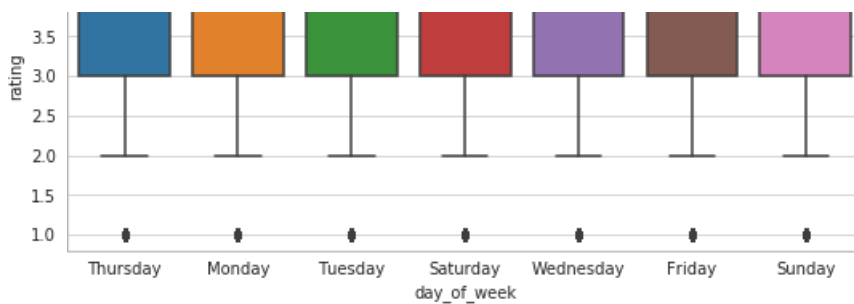
```
fig, ax = plt.subplots()
sns.countplot(x='day_of_week', data=train_df, ax=ax)
plt.title('No of ratings on each day...')
plt.ylabel('Total no of ratings')
plt.xlabel('')
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
plt.show()
```



In [0]:

```
start = datetime.now()
fig = plt.figure(figsize=plt.figaspect(.45))
sns.boxplot(y='rating', x='day_of_week', data=train_df)
plt.show()
print(datetime.now() - start)
```





0:00:00.457806

In [0]:

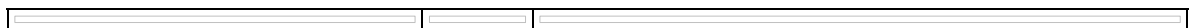
```
avg_week_df = train_df.groupby(by=['day_of_week'])['rating'].mean()
print(" AVerage ratings")
print("-"*30)
print(avg_week_df)
print("\n")
```

```

Average ratings
-----
day_of_week
Friday      3.603877
Monday      3.623158
Saturday    3.619205
Sunday      3.619197
Thursday    3.630634
Tuesday     3.636274
Wednesday   3.617269
Name: rating, dtype: float64

```

3.3.6 Creating sparse matrix from data frame



3.3.6.1 Creating sparse matrix from train data frame

In [0]:

```

start = datetime.now()
if os.path.isfile('drive/My Drive/train_sparse_matrix.npz'):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    train_sparse_matrix = sparse.load_npz('drive/My Drive/train_sparse_matrix.npz')
    print("DONE..")
else:
    print("We are creating sparse_matrix from the dataframe..")
    # create sparse_matrix and store it for after usage.
    # csr_matrix(data_values, (row_index, col_index), shape_of_matrix)
    # It should be in such a way that, MATRIX[row, col] = data
    train_sparse_matrix = sparse.csr_matrix((train_df.rating.values, (train_df.user.values,
                                                                    train_df.movie.values)),)

    print('Done. It\'s shape is : (user, movie) : ', train_sparse_matrix.shape)
    print('Saving it into disk for furthur usage..')
    # save it into disk
    sparse.save_npz("drive/My Drive/train_sparse_matrix.npz", train_sparse_matrix)
    print('Done..\n')

print(datetime.now() - start)

```

```

It is present in your pwd, getting it from disk....
DONE..
0:00:00.477774

```

The Sparsity of Train Sparse Matrix

In [0]:

```
us,mv = train_sparse_matrix.shape
elem = train_sparse_matrix.count_nonzero()

print("Sparsity Of Train matrix : {} % ".format( (1-(elem/(us*mv))) * 100) )
```

Sparsity Of Train matrix : 99.88386471516804 %

3.3.6.2 Creating sparse matrix from test data frame

In [0]:

```
start = datetime.now()
if os.path.isfile('drive/My Drive/test_sparse_matrix.npz'):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    test_sparse_matrix = sparse.load_npz('drive/My Drive/test_sparse_matrix.npz')
    print("DONE..")
else:
    print("We are creating sparse_matrix from the dataframe..")
    # create sparse_matrix and store it for after usage.
    # csr_matrix(data_values, (row_index, col_index), shape_of_matrix)
    # It should be in such a way that, MATRIX[row, col] = data
    test_sparse_matrix = sparse.csr_matrix((test_df.rating.values, (test_df.user.values,
                                                                    test_df.movie.values)))

    print('Done. It\'s shape is : (user, movie) : ',test_sparse_matrix.shape)
    print('Saving it into disk for furthur usage..')
    # save it into disk
    sparse.save_npz("drive/My Drive/test_sparse_matrix.npz", test_sparse_matrix)
    print('Done..\n')

print(datetime.now() - start)
```

It is present in your pwd, getting it from disk....

DONE..

0:00:00.407613

The Sparsity of Test data Matrix

In [0]:

```
us,mv = test_sparse_matrix.shape
elem = test_sparse_matrix.count_nonzero()

print("Sparsity Of Test matrix : {} % ".format( (1-(elem/(us*mv))) * 100) )
```

Sparsity Of Test matrix : 99.97096614591644 %

3.3.7 Finding Global average of all movie ratings, Average rating per user, and Average rating per movie

In [0]:

```
# get the user averages in dictionary (key: user_id/movie_id, value: avg rating)

def get_average_ratings(sparse_matrix, of_users):

    # average ratings of user/axes
    ax = 1 if of_users else 0 # 1 - User axes, 0 - Movie axes

    # ".A1" is for converting Column_Matrix to 1-D numpy array
    sum_of_ratings = sparse_matrix.sum(axis=ax).A1
```

```

# Boolean matrix of ratings ( whether a user rated that movie or not)
isRated = sparse_matrix!=0
# no of ratings that each user OR movie..
no_of_ratings = isRated.sum(axis=ax).A1

# max_user and max_movie ids in sparse matrix
u,m = sparse_matrix.shape
# create a dictionary of users and their average ratings..
average_ratings = { i : sum_of_ratings[i]/no_of_ratings[i]
                    for i in range(u if of_users else m)
                    if no_of_ratings[i] !=0}

# return that dictionary of average ratings
return average_ratings

```

3.3.7.1 finding global average of all movie ratings

In [0]:

```

train_averages = dict()
# get the global average of ratings in our train set.
train_global_average = train_sparse_matrix.sum()/train_sparse_matrix.count_nonzero()
train_averages['global'] = train_global_average
train_averages

```

Out[0]:

```
{'global': 3.622225}
```

In [0]:

```
train_averages
```

Out[0]:

```
{'global': 3.622225}
```

3.3.7.2 finding average rating per user

In [0]:

```

train_averages['user'] = get_average_ratings(train_sparse_matrix, of_users=True)
print('\nAverage rating of user 168 :',train_averages['user'][168])

```

Average rating of user 168 : 4.0

3.3.7.3 finding average rating per movie

In [0]:

```

train_averages['movie'] = get_average_ratings(train_sparse_matrix, of_users=False)
print('\n AVerage rating of movie 15 :',train_averages['movie'][15])

```

AVerage rating of movie 15 : 3.311284046692607

3.3.7.4 PDF's & CDF's of Avg.Ratings of Users & Movies (In Train Data)

In [0]:

```

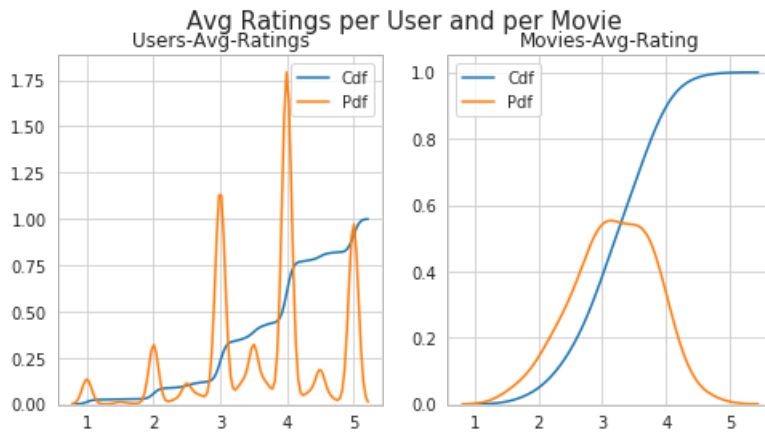
start = datetime.now()
# draw pdfs for average rating per user and average
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))
fig.suptitle('Avg Ratings per User and per Movie', fontsize=15)
ax1.set title('Users-Avg-Ratings')

```

```
# get the list of average user ratings from the averages dictionary..
user_averages = [rat for rat in train_averages['user'].values()]
sns.distplot(user_averages, ax=ax1, hist=False,
              kde_kws=dict(cumulative=True), label='Cdf')
sns.distplot(user_averages, ax=ax1, hist=False, label='Pdf')

ax2.set_title('Movies-Avg-Rating')
# get the list of movie average ratings from the dictionary..
movie_averages = [rat for rat in train_averages['movie'].values()]
sns.distplot(movie_averages, ax=ax2, hist=False,
              kde_kws=dict(cumulative=True), label='Cdf')
sns.distplot(movie_averages, ax=ax2, hist=False, label='Pdf')

plt.show()
print(datetime.now() - start)
```



0:00:12.545016

3.3.8 Cold Start problem

3.3.8.1 Cold Start problem with Users

In [0]:

```
total_users = len(np.unique(df.user))
users_train = len(train_averages['user'])
new_users = total_users - users_train

print('\nTotal number of Users :', total_users)
print('\nNumber of Users in Train data :', users_train)
print("\nNo of Users that didn't appear in train data: {} ({} %) \n ".format(new_users,
np.round((new_users/total_users)*100, 2)))
```

Total number of Users : 180721

Number of Users in Train data : 148982

No of Users that didn't appear in train data: 31739 (17.56 %)

We might have to handle **new users (75148)** who didn't appear in train data.

3.3.8.2 Cold Start problem with Movies

In [0]:

```
total_movies = len(np.unique(df.movie))
movies_train = len(train_averages['movie'])
```



```

movies_train = train_data[movies_train > 0]
new_movies = total_movies - movies_train

print('\nTotal number of Movies :', total_movies)
print('\nNumber of Users in Train data :', movies_train)
print("\nNo of Movies that didn't appear in train data: {} ({} %) \n".format(new_movies,
np.round((new_movies/total_movies)*100, 2)))

```

Total number of Movies : 77

Number of Users in Train data : 77

No of Movies that didn't appear in train data: 0 (0.0 %)

We might have to handle **346 movies** (small comparatively) in test data

3.4 Computing Similarity matrices

3.4.1 Computing User-User Similarity matrix

1. Calculating User User Similarity_Matrix is **not very easy** (unless you have huge Computing Power and lots of time) because of number of. users being large.

- You can try if you want to. Your system could crash or the program stops with **Memory Error**

3.4.1.1 Trying with all dimensions (17k dimensions per user)

In [0]:

```

from sklearn.metrics.pairwise import cosine_similarity

def compute_user_similarity(sparse_matrix, compute_for_few=False, top = 100, verbose=False, verb_for_n_rows = 20,
                           draw_time_taken=True):
    no_of_users, _ = sparse_matrix.shape
    # get the indices of non zero rows(users) from our sparse matrix
    row_ind, col_ind = sparse_matrix.nonzero()
    row_ind = sorted(set(row_ind)) # we don't have to
    time_taken = list() # time taken for finding similar users for an user..

    # we create rows, cols, and data lists.., which can be used to create sparse matrices
    rows, cols, data = list(), list(), list()
    if verbose: print("Computing top",top,"similarities for each user..")

    start = datetime.now()
    temp = 0

    for row in row_ind[:top] if compute_for_few else row_ind:
        temp = temp+1
        prev = datetime.now()

        # get the similarity row for this user with all other users
        sim = cosine_similarity(sparse_matrix.getrow(row), sparse_matrix).ravel()
        # We will get only the top 'top' most similar users and ignore rest of them..
        top_sim_ind = sim.argsort()[-top:]
        top_sim_val = sim[top_sim_ind]

        # add them to our rows, cols and data
        rows.extend([row]*top)
        cols.extend(top_sim_ind)
        data.extend(top_sim_val)
        time_taken.append(datetime.now().timestamp() - prev.timestamp())
    if verbose:

```

```

        if temp%verb_for_n_rows == 0:
            print("computing done for {} users [ time elapsed : {} ]".format(temp, datetime.now()-start))

# lets create sparse matrix out of these and return it
if verbose: print('Creating Sparse matrix from the computed similarities')
#return rows, cols, data

if draw_time_taken:
    plt.plot(time_taken, label = 'time taken for each user')
    plt.plot(np.cumsum(time_taken), label='Total time')
    plt.legend(loc='best')
    plt.xlabel('User')
    plt.ylabel('Time (seconds)')
    plt.show()

return sparse.csr_matrix((data, (rows, cols)), shape=(no_of_users, no_of_users)), time_taken

```

In [0]:

```

start = datetime.now()
u_u_sim_sparse, _ = compute_user_similarity(train_sparse_matrix, compute_for_few=True, top = 100,
                                           verbose=True)

print("-"*100)
print("Time taken :",datetime.now()-start)

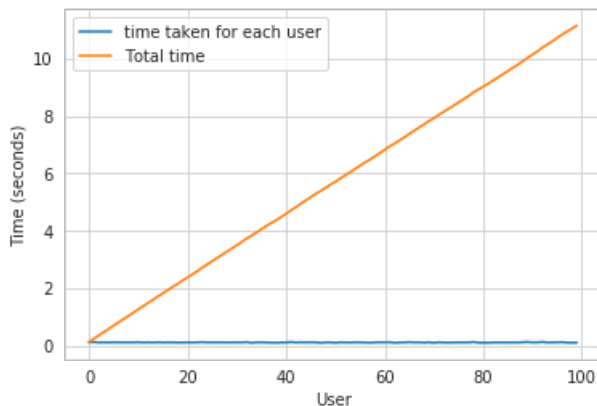
```

Computing top 100 similarities for each user..

```

computing done for 20 users [ time elapsed : 0:00:02.264690 ]
computing done for 40 users [ time elapsed : 0:00:04.485408 ]
computing done for 60 users [ time elapsed : 0:00:06.715259 ]
computing done for 80 users [ time elapsed : 0:00:08.934788 ]
computing done for 100 users [ time elapsed : 0:00:11.164707 ]
Creating Sparse matrix from the computed similarities

```



Time taken : 0:00:11.557404

3.4.1.2 Trying with reduced dimensions (Using TruncatedSVD for dimensionality reduction of user vector)

- We have **405,041 users** in our training set and computing similarities between them..(**17K dimensional vector**..) is time consuming..
- From above plot, It took roughly **8.88 sec** for computing similar users for **one user**
- We have **405,041 users** with us in training set.
- $405041 \times 8.88 = 3596764.08\text{sec} = 59946.068\text{ min}$
 - Even if we run on 4 cores parallelly (a typical system now a days), It will still take almost **10 and 1/2 days**.

IDEA: Instead, we will try to reduce the dimensions using SVD, so that **it might** speed up the process...

In [0]:

```
from datetime import datetime
from sklearn.decomposition import TruncatedSVD

start = datetime.now()

# initialize the algorithm with some parameters..
# All of them are default except n_components. n_iter is for Randomized SVD solver.
netflix_svd = TruncatedSVD(n_components=70, algorithm='randomized', random_state=15)
trunc_svd = netflix_svd.fit_transform(train_sparse_matrix)

print(datetime.now()-start)
```

0:01:47.515787

Here,

- $\sum \rightarrow (\text{netflix_svd.singular_values_})$
- $V^T \rightarrow (\text{netflix_svd.components_})$
- bigr is not returned. instead **Projection_of_X** onto the new vectorspace is returned.
- It uses **randomized svd** internally, which returns **All 3 of them separately**. Use that instead..

In [0]:

```
expl_var = np.cumsum(netflix_svd.explained_variance_ratio_)
```

In [0]:

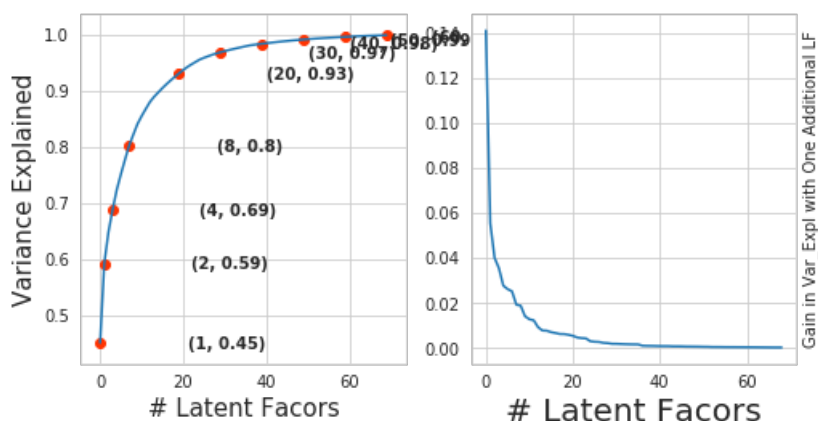
```
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))

ax1.set_ylabel("Variance Explained", fontsize=15)
ax1.set_xlabel("# Latent Facors", fontsize=15)
ax1.plot(expl_var)
# annotate some (latentfactors, expl_var) to make it clear
ind = [1, 2, 4, 8, 20, 30, 40, 50, 60, 70]
ax1.scatter(x = [i-1 for i in ind], y = expl_var[[i-1 for i in ind]], c='#ff3300')
for i in ind:
    ax1.annotate(s="({}, {})".format(i, np.round(expl_var[i-1], 2)), xy=(i-1, expl_var[i-1]),
                xytext = (i+20, expl_var[i-1] - 0.01), fontweight='bold')

change_in_expl_var = [expl_var[i+1] - expl_var[i] for i in range(len(expl_var)-1)]
ax2.plot(change_in_expl_var)

ax2.set_ylabel("Gain in Var_Expl with One Additional LF", fontsize=10)
ax2.yaxis.set_label_position("right")
ax2.set_xlabel("# Latent Facors", fontsize=20)

plt.show()
```



In [0]:

```
for i in ind:
    print("{} {}".format(i, np.round(expl_var[i-1], 2)))
```

```
(1, 0.45)
(2, 0.59)
(4, 0.69)
(8, 0.8)
(20, 0.93)
(30, 0.97)
(40, 0.98)
(50, 0.99)
(60, 1.0)
(70, 1.0)
```

italicized text

I think 70 dimensions is not bad

- By just taking **(20 to 30)** latent factors, explained variance that we could get is **20 %**.
- To take it to **60%**, we have to take **almost 400 latent factors**. It is not fare.
- It basically is the **gain of variance explained**, if we **add one additional latent factor to it**.
- By adding one by one latent factor too it, the **_gain in explained variance** with that addition is decreasing. (Obviously, because they are sorted that way).
- **LHS Graph:**
 - **x** --- (No of latent factors),
 - **y** --- (The variance explained by taking x latent factors)
- **More decrease in the line (RHS graph) :**
 - We are getting more explained variance than before.
- **Less decrease in that line (RHS graph) :**
 - We are not getting benefitted from adding latent factor further. This is what is shown in the plots.
- **RHS Graph:**
 - **x** --- (No of latent factors),
 - **y** --- (Gain n Expl_Var by taking one additional latent factor)

In [0]:

```
# Let's project our Original U_M matrix into into 500 Dimensional space...
start = datetime.now()
trunc_matrix = train_sparse_matrix.dot(netflix_svd.components_.T)
print(datetime.now() - start)
```

0:00:00.450511

In [0]:

```
type(trunc_matrix), trunc_matrix.shape
```

Out[0]:

(numpy.ndarray, (2649430, 70))

- Let's convert this to actual sparse matrix and store it for future purposes

In [0]:

```
if not os.path.isfile('drive/My Drive/trunc_sparse_matrix.npz'):
    # create that sparse matrix
    trunc_sparse_matrix = sparse.csr_matrix(trunc_matrix)
    # Save this truncated sparse matrix for later usage..
    sparse.save_npz('trunc_sparse_matrix', trunc_sparse_matrix)
else:
```

```
trunc_sparse_matrix = sparse.load_npz('drive/My Drive/trunc_sparse_matrix.npz')
```

```
In [0]:
```

```
trunc_sparse_matrix.shape
```

```
Out[0]:
```

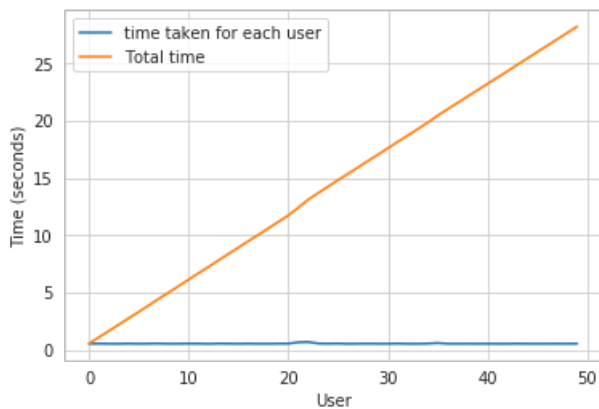
```
(2649430, 70)
```

```
In [0]:
```

```
start = datetime.now()
trunc_u_u_sim_matrix, _ = compute_user_similarity(trunc_sparse_matrix, compute_for_few=True, top=50,
, verbose=True,
                                                    verb_for_n_rows=10)

print("-"*50)
print("time:",datetime.now()-start)
```

```
Computing top 50 similarities for each user..
computing done for 10 users [   time elapsed : 0:00:05.590038   ]
computing done for 20 users [   time elapsed : 0:00:11.197443   ]
computing done for 30 users [   time elapsed : 0:00:17.072612   ]
computing done for 40 users [   time elapsed : 0:00:22.699134   ]
computing done for 50 users [   time elapsed : 0:00:28.262182   ]
Creating Sparse matrix from the computed similarities
```



```
-----
time: 0:00:30.187524
```

: This is taking more time for each user than Original one.

- from above plot, It took almost **12.18** for computing similar users for **one user**
- We have **405041 users** with us in training set.
- { $405041 \times 12.18 = 4933399.38 \text{ sec}$ } $= 82223.323 \text{ min} = 1370.388716667 \text{ hours} = 57.099529861 \text{ days}$...
 - Even we run on 4 cores parallelly (a typical system now a days), It will still take almost **(14 - 15) days**.
- **Why did this happen...??**
 - Just think about it. It's not that difficult.

----- (sparse & dense.....get it ??)-----

Is there any other way to compute user user similarity..??

-An alternative is to compute similar users for a particular user, whenever required (ie., Run time)

- We maintain a binary Vector for users, which tells us whether we already computed or

```

not..
- ***If not*** :
    - Compute top (let's just say, 1000) most similar users for this given user, and add
      this to our datastructure, so that we can just access it(similar users) without recomputing
      it again.
    -
- ***If It is already Computed***:
    - Just get it directly from our datastructure, which has that information.
    - In production time, We might have to recompute similarities, if it is computed a long
      time ago. Because user preferences changes over time. If we could maintain some kind of
      Timer, which when expires, we have to update it ( recompute it ).
    -
- ***Which datastructure to use***:
    - It is purely implementation dependant.
    - One simple method is to maintain a **Dictionary Of Dictionaries**.
      -
      - **key      **: _userid_
      - __value__: _Again a dictionary_
        - __key__ : _Similar User_
        - __value__: _Similarity Value_

```

3.4.2 Computing Movie-Movie Similarity matrix

In [0]:

```

start = datetime.now()
if not os.path.isfile('drive/My Drive/m_m_sim_sparse.npz'):
    print("It seems you don't have that file. Computing movie_movie similarity...")
    start = datetime.now()
    m_m_sim_sparse = cosine_similarity(X=train_sparse_matrix.T, dense_output=False)
    print("Done..")
    # store this sparse matrix in disk before using it. For future purposes.
    print("Saving it to disk without the need of re-computing it again.. ")
    sparse.save_npz("drive/My Drive/m_m_sim_sparse.npz", m_m_sim_sparse)
    print("Done..")
else:
    print("It is there, We will get it.")
    m_m_sim_sparse = sparse.load_npz("drive/My Drive/m_m_sim_sparse.npz")
    print("Done ...")

print("It's a ",m_m_sim_sparse.shape," dimensional matrix")

print(datetime.now() - start)

```

```

It is there, We will get it.
Done ...
It's a (78, 78) dimensional matrix
0:00:00.342763

```

In [0]:

```
m_m_sim_sparse.shape
```

Out[0]:

```
(78, 78)
```

- Even though we have similarity measure of each movie, with all other movies, We generally don't care much about least similar movies.
- Most of the times, only top_xxx similar items matters. It may be 10 or 100.
- We take only those top similar movie ratings and store them in a saperate dictionary.

In [0]:

```
movie_ids = np.unique(m_m_sim_sparse.nonzero()[1])
```

In [0]:

```
start = datetime.now()
similar_movies = dict()
for movie in movie_ids:
    # get the top similar movies and store them in the dictionary
    sim_movies = m_m_sim_sparse[movie].toarray().ravel().argsort()[::-1][1:]
    similar_movies[movie] = sim_movies[:100]
print(datetime.now() - start)

# just testing similar movies for movie_15
similar_movies[15]
```

0:00:00.025473

Out[0]:

```
array([31, 61, 49, 76, 27, 62, 73, 30, 50,  4, 12, 47, 35, 46, 58, 28, 55,
        7, 67, 23, 63, 41, 57, 18, 24, 44, 72, 66, 64, 60, 53,  1, 19, 17,
       45, 51, 77,  9, 29, 10, 16, 21, 36, 25, 48, 20, 75, 69, 13,  5, 26,
       14, 70, 65, 34, 38, 22, 54, 43, 39,  8, 42,  3, 59, 33, 11,  6, 71,
       37, 56,  2, 74, 40, 52, 68, 32,  0])
```

3.4.3 Finding most similar movies using similarity matrix

Does Similarity really works as the way we expected...?

Let's pick some random movie and check for its similar movies....

In [0]:

```
# First Let's load the movie details into soe dataframe..
# movie details are in 'netflix/movie_titles.csv'

movie_titles = pd.read_csv("drive/My Drive/data_folder/movie_titles.csv", sep=',', header = None,
                           names=['movie_id', 'year_of_release', 'title'], verbose=True,
                           index_col = 'movie_id', encoding = "ISO-8859-1")

movie_titles.head()
```

Tokenization took: 4.60 ms

Type conversion took: 10.48 ms

Parser memory cleanup took: 0.01 ms

Out[0]:

	year_of_release	title
movie_id		
1	2003.0	Dinosaur Planet
2	2004.0	Isle of Man TT 2004 Review
3	1997.0	Character
4	1994.0	Paula Abdul's Get Up & Dance
5	2004.0	The Rise and Fall of ECW

Similar Movies for 'Vampire Journals'

In [0]:

```
mv_id = 67

print("\nMovie ----->",movie_titles.loc[mv_id].values[1])

print("\nIt has {} Ratings from users.".format(train_sparse_matrix[:,mv_id].getnnz()))
```

```
print("\nWe have {} movies which are similar to this and we will get only top most..".format(m_m_sim_sparse[:,mv_id].getnnz()))
```

Movie ----> Vampire Journals

It has 267 Ratings from users.

We have 77 movies which are similar to this and we will get only top most..

In [0]:

```
similarities = m_m_sim_sparse[mv_id].toarray().ravel()

similar_indices = similarities.argsort()[::-1][1:]

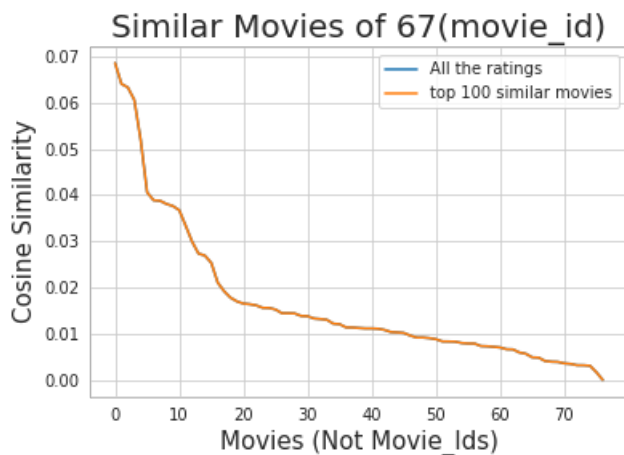
similarities[similar_indices]

sim_indices = similarities.argsort()[::-1][1:] # It will sort and reverse the array and ignore its similarity (ie.,1)

# and return its indices(movie_ids)
```

In [0]:

```
plt.plot(similarities[sim_indices], label='All the ratings')
plt.plot(similarities[sim_indices[:100]], label='top 100 similar movies')
plt.title("Similar Movies of {} (movie_id)".format(mv_id), fontsize=20)
plt.xlabel("Movies (Not Movie_Ids)", fontsize=15)
plt.ylabel("Cosine Similarity", fontsize=15)
plt.legend()
plt.show()
```



Top 10 similar movies

In [0]:

```
movie_titles.loc[sim_indices[:10]]
```

Out [0]:

	year_of_release	title
movie_id		
41	2000.0	Horror Vision
9	1991.0	Class of Nuke 'Em High 2
16	1996.0	Screamers
24	1981.0	My Bloody Valentine

movie_id	year_of_release	title
53	2005.0	The Bonesetter
36	1996.0	Dragonheart
55	1995.0	Jade
48	2001.0	Justice League
7	1992.0	8 Man
77	1995.0	Congo

Similarly, we can **find similar users** and compare how similar they are.

4. Machine Learning Models

In [0]:

```
def get_sample_sparse_matrix(sparse_matrix, no_users, no_movies, path, verbose = True):
    """
        It will get it from the 'path' if it is present or It will create
        and store the sampled sparse matrix in the path specified.
    """

    # get (row, col) and (rating) tuple from sparse_matrix...
    row_ind, col_ind, ratings = sparse.find(sparse_matrix)
    users = np.unique(row_ind)
    movies = np.unique(col_ind)

    print("Original Matrix : (users, movies) -- ({}) {}".format(len(users), len(movies)))
    print("Original Matrix : Ratings -- {}\n".format(len(ratings)))

    # It just to make sure to get same sample everytime we run this program..
    # and pick without replacement....
    np.random.seed(15)
    sample_users = np.random.choice(users, no_users)
    sample_movies = np.random.choice(movies, no_movies)
    # get the boolean mask or these sampled items in originl row/col_inds..
    mask = np.logical_and( np.isin(row_ind, sample_users),
                           np.isin(col_ind, sample_movies) )

    sample_sparse_matrix = sparse.csr_matrix((ratings[mask], (row_ind[mask], col_ind[mask])),
                                              shape=(max(sample_users)+1, max(sample_movies)+1))

    if verbose:
        print("Sampled Matrix : (users, movies) -- ({}) {}".format(len(sample_users), len(sample_movies)))
        print("Sampled Matrix : Ratings --", format(ratings[mask].shape[0]))

    print('Saving it into disk for furthur usage..')
    # save it into disk
    sparse.save_npz(path, sample_sparse_matrix)
    if verbose:
        print('Done..\n')

    return sample_sparse_matrix
```

4.1 Sampling Data

4.1.1 Build sample train data from the train data

In [0]:

```
start = datetime.now()
path = "drive/My Drive/sample_data/sample_train_sparse_matrix.npz"
if os.path.isfile(path):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    sample_train_sparse_matrix = sparse.load_npz(path)
    print("DONE..")
else:
    # get 25k users and 3k movies from available data
    sample_train_sparse_matrix = get_sample_sparse_matrix(train_sparse_matrix, no_users=25000, no_movies=3000,
                                                         path = path)

print(datetime.now() - start)
```

It is present in your pwd, getting it from disk....
DONE..
0:00:00.349492

4.1.2 Build sample test data from the test data

In [0]:

```
start = datetime.now()

path = "drive/My Drive/sample_data/sample_test_sparse_matrix.npz"
if os.path.isfile(path):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    sample_test_sparse_matrix = sparse.load_npz(path)
    print("DONE..")
else:
    # get 5k users and 1000 movies from available data
    sample_test_sparse_matrix = get_sample_sparse_matrix(test_sparse_matrix, no_users=5000, no_movies=1000,
                                                         path = path)

print(datetime.now() - start)
```

It is present in your pwd, getting it from disk....
DONE..
0:00:00.402101

4.2 Finding Global Average of all movie ratings, Average rating per User, and Average rating per Movie (from sampled train)

In [0]:

```
sample_train_averages = dict()
```

In [0]:

```
sample_train_averages
```

Out[0]:

```
{}
```

4.2.1 Finding Global Average of all movie ratings

In [0]:

```
# get the global average of ratings in our train set.
global_average = sample_train_sparse_matrix.sum()/sample_train_sparse_matrix.count_nonzero()
sample_train_averages['global'] = global_average
```

```
sample_train_averages['global'] = global_average
sample_train_averages
```

Out[0]:

```
{'global': 3.6180258664730713}
```

4.2.2 Finding Average rating per User

In [0]:

```
sample_train_averages['user'] = get_average_ratings(sample_train_sparse_matrix, of_users=True)
print('\nAverage rating of user 911 :',sample_train_averages['user'][911])
```

Average rating of user 911 : 4.0

4.2.3 Finding Average rating per Movie

In [0]:

```
sample_train_averages['movie'] = get_average_ratings(sample_train_sparse_matrix, of_users=False)
print('\n Average rating of movie 2 :',sample_train_averages['movie'][2])
```

Average rating of movie 2 : 3.8666666666666667

4.3 Featurizing data

In [0]:

```
print('\n No of ratings in Our Sampled train matrix is : {}'.format(sample_train_sparse_matrix.count_nonzero()))
print('\n No of ratings in Our Sampled test matrix is : {}'.format(sample_test_sparse_matrix.count_nonzero()))
```

No of ratings in Our Sampled train matrix is : 37191

No of ratings in Our Sampled test matrix is : 6022

4.3.1 Featurizing data for regression problem

4.3.1.1 Featurizing train data

In [0]:

```
# get users, movies and ratings from our samples train sparse matrix
sample_train_users, sample_train_movies, sample_train_ratings =
sparse.find(sample_train_sparse_matrix)
```

In [0]:

```
#####
# It took me almost 10 hours to prepare this train dataset. #
#####
start = datetime.now()
if os.path.isfile('drive/My Drive/sample_data/reg_train.csv'):
    print("File already exists you don't have to prepare again..." )
else:
    print('preparing {} tuples for the dataset..{}\n'.format(len(sample_train_ratings)))
    with open('drive/My Drive/sample_data/reg_train.csv', mode='w') as reg_data_file:
        count = 0
        for (user, movie, rating) in zip(sample_train_users, sample_train_movies, sample_train_ratings):
```

```

for (user, movie, rating) in zip(sample_train_users, sample_train_movies,
sample_train_ratings):
    st = datetime.now()
    # print(user, movie)
    #----- Ratings of "movie" by similar users of "user" -----
    --
    # compute the similar Users of the "user"
    user_sim = cosine_similarity(sample_train_sparse_matrix[user],
sample_train_sparse_matrix).ravel()
    top_sim_users = user_sim.argsort()[::-1][1:] # we are ignoring 'The User' from its simi
lar users.
    # get the ratings of most similar users for this movie
    top_ratings = sample_train_sparse_matrix[top_sim_users, movie].toarray().ravel()
    # we will make it's length "5" by adding movie averages to .
    top_sim_users_ratings = list(top_ratings[top_ratings != 0][:5])
    top_sim_users_ratings.extend([sample_train_averages['movie'][movie]]*(5 -
len(top_sim_users_ratings)))
    # print(top_sim_users_ratings, end=" ")

    #----- Ratings by "user" to similar movies of "movie" -----
    ----
    # compute the similar movies of the "movie"
    movie_sim = cosine_similarity(sample_train_sparse_matrix[:,movie].T,
sample_train_sparse_matrix.T).ravel()
    top_sim_movies = movie_sim.argsort()[::-1][1:] # we are ignoring 'The User' from its si
milar users.
    # get the ratings of most similar movie rated by this user..
    top_ratings = sample_train_sparse_matrix[user, top_sim_movies].toarray().ravel()
    # we will make it's length "5" by adding user averages to.
    top_sim_movies_ratings = list(top_ratings[top_ratings != 0][:5])
    top_sim_movies_ratings.extend([sample_train_averages['user']
[user]]*(5-len(top_sim_movies_ratings)))
    # print(top_sim_movies_ratings, end=" : -- ")

    #-----prepare the row to be stores in a file-----#
    row = list()
    row.append(user)
    row.append(movie)
    # Now add the other features to this data...
    row.append(sample_train_averages['global']) # first feature
    # next 5 features are similar_users "movie" ratings
    row.extend(top_sim_users_ratings)
    # next 5 features are "user" ratings for similar_movies
    row.extend(top_sim_movies_ratings)
    # Avg_user rating
    row.append(sample_train_averages['user'][user])
    # Avg_movie rating
    row.append(sample_train_averages['movie'][movie])

    # finalley, The actual Rating of this user-movie pair...
    row.append(rating)
    count = count + 1

    # add rows to the file opened..
    reg_data_file.write(', '.join(map(str, row)))
    reg_data_file.write('\n')
    if (count)%10000 == 0:
        # print(', '.join(map(str, row)))
        print("Done for {} rows----- {}".format(count, datetime.now() - start))

print(datetime.now() - start)

```

File already exists you don't have to prepare again...
0:00:00.001497

Reading from the file to make a Train_dataframe

In [0]:

```

reg_train = pd.read_csv('drive/My Drive/sample_data/reg_train.csv', names = ['user', 'movie', 'GAv
', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5', 'smr1', 'smr2', 'smr3', 'smr4', 'smr5', 'UAvg', 'MAvg', 'r
ating'], header=None)
reg_train.head()

```

Out [0]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating
0	2442	1	3.618026	4	4	4	4	5	3.0	3.0	3.0	3.00	3.00	3.00	3.640625	3
1	31913	1	3.618026	5	4	4	4	5	1.0	2.5	2.5	2.50	2.50	2.50	3.640625	4
2	42930	1	3.618026	4	4	3	2	3	3.0	4.0	5.0	3.00	3.60	3.60	3.640625	3
3	94565	1	3.618026	5	5	5	3	5	3.0	3.0	3.0	3.25	3.25	3.25	3.640625	4
4	145873	1	3.618026	5	5	4	2	3	5.0	4.0	4.0	4.00	4.00	4.00	3.640625	3

- **GAvg** : Average rating of all the ratings
- **Similar users rating of this movie:**
 - sur1, sur2, sur3, sur4, sur5 (top 5 similar users who rated that movie..)
- **Similar movies rated by this user:**
 - smr1, smr2, smr3, smr4, smr5 (top 5 similar movies rated by this movie..)
- **UAvg** : User's Average rating
- **MAvg** : Average rating of this movie
- **rating** : Rating of this movie by this user.

4.3.1.2 Featurizing test data

In [0]:

```
# get users, movies and ratings from the Sampled Test
sample_test_users, sample_test_movies, sample_test_ratings = sparse.find(sample_test_sparse_matrix
)
```

In [0]:

```
sample_train_averages['global']
```

Out[0]:

3.6180258664730713

In [0]:

```
start = datetime.now()

if os.path.isfile('drive/My Drive/sample_data/sample_test_sparse_matrix.npz'):
    print("It is already created...")
else:
    print('preparing {} tuples for the dataset...\n'.format(len(sample_test_ratings)))
    with open('drive/My Drive/sample_data/sample_test_sparse_matrix.npz', mode='w') as
reg_data_file:
        count = 0
        for (user, movie, rating) in zip(sample_test_users, sample_test_movies,
sample_test_ratings):
            st = datetime.now()

            #----- Ratings of "movie" by similar users of "user" -----
            #print(user, movie)
            try:
                # compute the similar Users of the "user"
                user_sim = cosine_similarity(sample_train_sparse_matrix[user]
```

```

        user_sim = cosine_similarity(sample_train_sparse_matrix[user],
sample_train_sparse_matrix).ravel()
        top_sim_users = user_sim.argsort()[::-1][1:] # we are ignoring 'The User' from its
similar users.
        # get the ratings of most similar users for this movie
        top_ratings = sample_train_sparse_matrix[top_sim_users, movie].toarray().ravel()
        # we will make it's length "5" by adding movie averages to .
        top_sim_users_ratings = list(top_ratings[top_ratings != 0][:5])
        top_sim_users_ratings.extend([sample_train_averages['movie'][movie]]*(5 -
len(top_sim_users_ratings)))
        # print(top_sim_users_ratings, end="--")

    except (IndexError, KeyError):
        # It is a new User or new Movie or there are no ratings for given user for top simi
lar movies...
        ##### Cold Start Problem #####
        top_sim_users_ratings.extend([sample_train_averages['global']]*(5 -
len(top_sim_users_ratings)))
        #print(top_sim_users_ratings)
    except:
        print(user, movie)
        # we just want KeyErrors to be resolved. Not every Exception...
        raise

#----- Ratings by "user" to similar movies of "movie" -----#
----
    try:
        # compute the similar movies of the "movie"
        movie_sim = cosine_similarity(sample_train_sparse_matrix[:,movie].T,
sample_train_sparse_matrix.T).ravel()
        top_sim_movies = movie_sim.argsort()[::-1][1:] # we are ignoring 'The User' from it
s similar users.
        # get the ratings of most similar movie rated by this user..
        top_ratings = sample_train_sparse_matrix[user, top_sim_movies].toarray().ravel()
        # we will make it's length "5" by adding user averages to.
        top_sim_movies_ratings = list(top_ratings[top_ratings != 0][:5])
        top_sim_movies_ratings.extend([sample_train_averages['user']
[user]]*(5-len(top_sim_movies_ratings)))
        #print(top_sim_movies_ratings)
    except (IndexError, KeyError):
        #print(top_sim_movies_ratings, end=" : -- ")

top_sim_movies_ratings.extend([sample_train_averages['global']]*(5-len(top_sim_movies_ratings)))
        #print(top_sim_movies_ratings)
    except :
        raise

#-----prepare the row to be stores in a file-----#
row = list()
# add usser and movie name first
row.append(user)
row.append(movie)
row.append(sample_train_averages['global']) # first feature
#print(row)
# next 5 features are similar_users "movie" ratings
row.extend(top_sim_users_ratings)
#print(row)
# next 5 features are "user" ratings for similar_movies
row.extend(top_sim_movies_ratings)
#print(row)
# Avg_user rating
try:
    row.append(sample_train_averages['user'][user])
except KeyError:
    row.append(sample_train_averages['global'])
except:
    raise
#print(row)
# Avg_movie rating
try:
    row.append(sample_train_averages['movie'][movie])
except KeyError:
    row.append(sample_train_averages['global'])
except:
    raise
#print(row)

```

```

# print(row)
# finalley, The actual Rating of this user-movie pair...
row.append(rating)
#print(row)
count = count + 1

# add rows to the file opened..
reg_data_file.write(','.join(map(str, row)))
#print(','.join(map(str, row)))
reg_data_file.write('\n')
if (count)%1000 == 0:
    #print(','.join(map(str, row)))
    print("Done for {} rows----- {}".format(count, datetime.now() - start))
print("",datetime.now() - start)

```

It is already created...

Reading from the file to make a test dataframe

In [0]:

```

reg_test_df = pd.read_csv('drive/My Drive/sample_data/reg_test.csv', names = ['user', 'movie', 'GAvg',
                                     'g', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
                                     'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
                                     'UAvg', 'MAvg', 'rating'], header=None)

reg_test_df.head(4)

```

Out[0]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MA
0	206115	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
1	314933	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
2	353369	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
3	389872	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406

In [0]:

```

reg_train=reg_train[:18000]
reg_test_df=reg_test_df[:5000]
print(reg_train.shape)
print(reg_test_df.shape)

```

```

(18000, 16)
(5000, 16)

```

- **GAvg** : Average rating of all the ratings
- **Similar users rating of this movie:**
 - sur1, sur2, sur3, sur4, sur5 (top 5 simiular users who rated that movie..)
- **Similar movies rated by this user:**
 - smr1, smr2, smr3, smr4, smr5 (top 5 simiular movies rated by this movie..)
- **UAvg** : User AVerage rating
- **MAvg** : Average rating of this movie
- **rating** : Rating of this movie by this user.

4.3.2 Transforming data for Surprise models

In [0]:

```
!pip install scikit-surprise
from surprise import Reader, Dataset
```

```
Collecting scikit-surprise
  Downloading
https://files.pythonhosted.org/packages/4d/fc/cd4210b247d1dca421c25994740cbbf03c5e980e31881f10eaddf
ab0/scikit-surprise-1.0.6.tar.gz (3.3MB)
  |████████████████████████████████████████| 3.3MB 2.8MB/s
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.6/dist-packages (from
scikit-surprise) (0.13.2)
Requirement already satisfied: numpy>=1.11.2 in /usr/local/lib/python3.6/dist-packages (from
scikit-surprise) (1.16.4)
Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.6/dist-packages (from
scikit-surprise) (1.3.0)
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.6/dist-packages (from scikit-
surprise) (1.12.0)
Building wheels for collected packages: scikit-surprise
  Building wheel for scikit-surprise (setup.py) ... done
  Stored in directory:
/root/.cache/pip/wheels/ec/c0/55/3a28eab06b53c220015063ebdb81213cd3dcbb72c088251ec
Successfully built scikit-surprise
Installing collected packages: scikit-surprise
Successfully installed scikit-surprise-1.0.6
```

4.3.2.1 Transforming train data

- We can't give raw data (movie, user, rating) to train the model in Surprise library.
- They have a separate format for TRAIN and TEST data, which will be useful for training the models like SVD, KNNBaseLineOnly....etc., in Surprise.
- We can form the trainset from a file, or from a Pandas DataFrame.
http://surprise.readthedocs.io/en/stable/getting_started.html#load-dom-dataframe-py

In [0]:

```
# It is to specify how to read the dataframe.
# for our dataframe, we don't have to specify anything extra..
reader = Reader(rating_scale=(1,5))

# create the traindata from the dataframe...
train_data = Dataset.load_from_df(reg_train[['user', 'movie', 'rating']], reader)

# build the trainset from traindata.. It is of dataset format from surprise library..
trainset = train_data.build_full_trainset()
```

4.3.2.2 Transforming test data

- Testset is just a list of (user, movie, rating) tuples. (Order in the tuple is important)

In [0]:

```
testset = list(zip(reg_test_df.user.values, reg_test_df.movie.values, reg_test_df.rating.values))
testset[:3]
```

Out[0]:

```
[(206115, 1, 4), (314933, 1, 3), (353369, 1, 3)]
```

4.4 Applying Machine Learning models

- Global dictionary that stores rmse and mape for all the models....

- It stores the metrics in a dictionary of dictionaries

keys : model names(string)

value: dict(**key** : metric, **value** : value)

In [0]:

```
models_evaluation_train = dict()
models_evaluation_test = dict()

models_evaluation_train, models_evaluation_test
```

Out[0]:

```
({}, {})
```

Utility functions for running regression models

In [0]:

```
# to get rmse and mape given actual and predicted ratings..
def get_error_metrics(y_true, y_pred):
    rmse = np.sqrt(np.mean([ (y_true[i] - y_pred[i])**2 for i in range(len(y_pred)) ]))
    mape = np.mean(np.abs( (y_true - y_pred)/y_true )) * 100
    return rmse, mape

#####
#####

def run_xgboost(algo, x_train, y_train, x_test, y_test, verbose=True):
    """
    It will return train_results and test_results
    """

    # dictionaries for storing train and test results
    train_results = dict()
    test_results = dict()

    # fit the model
    print('Training the model..')
    start =datetime.now()
    algo.fit(x_train, y_train, eval_metric = 'rmse')
    print('Done. Time taken : {} \n'.format(datetime.now()-start))
    print('Done \n')

    # from the trained model, get the predictions....
    print('Evaluating the model with TRAIN data...')
    start =datetime.now()
    y_train_pred = algo.predict(x_train)
    # get the rmse and mape of train data...
    rmse_train, mape_train = get_error_metrics(y_train.values, y_train_pred)

    # store the results in train_results dictionary..
    train_results = {'rmse': rmse_train,
                    'mape' : mape_train,
                    'predictions' : y_train_pred}

    #####
    # get the test data predictions and compute rmse and mape
    print('Evaluating Test data')
    y_test_pred = algo.predict(x_test)
    rmse_test, mape_test = get_error_metrics(y_true=y_test.values, y_pred=y_test_pred)
    # store them in our test results dictionary.
    test_results = {'rmse': rmse_test,
                   'mape' : mape_test,
                   'predictions':y_test_pred}

    if verbose:
        print('\nTEST DATA')
        print('-'*30)
        print('RMSE : ', rmse_test)
```

```

print('RMSE : ', rmse_test)
print('MAPE : ', mape_test)

# return these train and test results...
return train_results, test_results

```

Utility functions for Surprise modes

In [0]:

```

# it is just to make sure that all of our algorithms should produce same results
# everytime they run...

my_seed = 15
random.seed(my_seed)
np.random.seed(my_seed)

#####
# get (actual_list , predicted_list) ratings given list
# of predictions (prediction is a class in Surprise).
#####
def get_ratings(predictions):
    actual = np.array([pred.r_ui for pred in predictions])
    pred = np.array([pred.est for pred in predictions])

    return actual, pred

#####
# get 'rmse' and 'mape' , given list of prediction objects
#####
def get_errors(predictions, print_them=False):

    actual, pred = get_ratings(predictions)
    rmse = np.sqrt(np.mean((pred - actual)**2))
    mape = np.mean(np.abs(pred - actual)/actual)

    return rmse, mape*100

#####
# It will return predicted ratings, rmse and mape of both train and test data #
#####
def run_surprise(algo, trainset, testset, verbose=True):
    """
        return train_dict, test_dict

        It returns two dictionaries, one for train and the other is for test
        Each of them have 3 key-value pairs, which specify 'rmse', 'mape', and 'predicted ratings'.
    """
    start = datetime.now()
    # dictionaries that stores metrics for train and test..
    train = dict()
    test = dict()

    # train the algorithm with the trainset
    st = datetime.now()
    print('Training the model...')
    algo.fit(trainset)
    print('Done. time taken : {} \n'.format(datetime.now()-st))

    # ----- Evaluating train data-----#
    st = datetime.now()
    print('Evaluating the model with train data..')
    # get the train predictions (list of prediction class inside Surprise)
    train_preds = algo.test(trainset.build_testset())
    # get predicted ratings from the train predictions..
    train_actual_ratings, train_pred_ratings = get_ratings(train_preds)
    # get 'rmse' and 'mape' from the train predictions.
    train_rmse, train_mape = get_errors(train_preds)
    print('time taken : {}'.format(datetime.now()-st))

    if verbose:

```

```

print('-'*15)
print('Train Data')
print('-'*15)
print("RMSE : {}\n\nMAPE : {}\n".format(train_rmse, train_mape))

#store them in the train dictionary
if verbose:
    print('adding train results in the dictionary..')
train['rmse'] = train_rmse
train['mape'] = train_mape
train['predictions'] = train_pred_ratings

#----- Evaluating Test data-----#
st = datetime.now()
print('\nEvaluating for test data...')
# get the predictions( list of prediction classes) of test data
test_preds = algo.test(testset)
# get the predicted ratings from the list of predictions
test_actual_ratings, test_pred_ratings = get_ratings(test_preds)
# get error metrics from the predicted and actual ratings
test_rmse, test_mape = get_errors(test_preds)
print('time taken : {}'.format(datetime.now()-st))

if verbose:
    print('-'*15)
    print('Test Data')
    print('-'*15)
    print("RMSE : {}\n\nMAPE : {}\n".format(test_rmse, test_mape))
# store them in test dictionary
if verbose:
    print('storing the test results in test dictionary..')
test['rmse'] = test_rmse
test['mape'] = test_mape
test['predictions'] = test_pred_ratings

print('\n'+ '-'*45)
print('Total time taken to run this algorithm :', datetime.now() - start)

# return two dictionaries train and test
return train, test

```

4.4.1 XGBoost with initial 13 features

In [0]:

```
import xgboost as xgb
```

In [0]:

```

#splitting to x and y for computing hyperparameter tuning
x_train_data=reg_train.drop(['user', 'movie','rating'], axis=1)
y_train_data=reg_train['rating']

x_test_data=reg_test_df.drop(['user', 'movie','rating'], axis=1)
y_test_data=reg_test_df['rating']

```

In [0]:

```

def rmse(pred,actual):
    return np.sqrt(np.mean((pred - actual)**2))

```

In [185]:

```

import warnings
warnings.filterwarnings("ignore", category=DeprecationWarning)

train_rmse = []
test_rmse = []
max_depth = [1, 5, 10, 50, 100, 500,700,800,900 ,1000]
n_models =[10,20,30,40,50,60,70,80,90,100]

for i in max_depth:

```

```

xgb_model = xgb.XGBRegressor(max_depth=i,n_jobs=-1)
xgb_model.fit(x_train_data, y_train_data)
y_train_pred = xgb_model.predict(x_train_data)
y_test_pred = xgb_model.predict(x_test_data)
train_rmse_score=rmse(y_train_pred,y_train_data)
train_rmse.append(train_rmse_score)
test_rmse_score=rmse(y_test_pred ,y_test_data)
test_rmse.append(test_rmse_score)
print("Best_depth = ",i ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score\t:",
train_rmse_score)

best_depth = np.argmin(test_rmse)
print("best_depth :",max_depth[best_depth])

```

```

[12:30:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1 test_rmse_score : 1.118686519152436 train_rmse_score : 0.4742197391807634
[12:30:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 5 test_rmse_score : 1.1165563471004316 train_rmse_score : 0.20185295107101406
[12:30:59] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 10 test_rmse_score : 1.1853701654112845 train_rmse_score : 0.029068982365314522
[12:31:05] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 50 test_rmse_score : 1.2323091642703456 train_rmse_score : 0.0012429298136426526
[12:31:20] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 100 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:31:35] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 500 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:31:50] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 700 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:32:06] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 800 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:32:22] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 900 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:32:38] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1000 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
best_depth : 5

```

In [186]:

```

s_train_rmse=[]
s_test_rmse=[]

for j in n_models:
    xgb_model = xgb.XGBRegressor(max_depth=max_depth[best_depth],n_estimators=j,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    s_train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred ,y_test_data)
    s_test_rmse.append(test_rmse_score)
    print("best_estimators = ",j ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score
\t:",train_rmse_score)

best_estimators = np.argmin(s_test_rmse)
print("best_depth :",max_depth[best_depth])
print("best_estimators :",n_models[best_estimators])

```

```

[12:32:54] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror

```

```

n favor of reg:squarederror.
best_estimators = 10    test_rmse_score : 1.4804614168768708    train_rmse_score :
1.2152385499466214
[12:32:54] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 20    test_rmse_score : 1.1699152724516488    train_rmse_score :
0.5403831999609053
[12:32:55] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 30    test_rmse_score : 1.1209988437594627    train_rmse_score :
0.3543631578919108
[12:32:56] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 40    test_rmse_score : 1.114834216555102    train_rmse_score :
0.3015962804147239
[12:32:56] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 50    test_rmse_score : 1.1136777693704727    train_rmse_score :
0.27549062766421706
[12:32:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 60    test_rmse_score : 1.1139789327483343    train_rmse_score :
0.25435819750785754
[12:32:59] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 70    test_rmse_score : 1.1144393760518252    train_rmse_score :
0.2387096275277416
[12:33:00] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 80    test_rmse_score : 1.1150160301961656    train_rmse_score :
0.2252943871448911
[12:33:02] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 90    test_rmse_score : 1.1160731560467065    train_rmse_score :
0.21292301246595102
[12:33:04] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 100   test_rmse_score : 1.1165563471004316    train_rmse_score :
0.20185295107101406
best_depth : 5
best_estimators : 50

```

In [187]:

```

# prepare Train data
x_train = reg_train.drop(['user','movie','rating'], axis=1)
y_train = reg_train['rating']

# Prepare Test data
x_test = reg_test_df.drop(['user','movie','rating'], axis=1)
y_test = reg_test_df['rating']

# initialize Our first XGBoost model...
first_xgb = xgb.XGBRegressor(silent=1, n_jobs=13, random_state=15, n_estimators=n_models[best_estim
ators],max_depth=max_depth[best_depth])
train_results, test_results = run_xgboost(first_xgb, x_train, y_train, x_test, y_test)

# store the results in models_evaluations dictionaries
models_evaluation_train['first_algo'] = train_results
models_evaluation_test['first_algo'] = test_results

xgb.plot_importance(first_xgb)
plt.show()

```

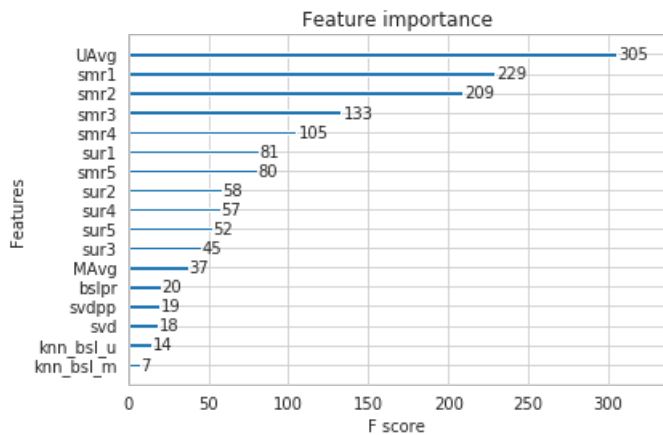
Training the model..
Done. Time taken : 0:00:01.120957

Done

Evaluating the model with TRAIN data...
Evaluating Test data

TEST DATA

RMSE : 1.1136777693704727
MAPE : 35.85766377810637



4.4.2 Surprise BaselineModel

In [0]:

```
from surprise import BaselineOnly
```

Predicted_rating : (baseline prediction)

-

http://surprise.readthedocs.io/en/stable/basic_algorithms.html#surprise.prediction_algorithmr seline_only.BaselineOnly

$$\widehat{r}_{ui} = b_{ui} = \mu + b_u + b_i$$

- μ : Average of all trainings in training data.
- b_u : User bias
- b_i : Item bias (movie biases)

Optimization function (Least Squares Problem)

- http://surprise.readthedocs.io/en/stable/prediction_algorithms.html#baselines-estimates-configuration

$$\sum_{r_{ui} \in R_{train}} \left(r_{ui} - (\mu + b_u + b_i) \right)^2 + \lambda (b_u^2 + b_i^2)$$

minimize $\{b_u, b_i\}$

In [189]:

```
# options are to specify..., how to compute those user and item biases
bsl_options = {'method': 'sgd',
               'learning_rate': .001
              }
bsl_algo = BaselineOnly(bsl_options=bsl_options)
# run this algorithm..., It will return the train and test results..
bsl_train_results, bsl_test_results = run_surprise(bsl_algo, trainset, testset, verbose=True)

# Just store these error metrics in our models_evaluation datastructure
models_evaluation_train['bsl_algo'] = bsl_train_results
models_evaluation_test['bsl_algo'] = bsl_test_results
```

Training the model...
Estimating biases using sgd...

```

Done. time taken : 0:00:00.095034

Evaluating the model with train data..
time taken : 0:00:00.140232
-----
Train Data
-----
RMSE : 0.9873944686749756

MAPE : 30.820930111024925

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.049481
-----
Test Data
-----
RMSE : 1.0657748808206107

MAPE : 34.17509113922799

storing the test results in test dictionary...

-----
Total time taken to run this algorithm : 0:00:00.288106

```

4.4.3 XGBoost with initial 13 features + Surprise Baseline predictor

Updating Train Data

In [190]:

```

# add our baseline_predicted value as our feature..
reg_train['bslpr'] = models_evaluation_train['bsl_algo']['predictions']
reg_train.head(2)

```

Out[190]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating	bslpr	kn
0	2442	1	3.618026	4	4	4	4	5	3.0	3.0	3.0	3.0	3.0	3.0	3.640625	3	3.628458	3.0
1	31913	1	3.618026	5	4	4	4	5	1.0	2.5	2.5	2.5	2.5	2.5	3.640625	4	3.593517	3.0

Updating Test Data

In [191]:

```

# add that baseline_predicted ratings with Surprise to the test data as well
reg_test_df['bslpr'] = models_evaluation_test['bsl_algo']['predictions']

reg_test_df.head(2)

```

Out[191]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MA
0	206115	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
1	314933	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406

In [0]:

```

#splitting to x and y for computing hyperparameter tuning
x_train_data=reg_train.drop(['user', 'movie', 'rating'], axis=1)
y_train_data=reg_train['rating']

x_test_data=reg_test_df.drop(['user', 'movie', 'rating'], axis=1)

```

```
y_test_data=reg_test_df['rating']
```

In [193]:

```
train_rmse = []
test_rmse = []
max_depth = [1, 5, 10, 50, 100, 500, 700, 800, 900, 1000]
n_models = [10, 20, 30, 40, 50, 60, 70, 80, 90, 100]

for i in max_depth:
    xgb_model = xgb.XGBRegressor(max_depth=i, n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score = rmse(y_train_pred, y_train_data)
    train_rmse.append(train_rmse_score)
    test_rmse_score = rmse(y_test_pred, y_test_data)
    test_rmse.append(test_rmse_score)
    print("Best_depth = ", i, "\t", "test_rmse_score\t:", test_rmse_score, "\t", "train_rmse_score\t:",
          train_rmse_score)

best_depth = np.argmin(test_rmse)
print("best_depth :", max_depth[best_depth])
```

```
[12:33:09] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1    test_rmse_score : 1.118686519152436    train_rmse_score : 0.4742197391807634
[12:33:09] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 5    test_rmse_score : 1.1165563471004316    train_rmse_score : 0.20185295107101406
[12:33:11] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 10   test_rmse_score : 1.1853701654112845    train_rmse_score : 0.029068982365314522
[12:33:17] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 50   test_rmse_score : 1.2323091642703456    train_rmse_score : 0.0012429298136426526
[12:33:33] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 100  test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
[12:33:49] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 500  test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
[12:34:05] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 700  test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
[12:34:22] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 800  test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
[12:34:39] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 900  test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
[12:34:56] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1000 test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
best_depth : 5
```

In [194]:

```
s_train_rmse=[]
s_test_rmse=[]

for j in n_models:
    xgb_model = xgb.XGBRegressor(max_depth=max_depth[best_depth], n_estimators=j, n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score = rmse(y_train_pred, y_train_data)
    s_train_rmse.append(train_rmse_score)
```



```

s_train_rmse.append(train_rmse_score)
test_rmse_score=rmse(y_test_pred ,y_test_data)
s_test_rmse.append(test_rmse_score)
print("best_estimators = ",j ,"\t", "test_rmse_score\t:",test_rmse_score, "\t", "train_rmse_score
\t:",train_rmse_score)

best_estimators = np.argmin(s_test_rmse)
print("best_estimators :",n_models[best_estimators])

```

```

[12:35:13] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 10    test_rmse_score : 1.4804614168768708    train_rmse_score :
1.2152385499466214
[12:35:13] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 20    test_rmse_score : 1.1699152724516488    train_rmse_score :
0.5403831999609053
[12:35:14] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 30    test_rmse_score : 1.1209988437594627    train_rmse_score :
0.3543631578919108
[12:35:14] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 40    test_rmse_score : 1.114834216555102    train_rmse_score :
0.3015962804147239
[12:35:15] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 50    test_rmse_score : 1.1136777693704727    train_rmse_score :
0.27549062766421706
[12:35:16] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 60    test_rmse_score : 1.1139789327483343    train_rmse_score :
0.25435819750785754
[12:35:18] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 70    test_rmse_score : 1.1144393760518252    train_rmse_score :
0.2387096275277416
[12:35:19] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 80    test_rmse_score : 1.1150160301961656    train_rmse_score :
0.2252943871448911
[12:35:21] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 90    test_rmse_score : 1.1160731560467065    train_rmse_score :
0.21292301246595102
[12:35:23] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 100   test_rmse_score : 1.1165563471004316    train_rmse_score :
0.20185295107101406
best_estimators : 50

```

In [195]:

```

# prepare train data
x_train = reg_train.drop(['user', 'movie', 'rating'], axis=1)
y_train = reg_train['rating']

# Prepare Test data
x_test = reg_test_df.drop(['user', 'movie', 'rating'], axis=1)
y_test = reg_test_df['rating']

# initialize Our first XGBoost model...
xgb_bsl = xgb.XGBRegressor(silent=1, n_jobs=13, random_state=15, n_estimators=n_models[best_estimat
ors], max_depth=max_depth[best_depth])
train_results, test_results = run_xgboost(xgb_bsl, x_train, y_train, x_test, y_test)

# store the results in models evaluations dictionaries
models_evaluation_train['xgb_bsl'] = train_results
models_evaluation_test['xgb_bsl'] = test_results

xgb.plot_importance(xgb_bsl)
plt.show()

```

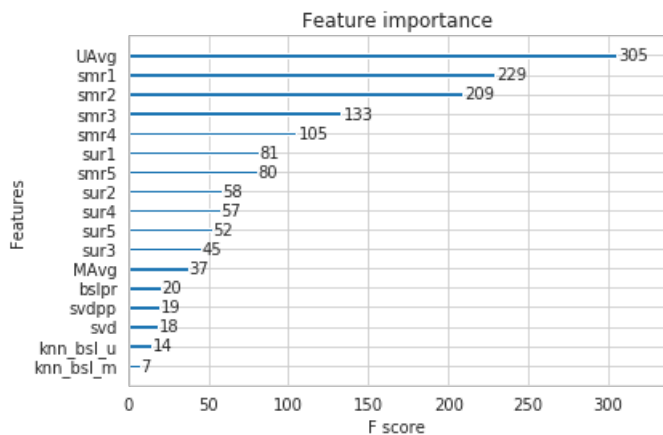
Training the model..
 Done. Time taken : 0:00:01.096367

Done

Evaluating the model with TRAIN data...
 Evaluating Test data

TEST DATA

 RMSE : 1.1136777693704727
 MAPE : 35.85766377810637



4.4.4 Surprise KNNBaseline predictor

In [0]:

```
from surprise import KNNBaseline
```

- KNN BASELINE
 - http://surprise.readthedocs.io/en/stable/knn_inspired.html#surprise.prediction_algorithms.knns.KNNBaseline
- PEARSON_BASELINE SIMILARITY
 - http://surprise.readthedocs.io/en/stable/similarities.html#surprise.similarities.pearson_baseline
- SHRINKAGE
 - 2.2 Neighborhood Models in <http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf>
- **predicted Rating : (based on User-User similarity)**

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{v \in N^k_i(u)} \text{sim}(u, v) \cdot (r_{vi} - b_{vi})}{\sum_{v \in N^k_i(u)} \text{sim}(u, v)}$$

- b_{ui} - Baseline prediction of (user, movie) rating
- $N^k_i(u)$ - Set of **K** similar users (neighbours) of **user (u)** who rated **movie(i)**
- $\text{sim}(u, v)$ - **Similarity** between users **u** and **v**
 - Generally, it will be cosine similarity or Pearson correlation coefficient.
 - But we use **shrunk Pearson-baseline correlation coefficient**, which is based on the pearsonBaseline similarity (we take base line predictions instead of mean rating of user/item)

- **Predicted rating (based on Item Item similarity):**
$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{j \in N^k_u(i)} \text{sim}(i, j) \cdot (r_{uj} - b_{uj})}{\sum_{j \in N^k_u(i)} \text{sim}(i, j)}$$
 - **Notations follows same as above (user user based predicted rating)**

4.4.4.1 Surprise KNNBaseline with user user similarities

In [197]:

```

# we specify , how to compute similarities and what to consider with sim_options to our algorithm
sim_options = {'user_based' : True,
               'name': 'pearson_baseline',
               'shrinkage': 100,
               'min_support': 2
              }

# we keep other parameters like regularization parameter and learning_rate as default values.
bsl_options = {'method': 'sgd'}

knn_bsl_u = KNNBaseline(k=40, sim_options = sim_options, bsl_options = bsl_options)
knn_bsl_u_train_results, knn_bsl_u_test_results = run_surprise(knn_bsl_u, trainset, testset,
verbose=True)

# Just store these error metrics in our models_evaluation datastructure
models_evaluation_train['knn_bsl_u'] = knn_bsl_u_train_results
models_evaluation_test['knn_bsl_u'] = knn_bsl_u_test_results

```

```

Training the model...
Estimating biases using sgd...
Computing the pearson_baseline similarity matrix...
Done computing similarity matrix.
Done. time taken : 0:00:13.587670

```

```

Evaluating the model with train data..
time taken : 0:00:53.970755

```

```

-----
Train Data
-----

```

```

RMSE : 0.05047943692266681

```

```

MAPE : 0.8026702832369823

```

```

adding train results in the dictionary..

```

```

Evaluating for test data...
time taken : 0:00:00.254883

```

```

-----
Test Data
-----

```

```

RMSE : 1.0673232737369938

```

```

MAPE : 34.26679323448809

```

```

storing the test results in test dictionary...

```

```

-----
Total time taken to run this algorithm : 0:01:07.817291

```

4.4.4.2 Surprise KNNBaseline with movie movie similarities

```

In [198]:

```

```

# we specify , how to compute similarities and what to consider with sim_options to our algorithm

# 'user_based' : Fals => this considers the similarities of movies instead of users

sim_options = {'user_based' : False,
               'name': 'pearson_baseline',
               'shrinkage': 100,
               'min_support': 2
              }

# we keep other parameters like regularization parameter and learning_rate as default values.
bsl_options = {'method': 'sgd'}

knn_bsl_m = KNNBaseline(k=40, sim_options = sim_options, bsl_options = bsl_options)

knn_bsl_m_train_results, knn_bsl_m_test_results = run_surprise(knn_bsl_m, trainset, testset,
verbose=True)

# Just store these error metrics in our models_evaluation datastructure
models_evaluation_train['knn_bsl_m'] = knn_bsl_m_train_results
models_evaluation_test['knn_bsl_m'] = knn_bsl_m_test_results

```


	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAVg	MA
	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAVg	MA
0	206115	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
1	314933	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406

In [0]:

```
#splitting to x and y for computing hyperparameter tuning
x_train_data=reg_train.drop(['user', 'movie','rating'], axis=1)
y_train_data=reg_train['rating']

x_test_data=reg_test_df.drop(['user', 'movie','rating'], axis=1)
y_test_data=reg_test_df['rating']
```

In [202]:

```
train_rmse = []
test_rmse = []
max_depth = [1, 5, 10, 50, 100, 500,700,800,900 ,1000]
n_models =[10,20,30,40,50,60,70,80,90,100]

for i in max_depth:
    xgb_model = xgb.XGBRegressor(max_depth=i,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred ,y_test_data)
    test_rmse.append(test_rmse_score)
    print("Best_depth = ",i ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score\t:",
train_rmse_score)

best_depth = np.argmin(test_rmse)
print("best_depth :",max_depth[best_depth])
```

```
[12:36:36] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1 test_rmse_score : 1.118686519152436 train_rmse_score : 0.4742197391807634
[12:36:36] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 5 test_rmse_score : 1.1165563471004316 train_rmse_score : 0.20185295107101406
[12:36:38] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 10 test_rmse_score : 1.1853701654112845 train_rmse_score : 0.029068982365314522
[12:36:44] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 50 test_rmse_score : 1.2323091642703456 train_rmse_score : 0.0012429298136426526
[12:37:00] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 100 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:37:16] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 500 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:37:32] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 700 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:37:49] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 800 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:38:06] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 900 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
[12:38:23] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1000 test_rmse_score : 1.2323091642703456 train_rmse_score :
0.0012429298136426526
best_depth : 5
```

In [203]:

```
s_train_rmse=[]
s_test_rmse=[]

for j in n_models:
    xgb_model = xgb.XGBRegressor(max_depth=max_depth[best_depth],n_estimators=j,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    s_train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred ,y_test_data)
    s_test_rmse.append(test_rmse_score)
    print("best_estimators = ",j ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score
\t:",train_rmse_score)

best_estimators = np.argmin(s_test_rmse)
print("best_depth :",max_depth[best_depth])
print("best_estimators :",n_models[best_estimators])
```

```
[12:38:40] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 10    test_rmse_score : 1.4804614168768708    train_rmse_score :
1.2152385499466214
[12:38:40] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 20    test_rmse_score : 1.1699152724516488    train_rmse_score :
0.5403831999609053
[12:38:40] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 30    test_rmse_score : 1.1209988437594627    train_rmse_score :
0.3543631578919108
[12:38:41] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 40    test_rmse_score : 1.114834216555102    train_rmse_score :
0.3015962804147239
[12:38:42] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 50    test_rmse_score : 1.1136777693704727    train_rmse_score :
0.27549062766421706
[12:38:43] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 60    test_rmse_score : 1.1139789327483343    train_rmse_score :
0.25435819750785754
[12:38:45] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 70    test_rmse_score : 1.1144393760518252    train_rmse_score :
0.2387096275277416
[12:38:46] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 80    test_rmse_score : 1.1150160301961656    train_rmse_score :
0.2252943871448911
[12:38:48] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 90    test_rmse_score : 1.1160731560467065    train_rmse_score :
0.21292301246595102
[12:38:50] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 100   test_rmse_score : 1.1165563471004316    train_rmse_score :
0.20185295107101406
best_depth : 5
best_estimators : 50
```

In [204]:

```
# prepare the train data....
x_train = reg_train.drop(['user', 'movie', 'rating'], axis=1)
y_train = reg_train['rating']

# prepare the train data....
x_test = reg_test_df.drop(['user','movie','rating'], axis=1)
y_test = reg_test_df['rating']
```

```
# declare the model
xgb_knn_bsl = xgb.XGBRegressor(n_jobs=10,
random_state=15,max_depth=max_depth[best_depth],n_estimators=n_models[best_estimators])
train_results, test_results = run_xgboost(xgb_knn_bsl, x_train, y_train, x_test, y_test)

# store the results in models_evaluations dictionaries
models_evaluation_train['xgb_knn_bsl'] = train_results
models_evaluation_test['xgb_knn_bsl'] = test_results

xgb.plot_importance(xgb_knn_bsl)
plt.show()
```

Training the model..

[12:38:52] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

Done. Time taken : 0:00:01.071529

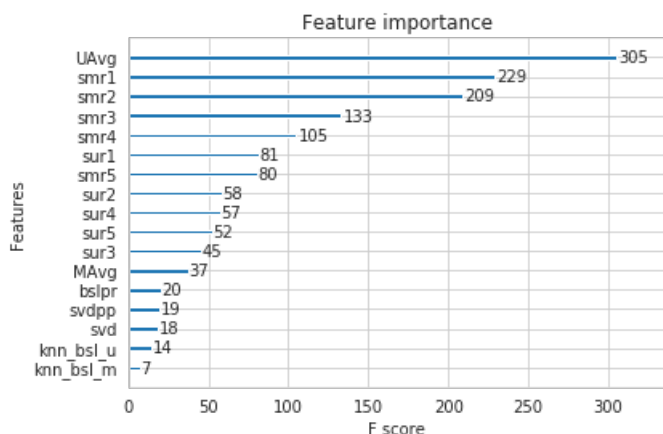
Done

Evaluating the model with TRAIN data...

Evaluating Test data

TEST DATA

```
-----
RMSE : 1.1136777693704727
MAPE : 35.85766377810637
```



4.4.6 Matrix Factorization Techniques

4.4.6.1 SVD Matrix Factorization User Movie interactions

In [0]:

```
from surprise import SVD
```

http://surprise.readthedocs.io/en/stable/matrix_factorization.html#surprise.prediction_algorithms.matrix_factorization.SVD

- Predicted Rating :

$$\hat{r}_{ui} = \mu + b_u + b_i + q_i^T p_u$$

- q_i - Representation of item(movie) in latent factor space

- p_u - Representation of user in new latent factor space

- A BASIC MATRIX FACTORIZATION MODEL in [https://datajobs.com/data-science-repo/Recommender-Systems-Netflix\].pdf](https://datajobs.com/data-science-repo/Recommender-Systems-Netflix.pdf)

- Optimization problem with user item interactions and regularization (to avoid overfitting)

$$- \sum_{\{r_{ui}\} \in R_{\{train\}}} \left(r_{ui} - \hat{r}_{ui} \right)^2 +$$

$$\lambda \left(b_u^2 + b_i^2 + \|q_i\|^2 + \|p_u\|^2 \right)$$

In [206]:

```
# initialize the model
svd = SVD(n_factors=100, biased=True, random_state=15)
svd_train_results, svd_test_results = run_surprise(svd, trainset, testset)

# Just store these error metrics in our models_evaluation datastructure
models_evaluation_train['svd'] = svd_train_results
models_evaluation_test['svd'] = svd_test_results
```

Training the model...

Done. time taken : 0:00:01.358997

Evaluating the model with train data..

time taken : 0:00:00.173458

Train Data

RMSE : 0.6638184366771201

MAPE : 19.883787468517028

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.141175

Test Data

RMSE : 1.0666587448560931

MAPE : 34.2032674271826

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:00:01.679097

4.4.6.2 SVD Matrix Factorization with implicit feedback from user (user rated movies)

In [0]:

```
from surprise import SVDpp
```

- -----> 2.5 Implicit Feedback in <http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf>

- Predicted Rating :

$$- \hat{r}_{ui} = \mu + b_u + b_i + q_i^T \left(p_u + \frac{1}{|I_u|} \sum_{j \in I_u} y_j \right)$$

- I_u --- the set of all items rated by user u
- y_j --- Our new set of item factors that capture implicit ratings.

- Optimization problem with user item interactions and regularization (to avoid overfitting)

$$- \sum_{ui \in R_{\text{train}}} (r_{ui} - \hat{r}_{ui})^2 +$$

$$\lambda (b_i^2 + b_u^2 + \|q_i\|^2 + \|p_u\|^2 + \|y_j\|^2)$$

In [208]:

```
# initialize the model
svdpp = SVDpp(n_factors=50, random_state=15)
svdpp_train_results, svdpp_test_results = run_surprise(svdpp, trainset, testset)

# Just store these error metrics in our models_evaluation datastructure
models_evaluation_train['svdpp'] = svdpp_train_results
models_evaluation_test['svdpp'] = svdpp_test_results
```

```
Training the model...
Done. time taken : 0:00:03.174091

Evaluating the model with train data..
time taken : 0:00:00.402840
-----
Train Data
-----
RMSE : 0.5669418677286878

MAPE : 17.01786548110549

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.061331
-----
Test Data
-----
RMSE : 1.0716797975344945

MAPE : 34.94221562143592

storing the test results in test dictionary...

-----
Total time taken to run this algorithm : 0:00:03.641526
```

4.4.7 XgBoost with 13 features + Surprise Baseline + Surprise KNNbaseline + MF Techniques

Preparing Train data

In [209]:

```
# add the predicted values from both knns to this dataframe
reg_train['svd'] = models_evaluation_train['svd']['predictions']
reg_train['svdpp'] = models_evaluation_train['svdpp']['predictions']

reg_train.head(2)
```

Out[209]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating	bslpr	kn
0	2442	1	3.618026	4	4	4	4	5	3.0	3.0	3.0	3.0	3.0	3.0	3.640625	3	3.628458	3.0
1	31913	1	3.618026	5	4	4	4	5	1.0	2.5	2.5	2.5	2.5	2.5	3.640625	4	3.593517	3.0

Preparing Test data

In [210]:

```
reg_test_df['svd'] = models_evaluation_test['svd']['predictions']
reg_test_df['svdpp'] = models_evaluation_test['svdpp']['predictions']

reg_test_df.head(2)
```

Out[210]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MA
0	206115	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406
1	314933	1	3.618026	4	3	4	3	5	3.618026	3.618026	3.618026	3.618026	3.618026	3.618026	3.6406

In [0]:

```
#splitting to x and y for computing hyperparameter tuning
x_train_data=reg_train.drop(['user', 'movie','rating'], axis=1)
y_train_data=reg_train['rating']

x_test_data=reg_test_df.drop(['user', 'movie','rating'], axis=1)
y_test_data=reg_test_df['rating']
```

In [212]:

```
train_rmse = []
test_rmse = []
max_depth = [1, 5, 10, 50, 100, 500,700,800,900 ,1000]
n_models = [10,20,30,40,50,60,70,80,90,100]

for i in max_depth:
    xgb_model = xgb.XGBRegressor(max_depth=i,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred ,y_test_data)
    test_rmse.append(test_rmse_score)
    print("Best_depth = ",i ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score\t:",
train_rmse_score)

best_depth = np.argmin(test_rmse)
print("best_depth :",max_depth[best_depth])
```

```
[12:39:00] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1   test_rmse_score : 1.118686519152436   train_rmse_score : 0.4742197391807634
[12:39:00] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 5   test_rmse_score : 1.1165563471004316   train_rmse_score : 0.20185295107101406
[12:39:02] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 10  test_rmse_score : 1.1853701654112845   train_rmse_score : 0.029068982365314522
[12:39:08] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 50  test_rmse_score : 1.2323091642703456   train_rmse_score : 0.0012429298136426526
[12:39:23] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 100 test_rmse_score : 1.2323091642703456   train_rmse_score :
0.0012429298136426526
[12:39:38] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 500 test_rmse_score : 1.2323091642703456   train_rmse_score :
0.0012429298136426526
[12:39:53] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 700 test_rmse_score : 1.2323091642703456   train_rmse_score :
0.0012429298136426526
[12:40:09] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 800 test_rmse_score : 1.2323091642703456   train_rmse_score :
0.0012429298136426526
[12:40:25] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 900 test_rmse_score : 1.2323091642703456   train_rmse_score :
0.0012429298136426526
```

```
[12:40:41] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1000    test_rmse_score : 1.2323091642703456    train_rmse_score :
0.0012429298136426526
best_depth : 5
```

In [213]:

```
s_train_rmse=[]
s_test_rmse=[]

for j in n_models:
    xgb_model = xgb.XGBRegressor(max_depth=max_depth[best_depth],n_estimators=j,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    s_train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred,y_test_data)
    s_test_rmse.append(test_rmse_score)
    print("best_estimators = ",j,"\\t","test_rmse_score\\t:",test_rmse_score, "\\t","train_rmse_score
\\t:",train_rmse_score)

best_estimators = np.argmin(s_test_rmse)
print("best_depth :",max_depth[best_depth])
print("best_estimators :",n_models[best_estimators])
```

```
[12:40:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 10    test_rmse_score : 1.4804614168768708    train_rmse_score :
1.2152385499466214
[12:40:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 20    test_rmse_score : 1.1699152724516488    train_rmse_score :
0.5403831999609053
[12:40:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 30    test_rmse_score : 1.1209988437594627    train_rmse_score :
0.3543631578919108
[12:40:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 40    test_rmse_score : 1.114834216555102    train_rmse_score :
0.3015962804147239
[12:40:59] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 50    test_rmse_score : 1.1136777693704727    train_rmse_score :
0.27549062766421706
[12:41:00] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 60    test_rmse_score : 1.1139789327483343    train_rmse_score :
0.25435819750785754
[12:41:02] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 70    test_rmse_score : 1.1144393760518252    train_rmse_score :
0.2387096275277416
[12:41:03] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 80    test_rmse_score : 1.1150160301961656    train_rmse_score :
0.2252943871448911
[12:41:05] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 90    test_rmse_score : 1.1160731560467065    train_rmse_score :
0.21292301246595102
[12:41:07] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 100    test_rmse_score : 1.1165563471004316    train_rmse_score :
0.20185295107101406
best_depth : 5
best_estimators : 50
```

In [214]:

```
# prepare x_train and y_train
x_train = reg_train.drop(['user', 'movie', 'rating'], axis=1)
y_train = reg_train['rating']
```

```

y_train = reg_train['rating']

# prepare test data
x_test = reg_test_df.drop(['user', 'movie', 'rating'], axis=1)
y_test = reg_test_df['rating']

xgb_final = xgb.XGBRegressor(n_jobs=10,
random_state=15,max_depth=max_depth[best_depth],n_estimators=n_models[best_estimators])
train_results, test_results = run_xgboost(xgb_final, x_train, y_train, x_test, y_test)

# store the results in models_evaluations dictionaries
models_evaluation_train['xgb_final'] = train_results
models_evaluation_test['xgb_final'] = test_results

xgb.plot_importance(xgb_final)
plt.show()

```

Training the model..

[12:41:09] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

Done. Time taken : 0:00:01.052774

Done

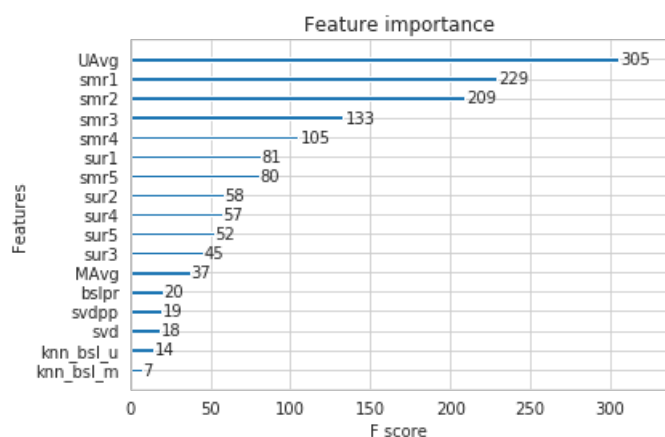
Evaluating the model with TRAIN data...

Evaluating Test data

TEST DATA

RMSE : 1.1136777693704727

MAPE : 35.85766377810637



In [215]:

```
reg_train.head()
```

Out[215]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating	bslpr	k
0	2442	1	3.618026	4	4	4	4	5	3.0	3.0	3.0	3.00	3.00	3.00	3.640625	3	3.628458	3
1	31913	1	3.618026	5	4	4	4	5	1.0	2.5	2.5	2.50	2.50	2.50	3.640625	4	3.593517	3
2	42930	1	3.618026	4	4	3	2	3	3.0	4.0	5.0	3.00	3.60	3.60	3.640625	3	3.748415	1
3	94565	1	3.618026	5	5	5	3	5	3.0	3.0	3.0	3.25	3.25	3.25	3.640625	4	3.613201	3
4	145873	1	3.618026	5	5	4	2	3	5.0	4.0	4.0	4.00	4.00	4.00	3.640625	3	3.768099	2

4.4.8 XgBoost with Surprise Baseline + Surprise KNNbaseline + MF Techniques

In [0]:

```
#splitting to x and y for computing hyperparameter tuning
x_train_data=reg_train[['bslpr','knn_bsl_u', 'knn_bsl_m', 'svd', 'svdpp']]
y_train_data=reg_train['rating']

x_test_data=reg_test_df[['bslpr','knn_bsl_u', 'knn_bsl_m', 'svd', 'svdpp']]
y_test_data=reg_test_df['rating']
```

In [217]:

```
train_rmse = []
test_rmse = []
max_depth = [1, 5, 10, 50, 100, 500,700,800,900 ,1000]
n_models =[10,20,30,40,50,60,70,80,90,100]

for i in max_depth:
    xgb_model = xgb.XGBRegressor(max_depth=i,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y_test_pred = xgb_model.predict(x_test_data)
    train_rmse_score=rmse(y_train_pred,y_train_data)
    train_rmse.append(train_rmse_score)
    test_rmse_score=rmse(y_test_pred ,y_test_data)
    test_rmse.append(test_rmse_score)
    print("Best_depth = ",i ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score\t:",
train_rmse_score)

best_depth = np.argmin(test_rmse)
print("best_depth :",max_depth[best_depth])
```

```
[12:41:11] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1   test_rmse_score : 1.095765243854484   train_rmse_score : 1.0473823456707518
[12:41:11] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 5   test_rmse_score : 1.1063392084983343   train_rmse_score : 0.988656044557846
[12:41:12] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 10  test_rmse_score : 1.0900557295214153   train_rmse_score : 0.8314866625392264
[12:41:15] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 50  test_rmse_score : 1.121940708750139   train_rmse_score : 0.0546725478338043
[12:41:31] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 100 test_rmse_score : 1.1499376035083695   train_rmse_score : 0.05185161269123567
[12:41:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 500 test_rmse_score : 1.1304437141210129   train_rmse_score : 0.041706031995373184
[12:43:17] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 700 test_rmse_score : 1.1304437141210129   train_rmse_score : 0.041706031995373184
[12:45:01] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 800 test_rmse_score : 1.1304437141210129   train_rmse_score : 0.041706031995373184
[12:46:57] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 900 test_rmse_score : 1.1304437141210129   train_rmse_score : 0.041706031995373184
[12:49:06] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
Best_depth = 1000 test_rmse_score : 1.1304437141210129   train_rmse_score :
0.041706031995373184
best_depth : 10
```

In [218]:

```
s_train_rmse=[]
s_test_rmse=[]

for j in n_models:
    xgb_model = xgb.XGBRegressor(max_depth=max_depth[best_depth],n_estimators=j,n_jobs=-1)
    xgb_model.fit(x_train_data, y_train_data)
    y_train_pred = xgb_model.predict(x_train_data)
    y test pred = xqb model.predict(x test data)
```

```

train_rmse_score=rmse(y_train_pred,y_train_data)
s_train_rmse.append(train_rmse_score)
test_rmse_score=rmse(y_test_pred ,y_test_data)
s_test_rmse.append(test_rmse_score)
print("best_estimators = ",j ,"\t","test_rmse_score\t:",test_rmse_score, "\t","train_rmse_score
\t:",train_rmse_score)

best_estimators = np.argmin(s_test_rmse)
print("best_depth :",max_depth[best_depth])
print("best_estimators :",n_models[best_estimators])

```

```

[12:51:27] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 10    test_rmse_score : 1.5572833587721402    train_rmse_score :
1.4837894288555062
[12:51:27] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 20    test_rmse_score : 1.1667471790252508    train_rmse_score :
1.0390805120986386
[12:51:28] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 30    test_rmse_score : 1.1128331780492484    train_rmse_score :
0.9505227319707155
[12:51:29] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 40    test_rmse_score : 1.0919243845131785    train_rmse_score :
0.9253313283090095
[12:51:30] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 50    test_rmse_score : 1.08789838805482    train_rmse_score : 0.908762802443123
[12:51:31] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 60    test_rmse_score : 1.0895836164765629    train_rmse_score :
0.8913346190763096
[12:51:33] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 70    test_rmse_score : 1.0894324325642135    train_rmse_score :
0.8777394327273478
[12:51:35] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 80    test_rmse_score : 1.0894915137075443    train_rmse_score :
0.8618544788974313
[12:51:37] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 90    test_rmse_score : 1.0897862085578762    train_rmse_score :
0.8467641258153972
[12:51:40] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated i
n favor of reg:squarederror.
best_estimators = 100   test_rmse_score : 1.0900557295214153    train_rmse_score :
0.8314866625392264
best_depth : 10
best_estimators : 50

```

In [219]:

```

# prepare train data
x_train = reg_train[['bslpr', 'knn_bsl_u', 'knn_bsl_m', 'svd', 'svdpp']]
y_train = reg_train['rating']

# test data
x_test = reg_test_df[['bslpr', 'knn_bsl_u', 'knn_bsl_m', 'svd', 'svdpp']]
y_test = reg_test_df['rating']

xgb_all_models = xgb.XGBRegressor(n_jobs=10,
random_state=15,max_depth=max_depth[best_depth],n_estimators=n_models[best_estimators])
train_results, test_results = run_xgboost(xgb_all_models, x_train, y_train, x_test, y_test)

# store the results in models_evaluations dictionaries
models_evaluation_train['xgb_all_models'] = train_results
models_evaluation_test['xgb_all_models'] = test_results

xgb.plot_importance(xgb_all_models)

```

```
plt.show()
```

Training the model..

```
[12:51:43] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
```

Done. Time taken : 0:00:01.326331

Done

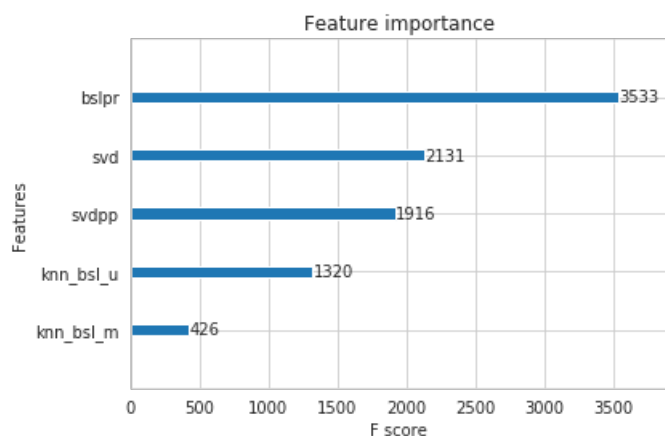
Evaluating the model with TRAIN data...

Evaluating Test data

TEST DATA

RMSE : 1.0878983880548427

MAPE : 33.948453881343205



4.5 Comparison between all models

In [229]:

```
# Saving our TEST_RESULTS into a dataframe so that you don't have to run it again
pd.DataFrame(models_evaluation_test).to_csv('drive/My Drive/sample_data/small_sample_results.csv')
models = pd.read_csv('drive/My Drive/sample_data/small_sample_results.csv', index_col=0)
models = models.drop(['best_depth', 'best_estimator'], axis=1)
models.loc['rmse'].sort_values()
```

Out[229]:

```
bsl_algo      1.0657748808206107
svd           1.0666587448560931
knn_bsl_u     1.0673232737369938
knn_bsl_m     1.0703781928309377
svdpp         1.0716797975344945
xgb_all_models 1.0878983880548427
first_algo    1.1136777693704727
xgb_bsl       1.1136777693704727
xgb_knn_bsl   1.1136777693704727
xgb_final     1.1136777693704727
Name: rmse, dtype: object
```

In [230]:

```
print("time taken to run this whole notebook is ",datetime.now() - globalstart)
```

time taken to run this whole notebook is 2:38:36.928600

5. Observations

1. In this Netflix movie recommendation case study we use 0.3 million data points for data-preprocessing and EDA due to less computational resources

2. For computing quick we convert our data into sparse representation because it takes less time to executing compare to raw data
3. We seen that average ratings for each movie and user average ratings to decide that user overall ratings status and overall movie status
4. We limited to compute user similarity because it takes days to implement with original data for that we take sample data 18k points for computing similarity with user and movies with surprise models and base line models
5. We extracted some features like SVD knn baseline algorithms using surprise models we also done feature extraction by simple manual created features which are some similar with base line models
6. We done every surprise model features and for every feature we combine those features with our baseline features and handcraft features we computed xgboost regressor model with optimal hyper parameter tuning
7. After done models for every surprise models and all xg boost models rmse score for some models get slightly lower and some models get slightly higher results compare to older results with out hyper parameter tuning
8. We can reduce rmse more by taking larger sample but due to less computational resources we got results based on small sampled based results
9. After observing the results base_line models and svd and knn user-user similarity gives similar and less rmse score compare to remaining models