



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 :

- 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
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 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 [1]: !pip install dask
```

```
Requirement already satisfied: dask in c:\users\user\anaconda3\lib\site-packages (0.16.1)
```

```
You are using pip version 18.0, however version 19.0.3 is available.
```

```
You should consider upgrading via the 'python -m pip install --upgrade pip' command.
```

```
In [7]: # this is just to know how much time will it take to run this entire ipython notebook
from datetime import datetime
# globalstart = datetime.now()
import pandas as pd
import dask.dataframe as dd
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 [3]: start = datetime.now()
if not os.path.isfile('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('data.csv', mode='w')

    row = list()
    files=['data_folder/combined_data_1.txt','data_folder/combined_data_2.txt',
           'data_folder/combined_data_3.txt', '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 : 0:00:00.000996


```
In [11]: print("creating the dataframe from data.csv file..")
df = pd.read_csv('data.csv', sep=',', 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 [12]: df.head()
```

Out[12]:

	movie	user	rating	date
56431994	10341	510180	4	1999-11-11
9056171	1798	510180	5	1999-11-11
58698779	10774	510180	3	1999-11-11
48101611	8651	510180	2	1999-11-11
81893208	14660	510180	2	1999-11-11

```
In [13]: df.describe()['rating']
```

```
Out[13]: count    1.004805e+08  
         mean     3.604290e+00  
         std      1.085219e+00  
         min      1.000000e+00  
         25%      3.000000e+00  
         50%      4.000000e+00  
         75%      4.000000e+00  
         max      5.000000e+00  
         Name: rating, dtype: float64
```

3.1.2 Checking for NaN values

```
In [6]: # 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 [7]: 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 [8]: 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 : 100480507
Total No of Users   : 480189
Total No of movies  : 17770
```

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

```
In [14]: if not os.path.isfile('train.csv'):
# create the dataframe and store it in the disk for offline purposes..
df.iloc[:int(df.shape[0]*0.80)].to_csv("train.csv", index=False)

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

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

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

```
In [10]: # 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 : 80384405
Total No of Users   : 405041
Total No of movies  : 17424
```

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

```
In [11]: 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 : 20096102
Total No of Users   : 349312
Total No of movies  : 17757
```

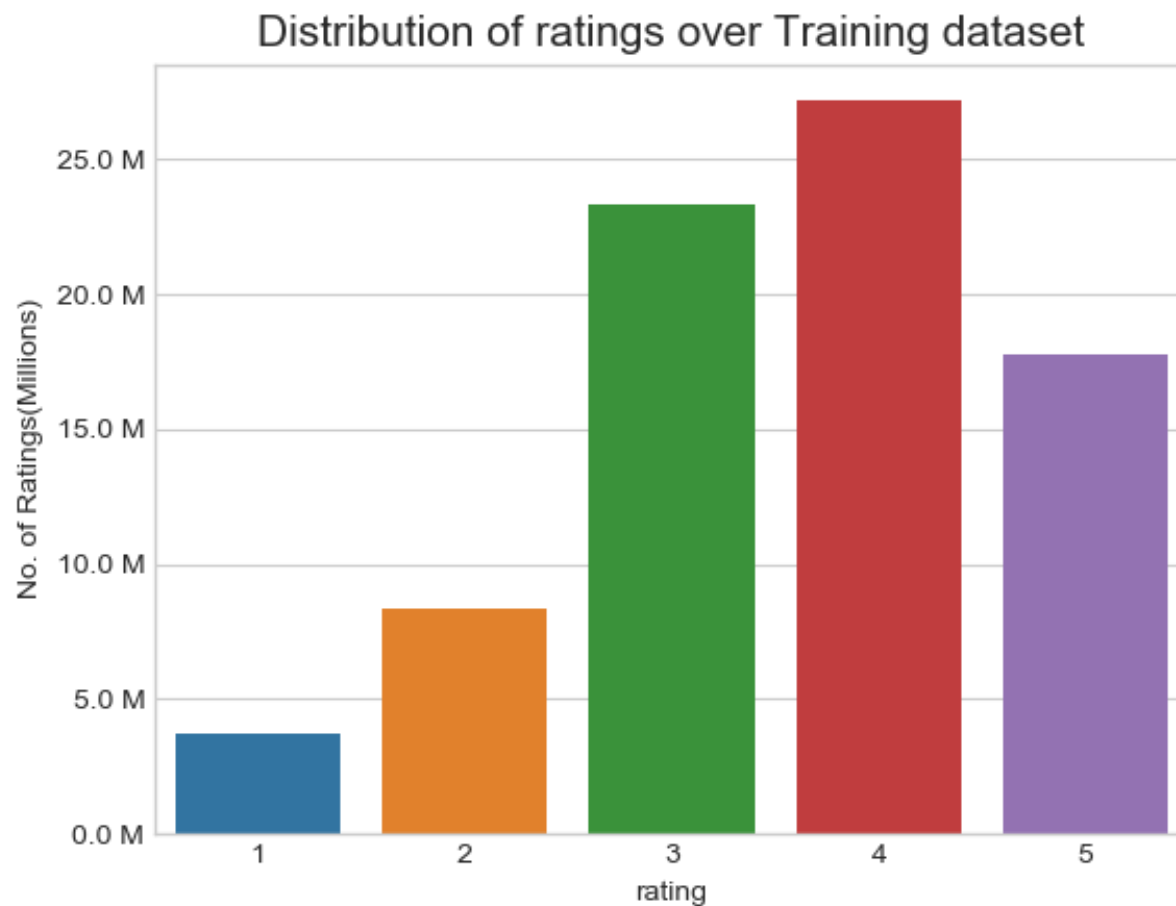
3.3 Exploratory Data Analysis on Train data

```
In [15]: # 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

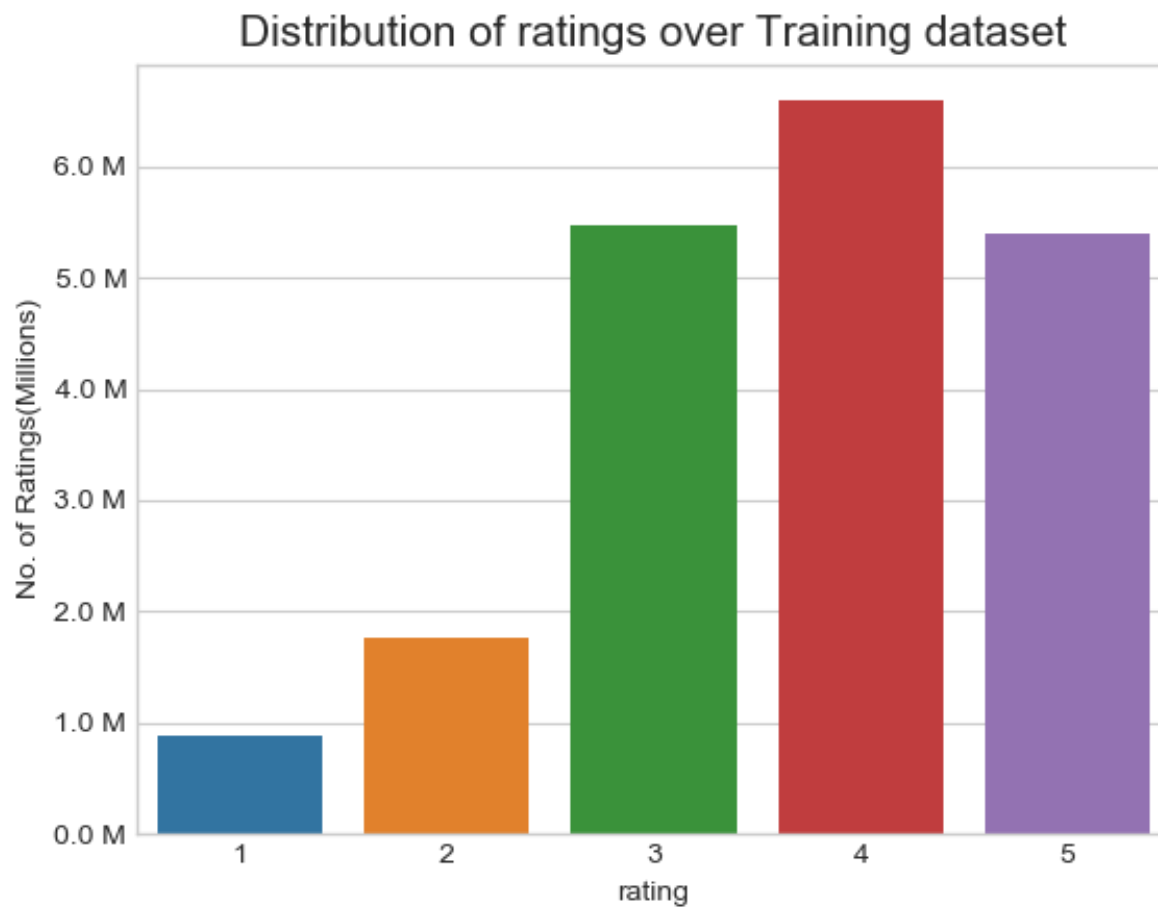
```
In [16]: fig, ax = plt.subplots()
plt.title('Distribution of ratings over Training dataset', fontsize=15)
sns.countplot(train_df.rating)
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
ax.set_ylabel('No. of Ratings(Millions)')

plt.show()
```



```
In [30]: fig, ax = plt.subplots()
plt.title('Distribution of ratings over Training dataset', fontsize=15)
sns.countplot(test_df.rating)
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
ax.set_ylabel('No. of Ratings(Millions)')

plt.show()
```



Obseration:

- The train and the test data is some what is not equally distributed. From the above plots we can see that in test data ratings 3,5 are likely the same where as in the train data rating 5 is lower than rating 3.

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

```
In [17]: # It is used to skip the warning ''SettingWithCopyWarning' '..
pd.options.mode.chained_assignment = None # default='warn'

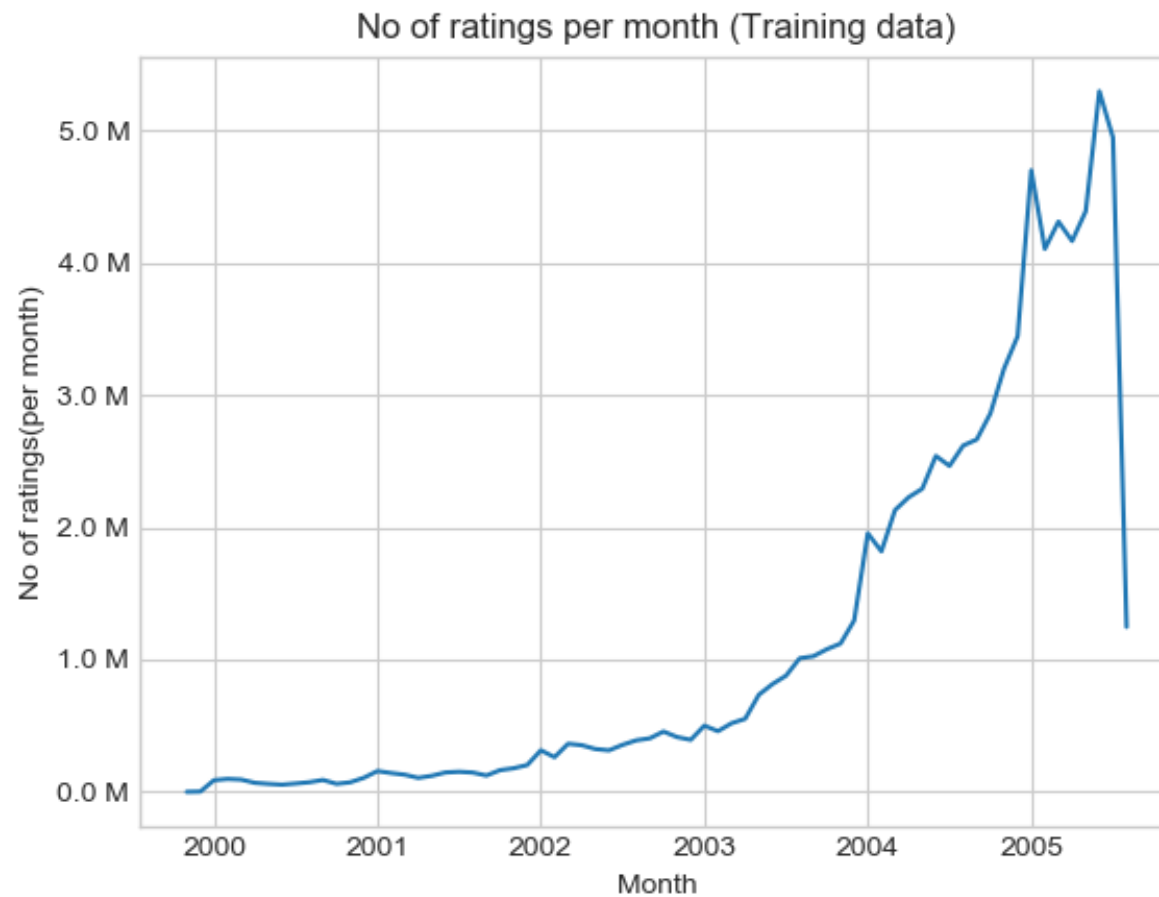
train_df['day_of_week'] = train_df.date.dt.weekday_nam
train_df.tail()
```

Out[17]:

	movie	user	rating	date	day_of_week
80384400	12074	2033618	4	2005-08-08	Monday
80384401	862	1797061	3	2005-08-08	Monday
80384402	10986	1498715	5	2005-08-08	Monday
80384403	14861	500016	4	2005-08-08	Monday
80384404	5926	1044015	5	2005-08-08	Monday

3.3.2 Number of Ratings per a month


```
In [18]: 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()
```



Observation:

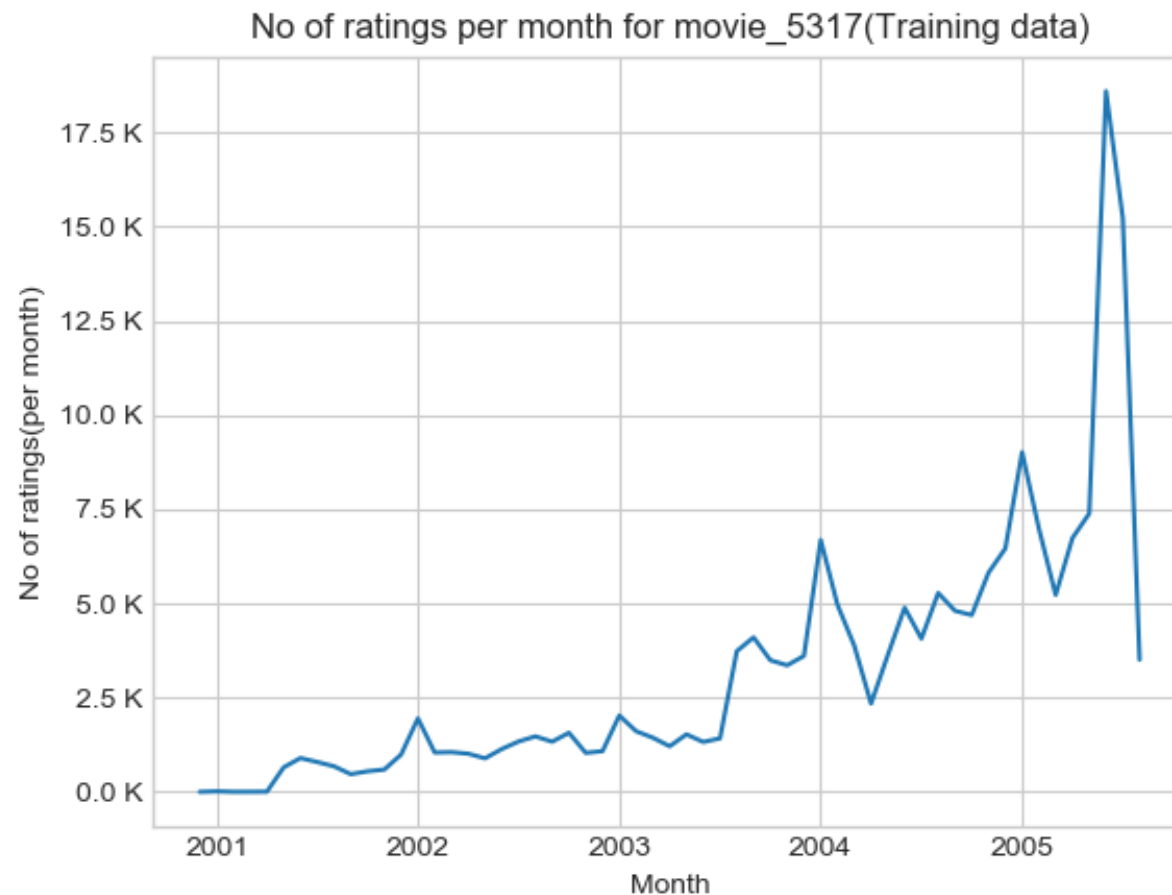
- The no.of ratings increased from 2003 may be because of the growth of the stated taking off from that same year.
- let's try plot the distribution for individual user.

```
In [46]: no_of_movies Rated = train_df.groupby(by='movie')['rating'].count().sort_values(ascending=False)
no_of_movies Rated.head()
```

```
Out[46]: movie
5317      179684
15124     176811
1905      160062
6287      155787
14313     153899
Name: rating, dtype: int64
```

```
In [50]: grouped = train_df.groupby('movie')
movie_5317 = grouped.get_group(5317)

ax = movie_5317.resample('m', on='date')['rating'].count().plot()
ax.set_title('No of ratings per month for movie_5317(Training data)')
plt.xlabel('Month')
plt.ylabel('No of ratings(per month)')
ax.set_yticklabels([human(item, 'k') for item in ax.get_yticks()])
plt.show()
```

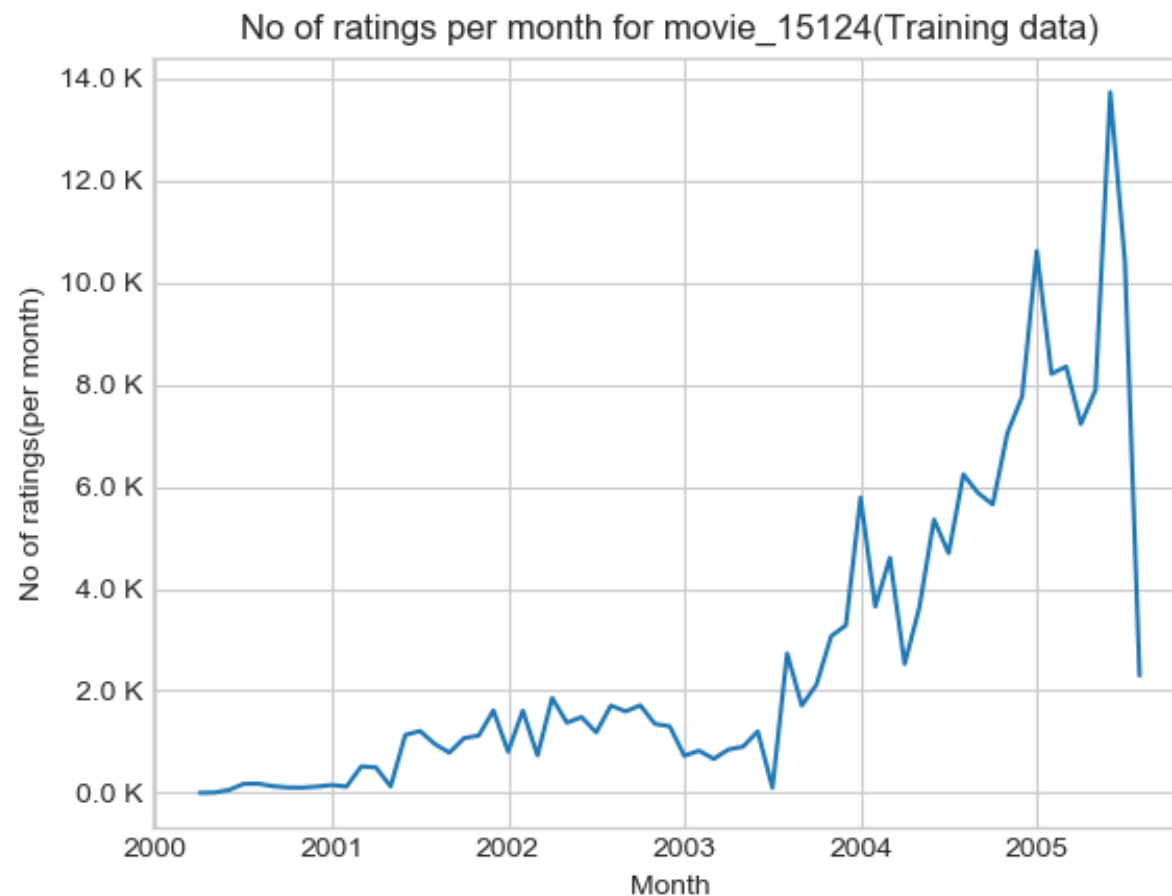


Observation:

- movie_5317 is similar to the total movie ratings.

```
In [48]: grouped = train_df.groupby('movie')
movie_15124 = grouped.get_group(15124)

ax = movie_15124.resample('m', on='date')['rating'].count().plot()
ax.set_title('No of ratings per month for movie_15124(Training data)')
plt.xlabel('Month')
plt.ylabel('No of ratings(per month)')
ax.set_yticklabels([human(item, 'k') for item in ax.get_yticks()])
plt.show()
```



Observation:

- movie_15124 is also similar to the total rating. so we can say that the movie which is added to the netflix has been show a steady growth since 2003.

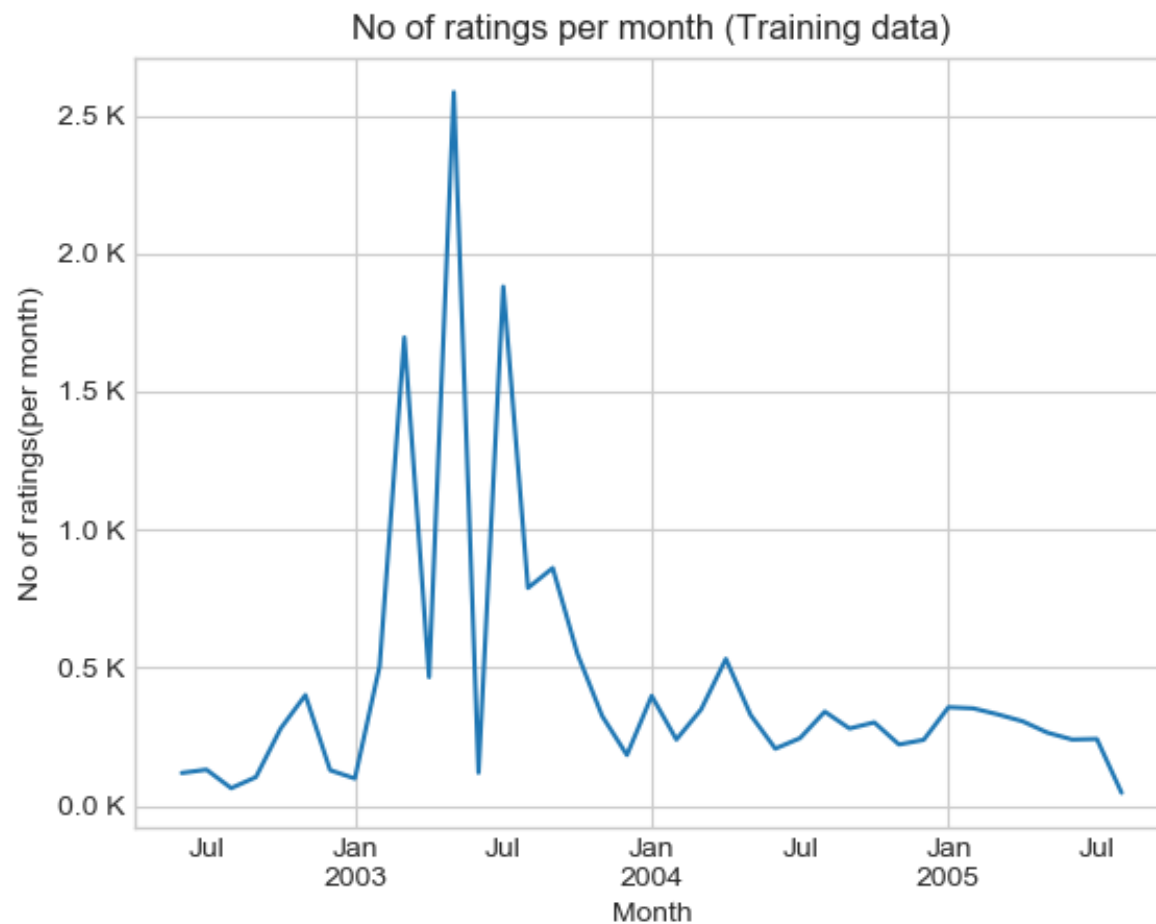
3.3.3 Analysis on the Ratings given by user

```
In [19]: 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[19]: user
         305344      17112
         2439493      15896
         387418       15402
         1639792        9767
         1461435        9447
         Name: rating, dtype: int64
```

```
In [39]: grouped = train_df.groupby('user')
         user_305344 = grouped.get_group(305344)
```

```
In [42]: ax = user_305344.resample('m', on='date')['rating'].count().plot()  
ax.set_title('No of ratings per month by user_305344(Training data)')  
plt.xlabel('Month')  
plt.ylabel('No of ratings(per month)')  
ax.set_yticklabels([human(item, 'k') for item in ax.get_yticks()])  
plt.show()
```



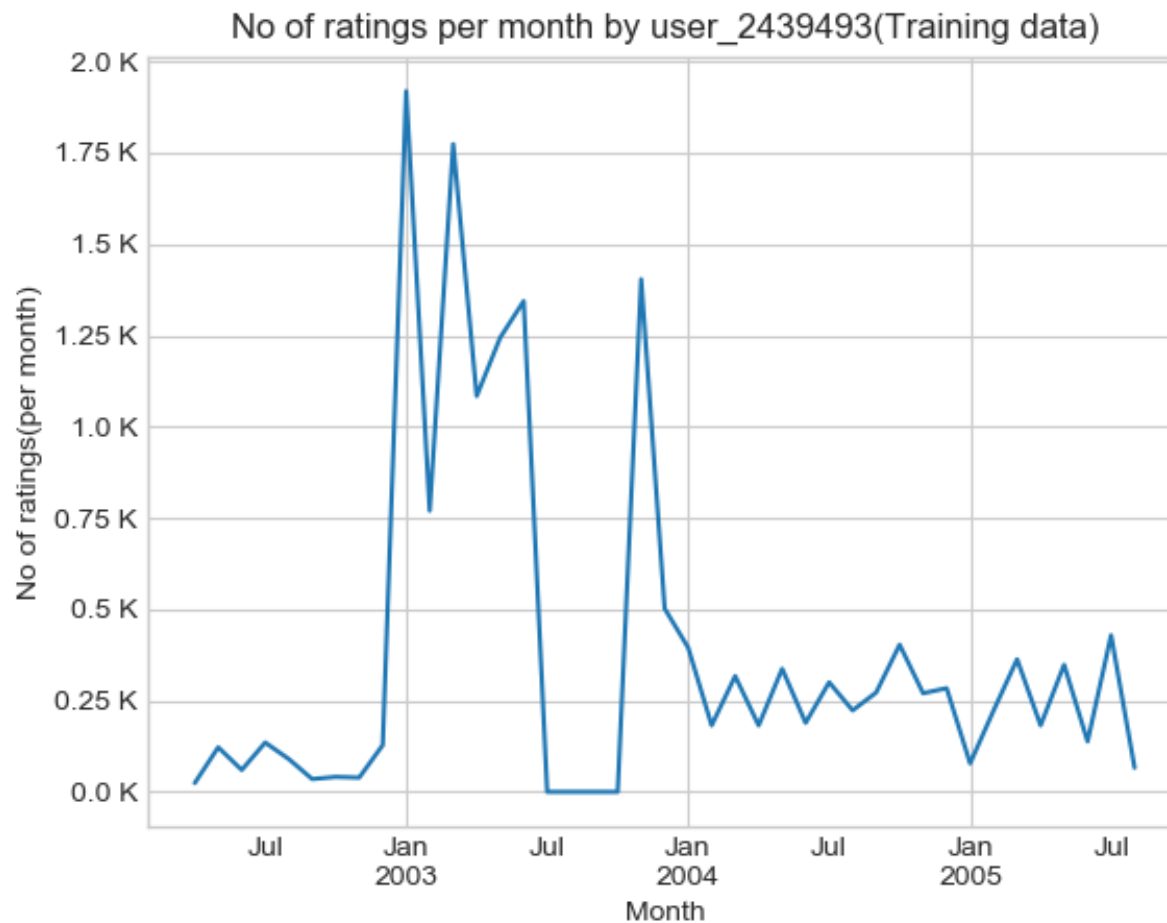
Observation:

- user_305344 has rated approx 7.5k movies in the span of one year their may be a chance that the user may rated some movies randomly.

```
In [43]: grouped = train_df.groupby('user')  
user_2439493 = grouped.get_group(2439493)
```



```
In [44]: ax = user_2439493.resample('m', on='date')['rating'].count().plot()  
ax.set_title('No of ratings per month by user_2439493(Training data)')  
plt.xlabel('Month')  
plt.ylabel('No of ratings(per month)')  
ax.set_yticklabels([human(item, 'k') for item in ax.get_yticks()])  
plt.show()
```

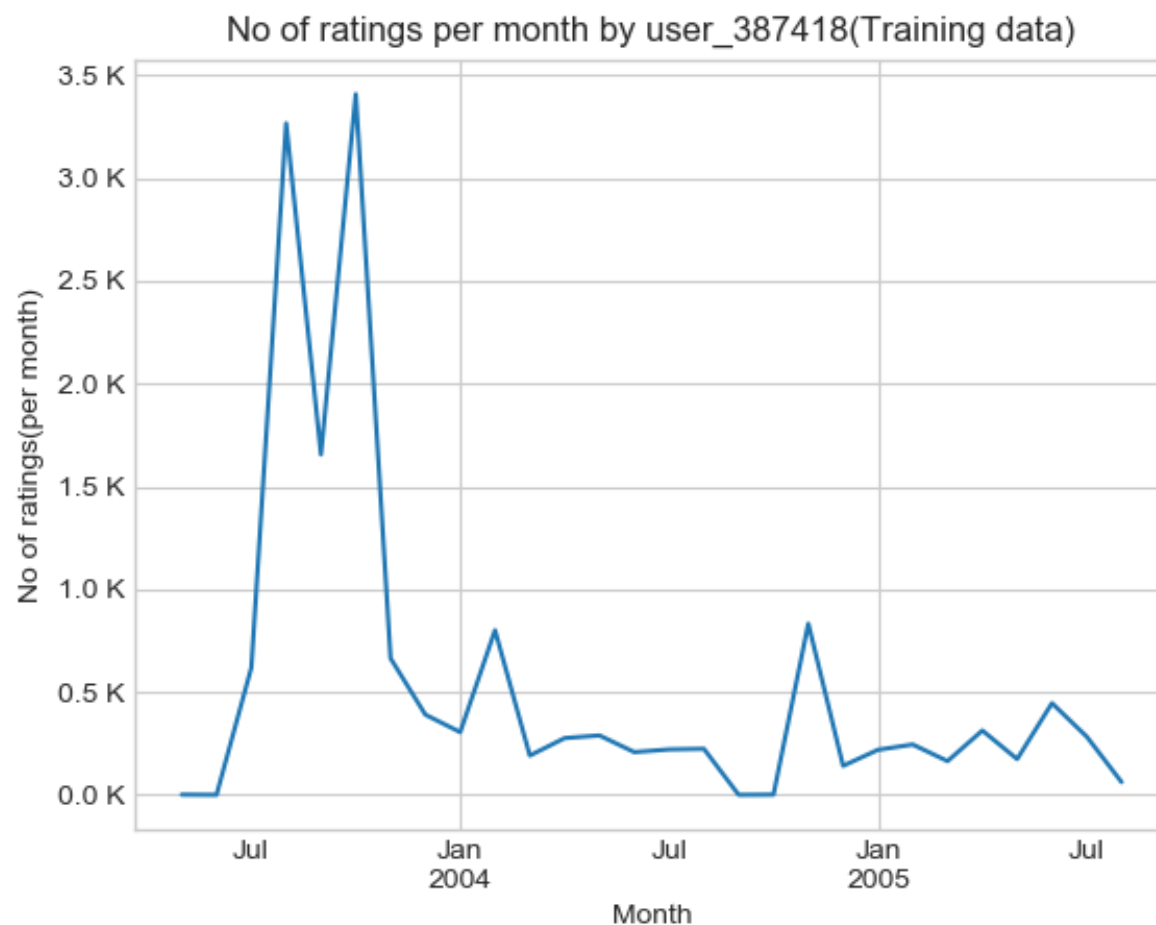


Observation:

- user_2439493 also rated movies heavily. So, there may be a chance that this user may have rated movies randomly.

```
In [45]: grouped = train_df.groupby('user')
user_387418 = grouped.get_group(387418)

ax = user_387418.resample('m', on='date')['rating'].count().plot()
ax.set_title('No of ratings per month by user_387418(Training data)')
plt.xlabel('Month')
plt.ylabel('No of ratings(per month)')
ax.set_yticklabels([human(item, 'k') for item in ax.get_yticks()])
plt.show()
```



Observation:

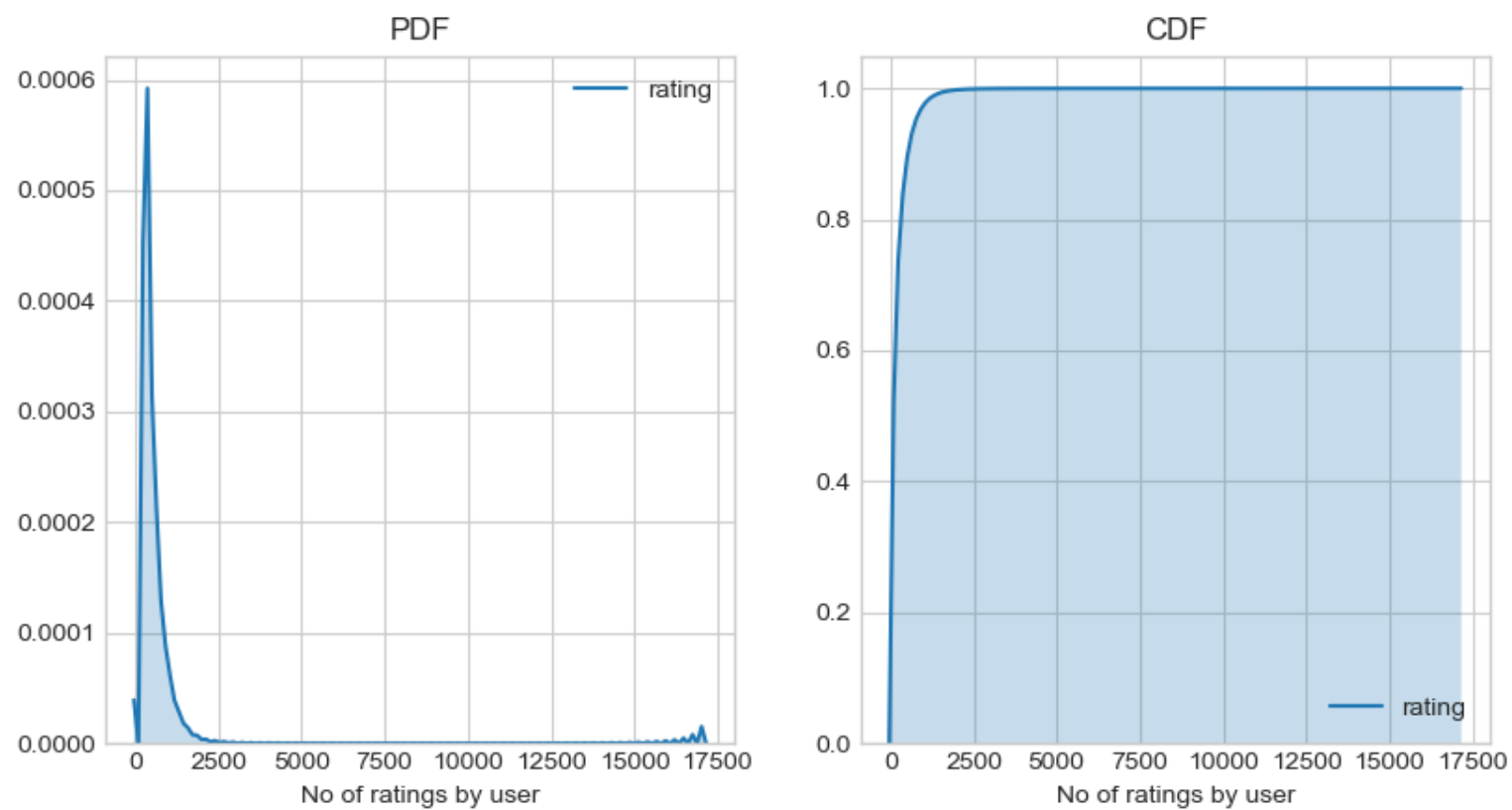
By seeing the above plot we can say that there is a chance that the rating given by top 3 users can be outliers because these users in the span of one year have watched approx 6k movies which is unconventional. Let's plot pdf and cdf if we can come with any insight.

```
In [20]: 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()
```



- looks like only few user have been rated 100s of movies.

```
In [21]: no_of Rated_movies_per_user.describe()
```

```
Out[21]: count      405041.000000  
         mean        198.459921  
         std         290.793238  
         min          1.000000  
         25%         34.000000  
         50%         89.000000  
         75%        245.000000  
         max        17112.000000  
         Name: rating, dtype: float64
```

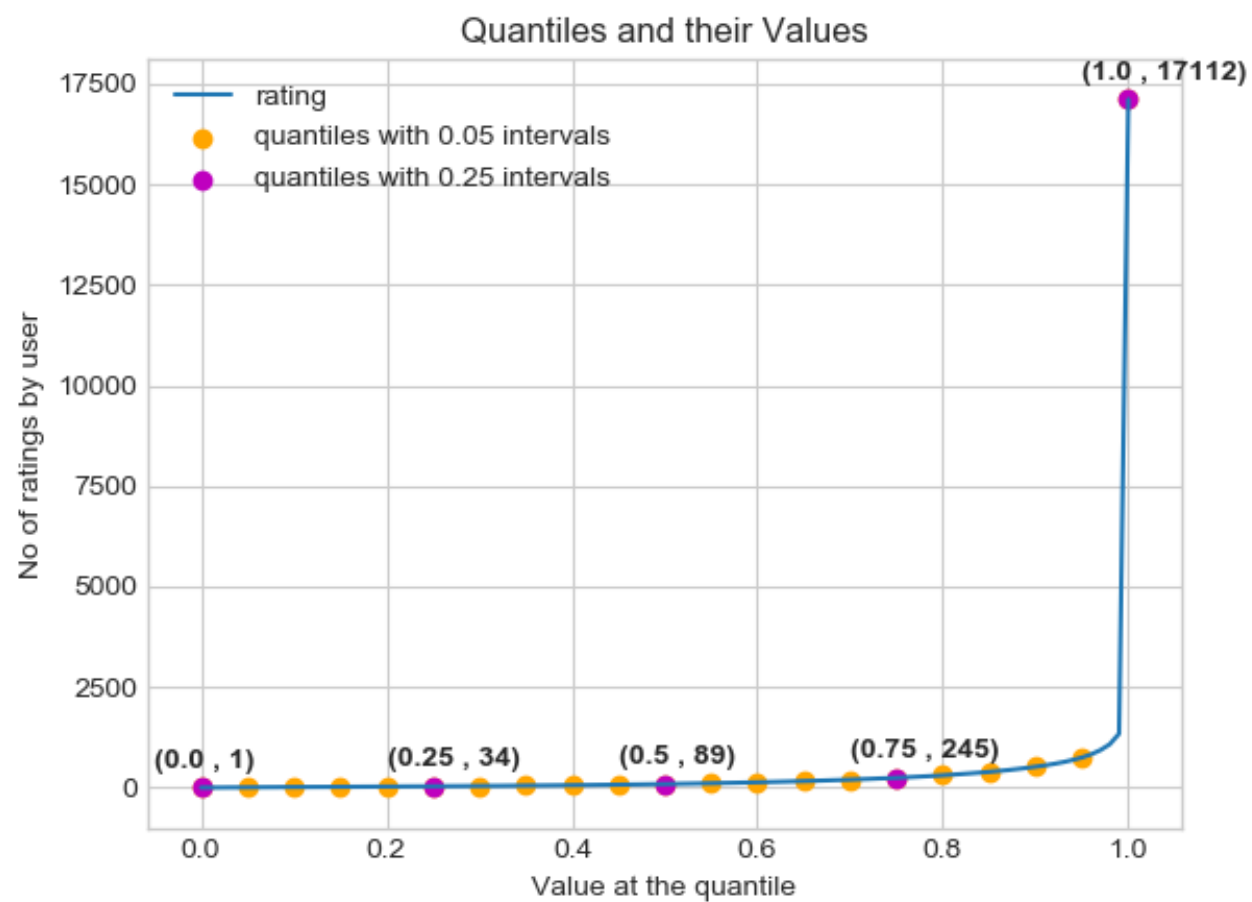
There, is something interesting going on with the quantiles..

```
In [22]: quantiles = no_of Rated_movies_per_user.quantile(np.arange(0,1.01,0.01), interpolation='higher')
```

```
In [23]: 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()
```

```
In [24]: quantiles[::5]
```

```
Out[24]: 0.00      1
          0.05      7
          0.10     15
          0.15     21
          0.20     27
          0.25     34
          0.30     41
          0.35     50
          0.40     60
          0.45     73
          0.50     89
          0.55    109
          0.60    133
          0.65    163
          0.70    199
          0.75    245
          0.80    307
          0.85    392
          0.90    520
          0.95    749
          1.00   17112
          Name: rating, dtype: int64
```

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

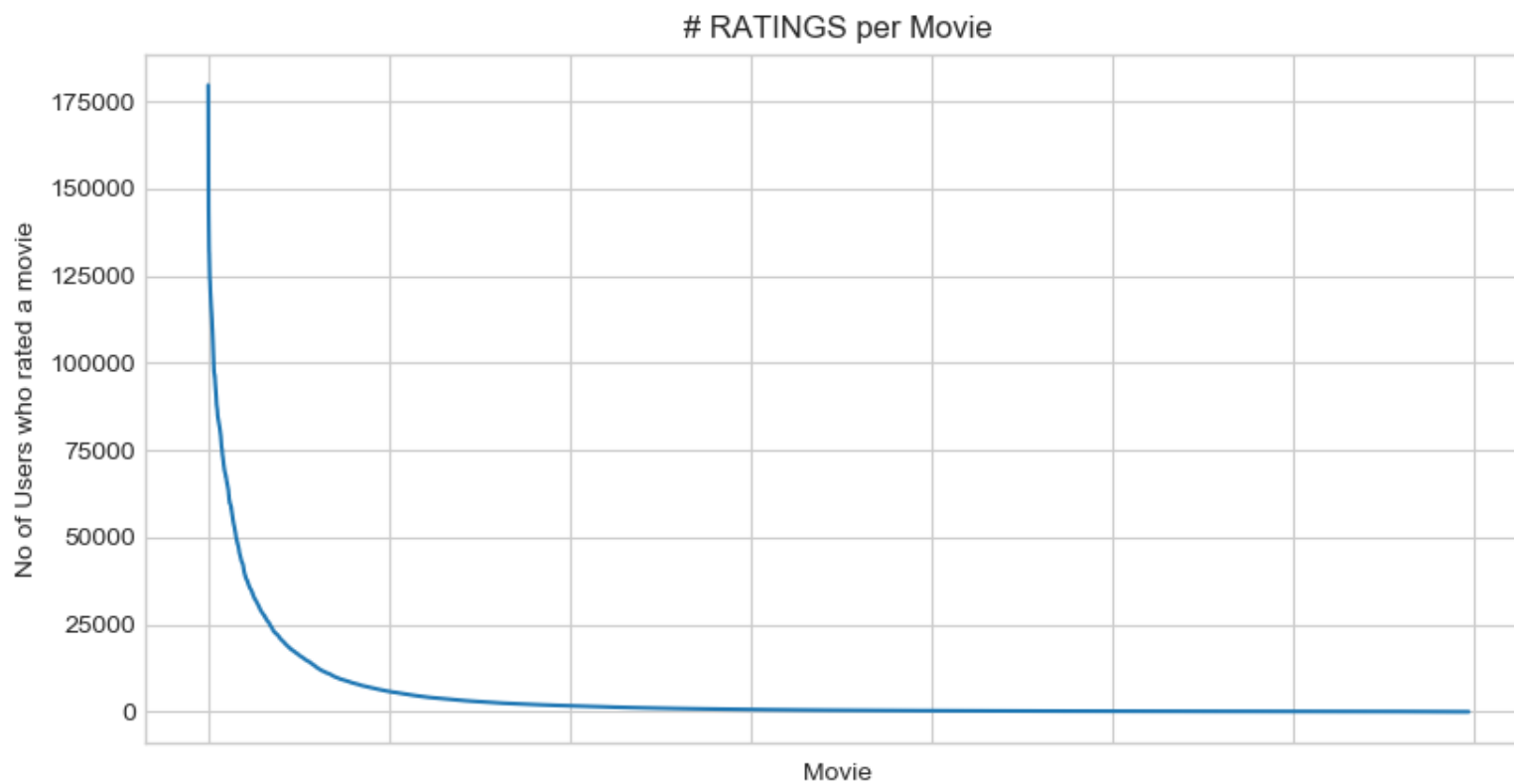
```
In [25]: print('\n No of ratings at last 5 percentile : {}'.format(sum(no_of Rated_movies_per_user>= 749)) )
          No of ratings at last 5 percentile : 20305
```

3.3.4 Analysis of ratings of a movie given by a user

```
In [26]: 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([])

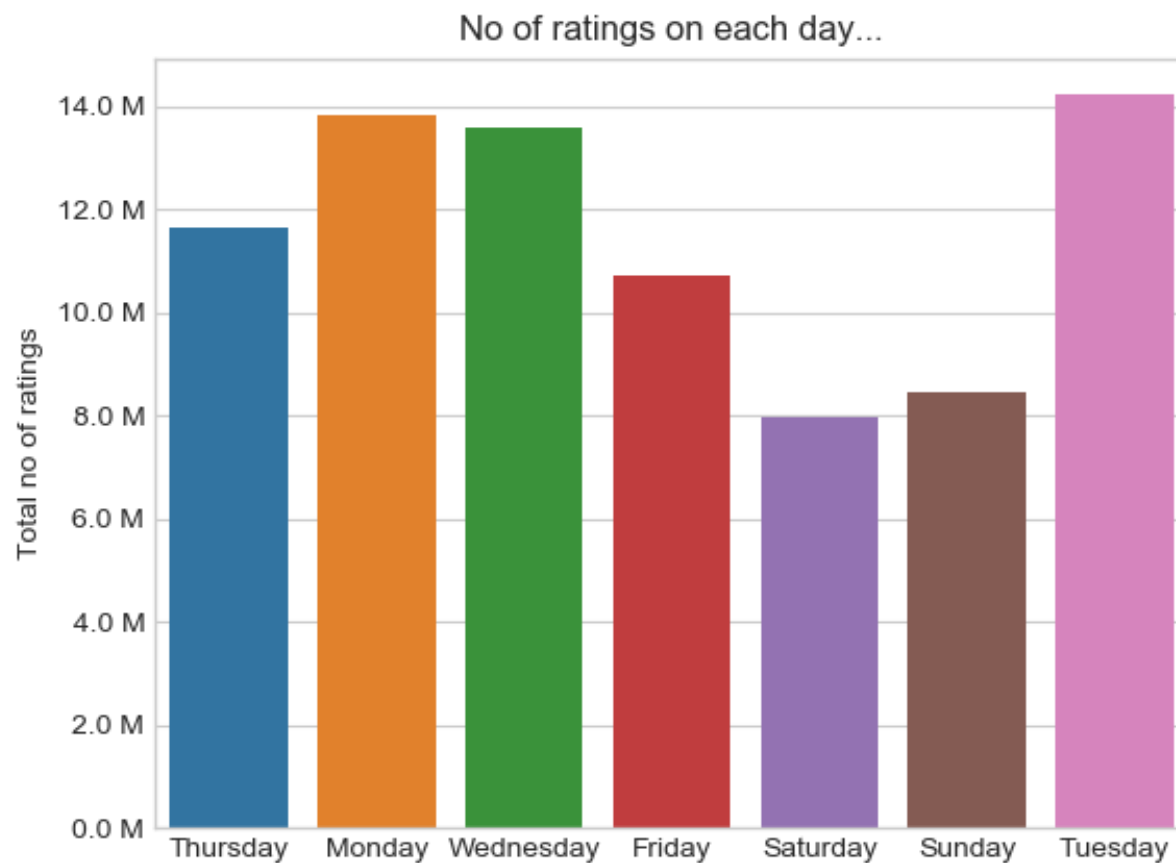
plt.show()
```



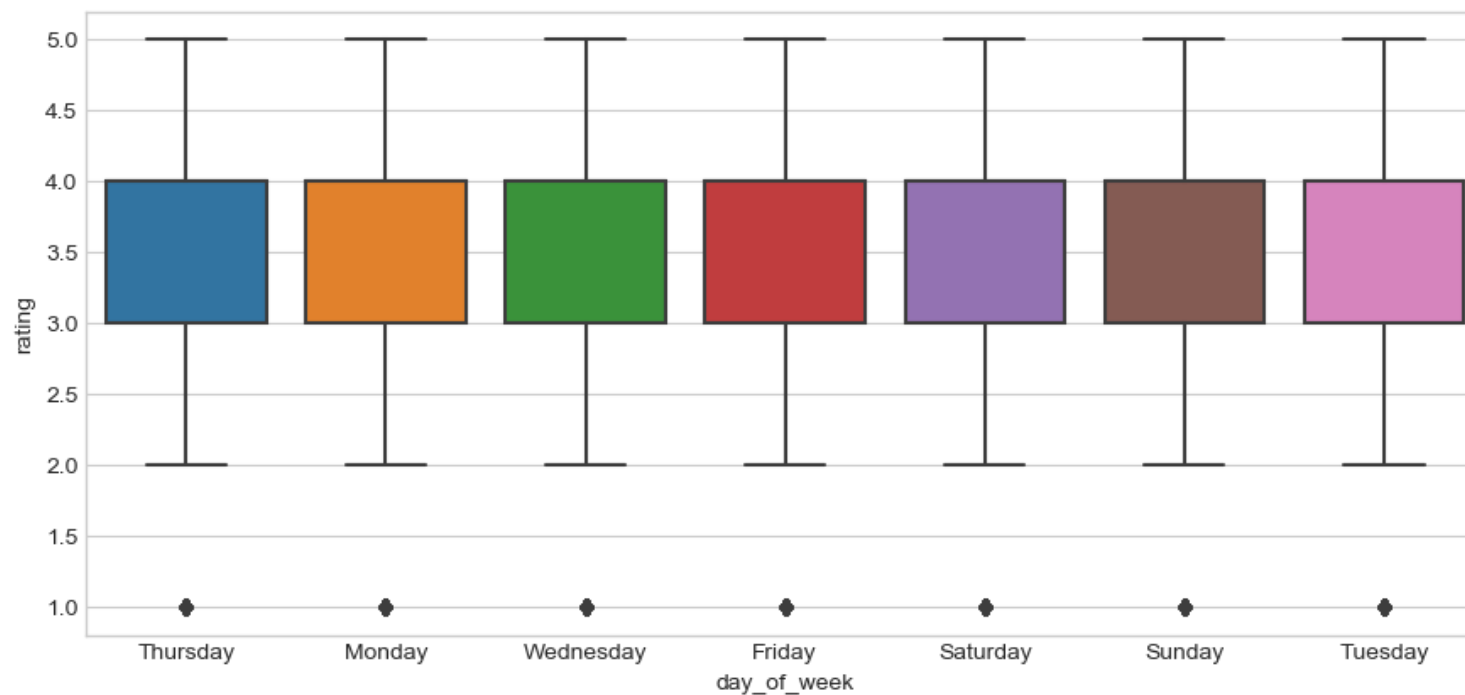
- **It is very skewed.. just like number of ratings given per user.**
 - There are some movies (which are very popular) which are rated by huge number of users.
 - But most of the movies(like 90%) got some hundreds of ratings.

3.3.5 Number of ratings on each day of the week

```
In [27]: 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 [28]: 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)
```



0:02:16.434910

```
In [29]: 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 3.585274

Monday 3.577250

Saturday 3.591791

Sunday 3.594144

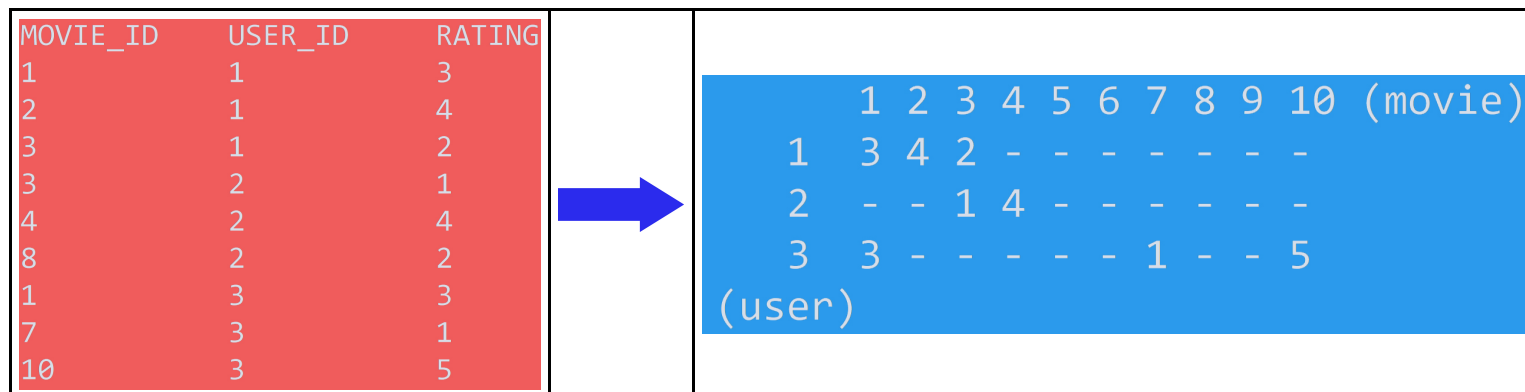
Thursday 3.582463

Tuesday 3.574438

Wednesday 3.583751

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 [27]: start = datetime.now()
if os.path.isfile('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('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("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..
0:00:12.156819
```

The Sparsity of Train Sparse Matrix


```
In [28]: 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 : 99.8292709259195 %
```

3.3.6.2 Creating sparse matrix from test data frame

```
In [29]: start = datetime.now()
         if os.path.isfile('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('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 further usage..')
             # save it into disk
             sparse.save_npz("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..
0:00:02.681418
```

The Sparsity of Test data Matrix

```
In [30]: 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 : 99.95731772988694 %
```

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

```
In [17]: # 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)
    is Rated = sparse_matrix!=0
    # no of ratings that each user OR movie..
    no_of_ratings = is Rated.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 [32]: 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[32]: {'global': 3.582890686321557}
```

3.3.7.2 finding average rating per user

```
In [33]: train_averages['user'] = get_average_ratings(train_sparse_matrix, of_users=True)
print('\nAverage rating of user 10 :',train_averages['user'][10])
```

```
Average rating of user 10 : 3.3781094527363185
```

3.3.7.3 finding average rating per movie

```
In [34]: 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 : 3.3038461538461537
```

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

```
In [35]: 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)
```

0:02:09.393514

3.3.8 Cold Start problem

3.3.8.1 Cold Start problem with Users

```
In [36]: 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 : 480189

Number of Users in Train data : 405041

No of Users that didn't appear in train data: 75148(15.65 %)

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 [37]: total_movies = len(np.unique(df.movie))
movies_train = len(train_averages['movie'])
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 : 17770

Number of Users in Train data : 17424

No of Movies that didn't appear in train data: 346(1.95 %)

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 [16]: **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 [39]: 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 : 0:03:40.371371]

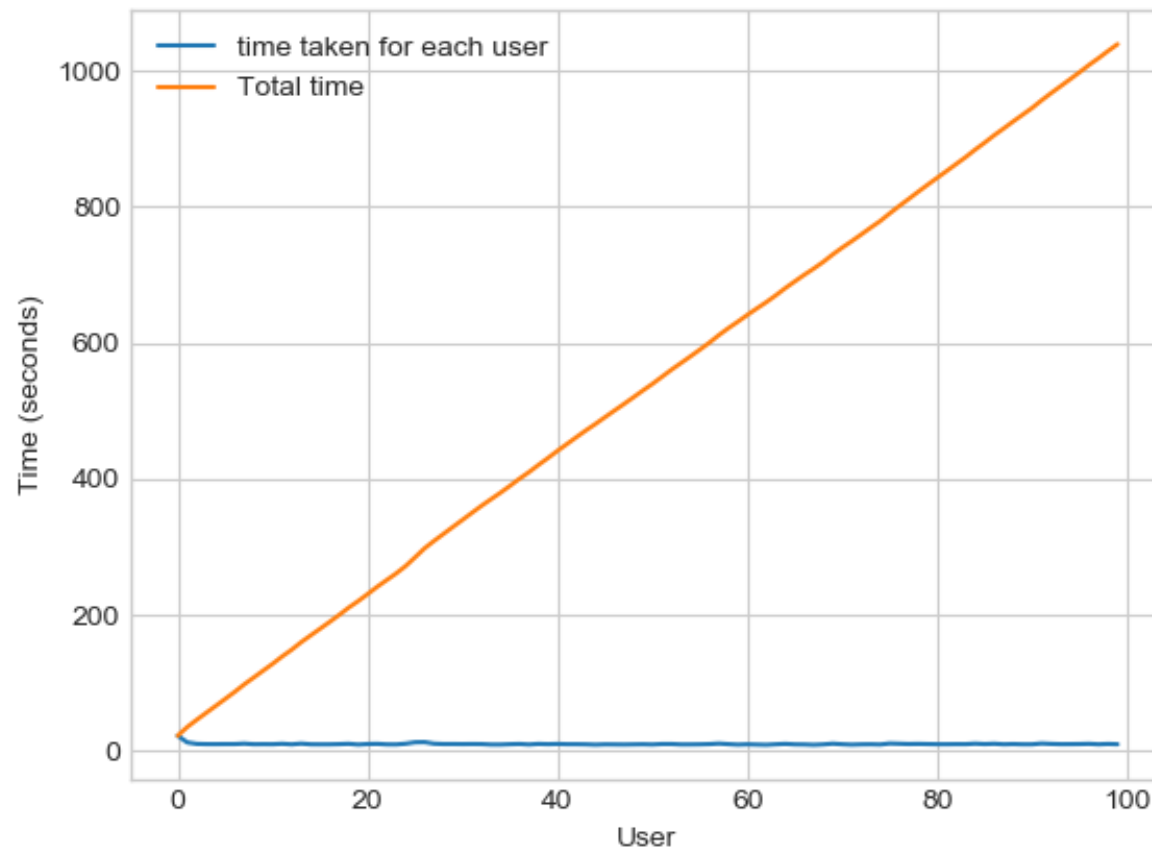
computing done for 40 users [time elapsed : 0:07:10.810329]

computing done for 60 users [time elapsed : 0:10:31.167230]

computing done for 80 users [time elapsed : 0:13:52.550700]

computing done for 100 users [time elapsed : 0:17:19.174500]

Creating Sparse matrix from the computed similarities



Time taken : 0:17:46.390183

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.
- $405041 \times 8.88 = 3596764.08 \text{ sec} = 59946.068 \text{ min} = 999.101133333 \text{ hours} = 41.629213889 \text{ days} \dots$
 - 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 [40]: train_sparse_matrix.shape
```

```
Out[40]: (2649430, 17771)
```

```
In [41]: 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=50, algorithm='randomized', random_state=15)
         trunc_svd = netflix_svd.fit_transform(train_sparse_matrix)

         print(datetime.now()-start)

0:03:59.891268
```

Here,

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

```
In [42]: expl_var = np.cumsum(netflix_svd.explained_variance_ratio_)
```

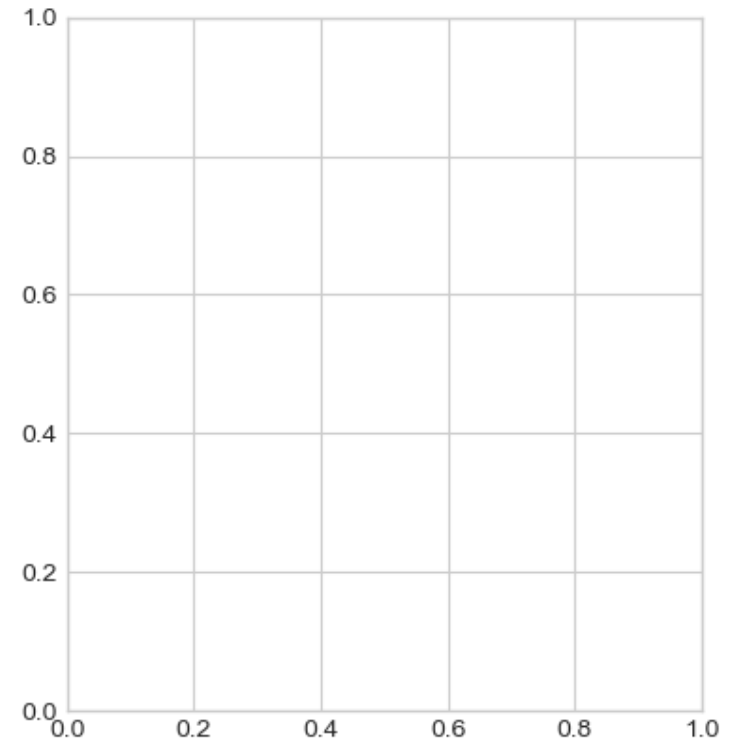
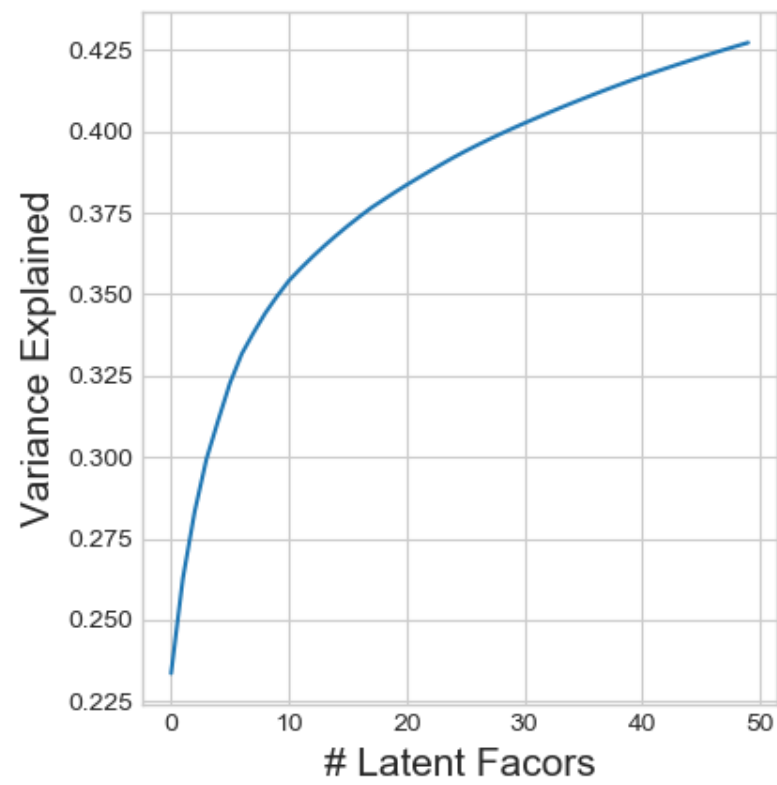
```
In [44]: fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))

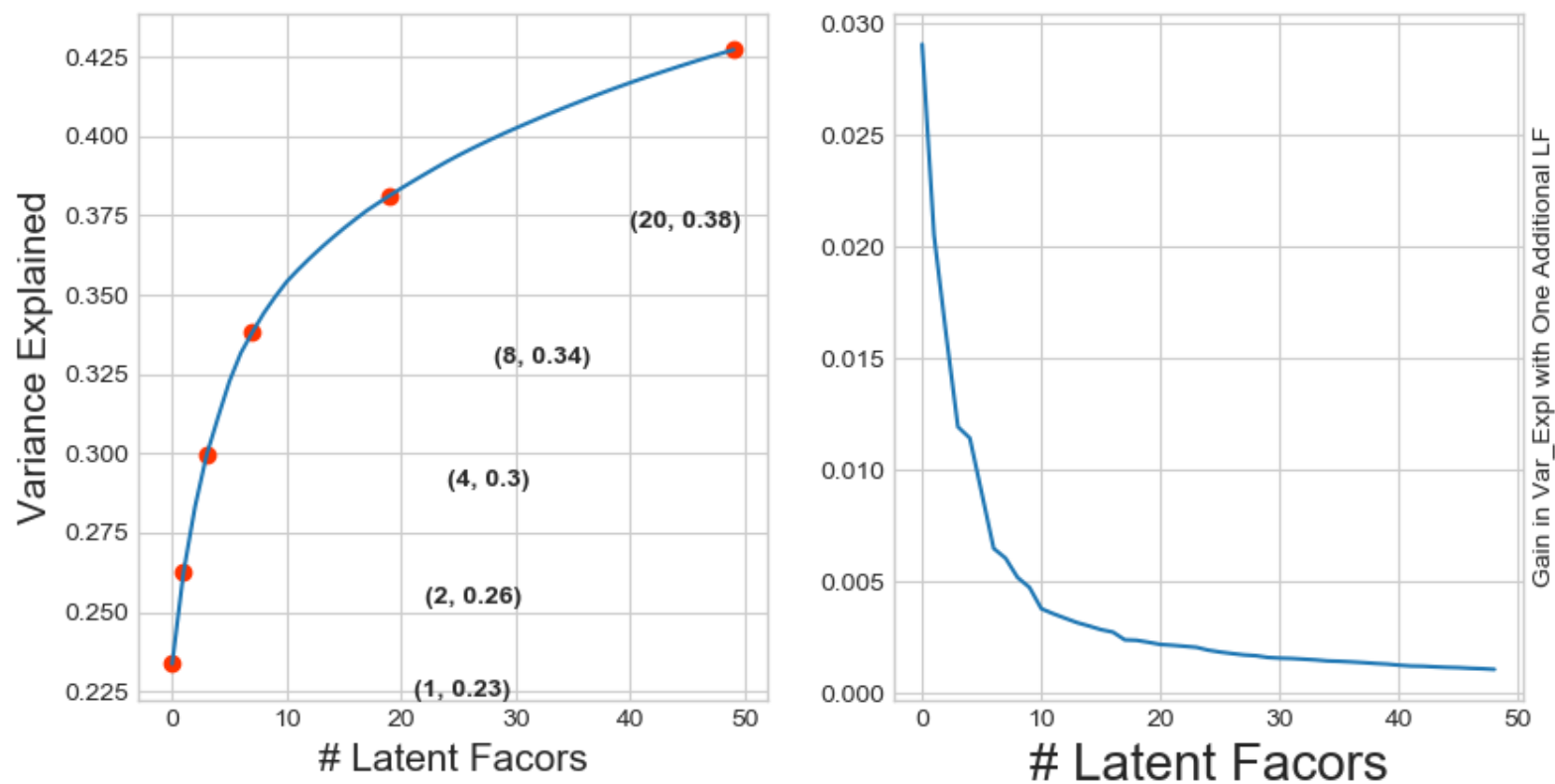
ax1.set_ylabel("Variance Explained", fontsize=15)
ax1.set_xlabel("# Latent Facors", fontsize=15)
ax1.plot(expl_var)
# annotate some (latentfactors, expl_var) to make it clear
ind = [1, 2,4,8,20, 50]
ax1.scatter(x = [i-1 for i in ind], y = expl_var[[i-1 for i in ind]], c='#ff3300')
for i in ind:
    ax1.annotate(s ="({}, {})".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 [45]: for i in ind:
          print("{} , {}".format(i, np.round(expl_var[i-1], 2)))
```

```
(1, 0.23)
(2, 0.26)
(4, 0.3)
(8, 0.34)
(20, 0.38)
(50, 0.43)
```


I think 500 dimensions is good enough

- 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 factore too it, the **_gain in expained variance** with that addition is decreasing. (Obviously, because they are sorted that way).
- **LHS Graph:**
 - **x** --- (No of latent factos),
 - **y** --- (The variance explained by taking x latent factors)
- **More decrease in the line (RHS graph) :**
 - We are getting more expained variance than before.
- **Less decrease in that line (RHS graph) :**
 - We are not getting benifitted from adding latent factor furthur. 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 [46]: # 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)

0:00:09.603880
```

```
In [47]: type(trunc_matrix), trunc_matrix.shape
```

```
Out[47]: (numpy.ndarray, (2649430, 50))
```

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

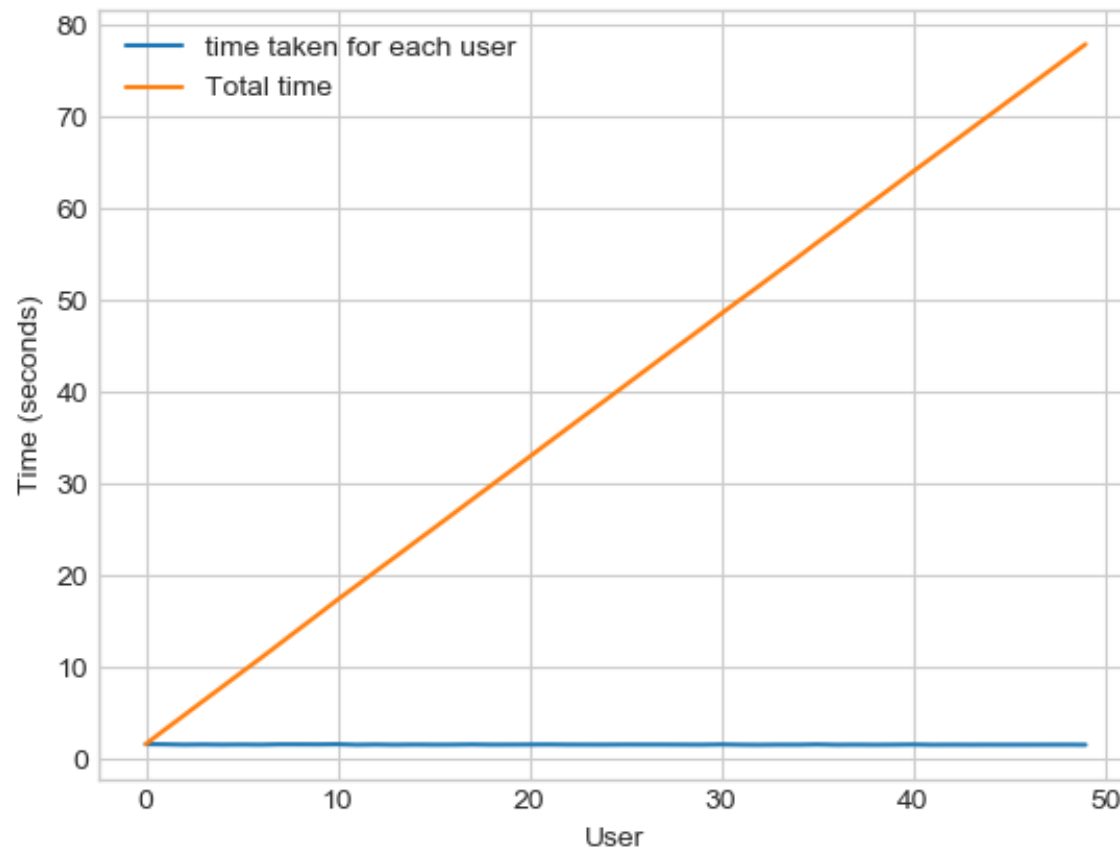
```
In [48]: if not os.path.isfile('trunc_sparse_matrix.npz'):
          # create that sparse matrix
          trunc_sparse_matrix = sparse.csr_matrix(trunc_matrix)
          # Save this truncated sparse matrix for later usage..
          sparse.save_npz('trunc_sparse_matrix', trunc_sparse_matrix)
        else:
          trunc_sparse_matrix = sparse.load_npz('trunc_sparse_matrix.npz')
```

```
In [49]: trunc_sparse_matrix.shape
```

```
Out[49]: (2649430, 50)
```

```
In [50]: 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 : 0:00:15.758731 ]  
computing done for 20 users [ time elapsed : 0:00:31.363822 ]  
computing done for 30 users [ time elapsed : 0:00:46.958216 ]  
computing done for 40 users [ time elapsed : 0:01:02.470898 ]  
computing done for 50 users [ time elapsed : 0:01:17.932740 ]  
Creating Sparse matrix from the computed similarities
```



time: 0:01:23.011698

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

- from above plot, It took almost **12.18** for computing simlilar users for **one user**
- We have **405041 users** with us in training set.
- $405041 \times 12.18 \text{ sec} \text{ ===== } 4933399.38 \text{ sec} \text{ ===== } 82223.323 \text{ min} \text{ ===== } 1370.388716667 \text{ hours} \text{ ===== } 57.099529861 \text{ days} \dots$
 - Even we run on 4 cores parallely (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 [51]: start = datetime.now()
if not os.path.isfile('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("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("m_m_sim_sparse.npz")
    print("Done ...")

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

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

```
It seems you don't have that file. Computing movie_movie similarity...
Done..
Saving it to disk without the need of re-computing it again..
Done..
It's a (17771, 17771) dimensional matrix
0:15:57.259760
```

```
In [52]: m_m_sim_sparse.shape
```

```
Out[52]: (17771, 17771)
```

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

```
In [54]: 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]
```

```
0:02:49.335626
```

```
Out[54]: array([ 8279,  8013, 16528,  5927, 13105, 12049,  4424, 10193, 17590,
                4549,  3755,   590, 14059, 15144, 15054,  9584,  9071,  6349,
               16402,  3973,  1720,  5370, 16309,  9376,  6116,  4706,  2818,
                 778, 15331,  1416, 12979, 17139, 17710,  5452,  2534,   164,
               15188,  8323,  2450, 16331,  9566, 15301, 13213, 14308, 15984,
               10597,  6426,  5500,  7068,  7328,  5720,  9802,   376, 13013,
                 8003, 10199,  3338, 15390,  9688, 16455, 11730,  4513,   598,
               12762,  2187,   509,  5865,  9166, 17115, 16334,  1942,  7282,
               17584,  4376,  8988,  8873,  5921,  2716, 14679, 11947, 11981,
                 4649,   565, 12954, 10788, 10220, 10963,  9427,  1690,  5107,
               7859,  5969,  1510,  2429,   847,  7845,  6410, 13931,  9840,
               3706], dtype=int64)
```


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 [55]: # First Let's load the movie details into soe dataframe..  
# movie details are in 'netflix/movie_titles.csv'  
  
movie_titles = pd.read_csv("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: 18.59 ms

Type conversion took: 500.36 ms

Parser memory cleanup took: 0.00 ms

Out[55]:

	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 [56]: 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 similarto this  and we will get only top most..".format(m_m_sim_sparse[:,mv_id].getnnz()))
```

Movie -----> Vampire Journals

It has 270 Ratings from users.

We have 17284 movies which are similarto this and we will get only top most..

```
In [57]: similarities = m_m_sim_sparse[mv_id].toarray().ravel()

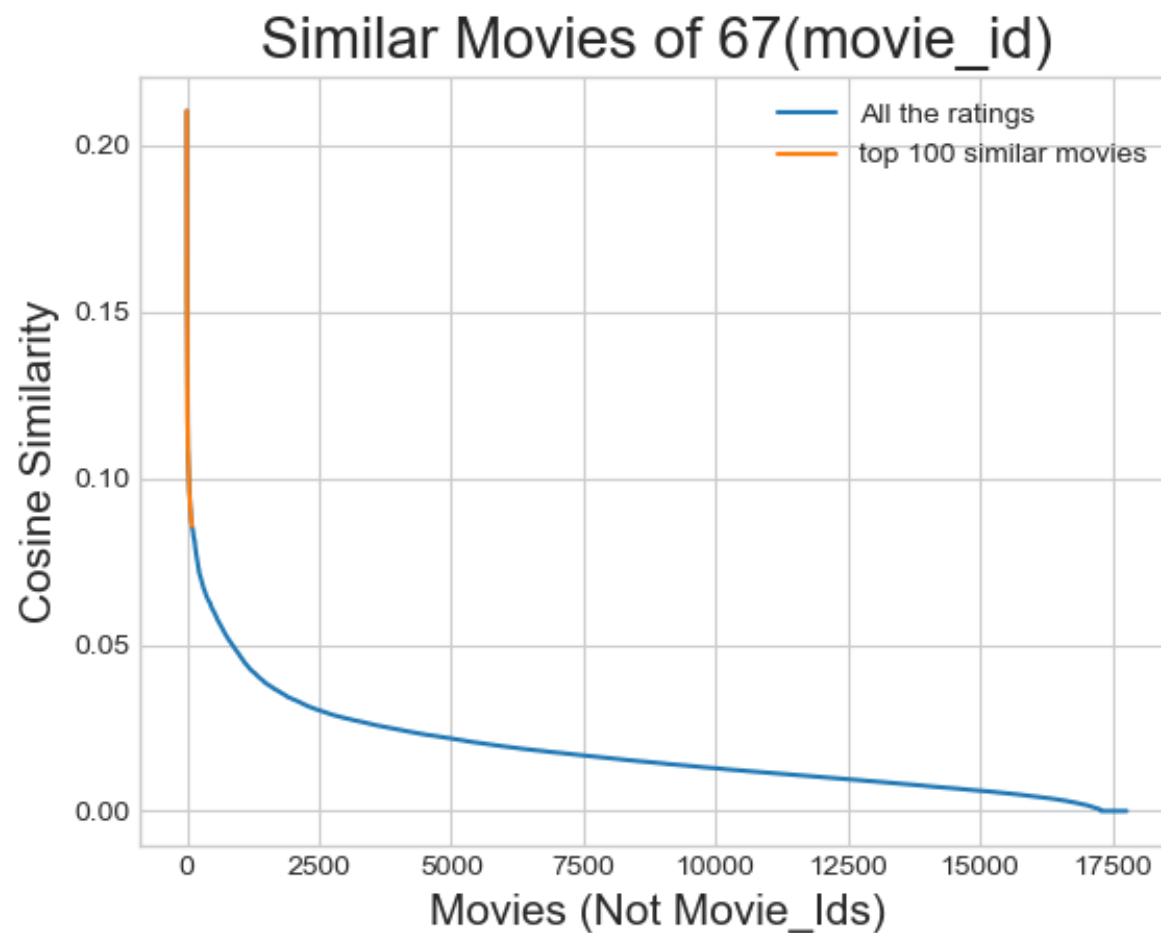
similar_indices = similarities.argsort()[::-1][1:]

similarities[similar_indices]

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

# and return its indices(movie_ids)
```

```
In [58]: 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 [59]: movie_titles.loc[sim_indices[:10]]
```

Out[59]:

	year_of_release	title
movie_id		
323	1999.0	Modern Vampires
4044	1998.0	Subspecies 4: Bloodstorm
1688	1993.0	To Sleep With a Vampire
13962	2001.0	Dracula: The Dark Prince
12053	1993.0	Dracula Rising
16279	2002.0	Vampires: Los Muertos
4667	1996.0	Vampirella
1900	1997.0	Club Vampire
13873	2001.0	The Breed
15867	2003.0	Dracula II: Ascension

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

4. Machine Learning Models



```

In [62]: 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("Origin
              al Matrix : Ratings -- {}".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, replace=False)
        sample_movies = np.random.choice(movies, no_movies, replace=False)
        # 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 [3]: start = datetime.now()
path = "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 12k users and 1.2k movies from available data
    sample_train_sparse_matrix = get_sample_sparse_matrix(train_sparse_matrix, no_users=12000, no_movies=1200,
                                                         path = path)

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

It is present in your pwd, getting it from disk....
DONE..
0:00:00.438861

4.1.2 Build sample test data from the test data

```
In [4]: start = datetime.now()

path = "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 500 movies from available data
    sample_test_sparse_matrix = get_sample_sparse_matrix(test_sparse_matrix, no_users=6000, no_movies=600,
                                                         path = path)

print(datetime.now() - start)

It is present in your pwd, getting it from disk....
DONE..
0:00:00.209952
```

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

```
In [5]: sample_train_averages = dict()
```

4.2.1 Finding Global Average of all movie ratings

```
In [6]: # 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
```

```
Out[6]: {'global': 3.5506507537264302}
```


4.2.2 Finding Average rating per User

```
In [19]: sample_train_averages['user'] = get_average_ratings(sample_train_sparse_matrix, of_users=True)
print('\nAverage rating of user 1515220 : ',sample_train_averages['user'][1515220])
```

Average rating of user 1515220 : 3.8857142857142857

4.2.3 Finding Average rating per Movie

```
In [20]: sample_train_averages['movie'] = get_average_ratings(sample_train_sparse_matrix, of_users=False)
print('\n AVerage rating of movie 15153 : ',sample_train_averages['movie'][15153])
```

AVerage rating of movie 15153 : 2.603448275862069

4.3 Featurizing data

```
In [21]: 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 : 177717

No of ratings in Our Sampled test matrix is : 11207

4.3.1 Featurizing data for regression problem

4.3.1.1 Featurizing train data

```
In [22]: # 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 [24]: #####
# It took me almost 10 hours to prepare this train dataset.#
#####
start = datetime.now()
import tqdm
if os.path.isfile('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('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):
            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 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=" ")

            #----- 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).r
            avel()
            top_sim_movies = movie_sim.argsort()[::-1][1:] # we are ignoring 'The User' from its 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, 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)%1000 == 0:
    # print(','.join(map(str, row)))
    print("Done for {} rows----- {}".format(count, datetime.now() - start))

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

preparing 177717 tuples for the dataset..

Done for 1000 rows----- 0:15:51.333448
Done for 2000 rows----- 0:29:47.709937
Done for 3000 rows----- 0:48:33.199195
Done for 4000 rows----- 1:01:27.725739
Done for 5000 rows----- 1:14:55.665047
Done for 6000 rows----- 1:28:00.988104
Done for 7000 rows----- 1:40:58.425156
Done for 8000 rows----- 1:53:59.651670
Done for 9000 rows----- 2:06:51.752968
Done for 10000 rows----- 2:19:47.660780
Done for 11000 rows----- 2:32:38.238940
Done for 12000 rows----- 2:45:39.208825
Done for 13000 rows----- 2:58:37.654528
Done for 14000 rows----- 3:11:31.053123
Done for 15000 rows----- 3:24:32.750318
Done for 16000 rows----- 3:37:33.222882
Done for 17000 rows----- 3:50:26.922369
Done for 18000 rows----- 4:03:25.262373
Done for 19000 rows----- 4:16:23.045653
Done for 20000 rows----- 4:29:19.174863
Done for 21000 rows----- 4:42:12.674314
Done for 22000 rows----- 4:55:13.308476
Done for 23000 rows----- 5:07:22.048971
Done for 24000 rows----- 5:19:25.435250
Done for 25000 rows----- 5:31:19.005603
Done for 26000 rows----- 5:43:20.817872
Done for 27000 rows----- 5:55:25.292839
Done for 28000 rows----- 6:07:30.290312
Done for 29000 rows----- 6:19:27.018334
Done for 30000 rows----- 6:31:19.830067
Done for 31000 rows----- 6:43:11.876775
Done for 32000 rows----- 6:55:06.587002
Done for 33000 rows----- 7:06:58.309552
Done for 34000 rows----- 7:18:51.481510
Done for 35000 rows----- 7:30:48.704976
Done for 36000 rows----- 7:42:47.704534
Done for 37000 rows----- 7:54:32.548470
Done for 38000 rows----- 8:06:14.982600

Done for 39000 rows----- 8:17:57.930539
Done for 40000 rows----- 8:29:51.736873
Done for 41000 rows----- 8:41:46.602546
Done for 42000 rows----- 8:53:31.511222
Done for 43000 rows----- 9:05:28.783986
Done for 44000 rows----- 9:17:29.853104
Done for 45000 rows----- 9:29:21.459987
Done for 46000 rows----- 9:41:15.880478
Done for 47000 rows----- 9:53:16.086942
Done for 48000 rows----- 10:05:04.990314
Done for 49000 rows----- 10:16:48.717979
Done for 50000 rows----- 10:28:52.160434
Done for 51000 rows----- 10:40:42.928957
Done for 52000 rows----- 10:52:44.318534
Done for 53000 rows----- 11:04:41.916312
Done for 54000 rows----- 11:16:23.674384
Done for 55000 rows----- 11:29:13.092288
Done for 56000 rows----- 11:42:12.394283
Done for 57000 rows----- 11:54:26.518780
Done for 58000 rows----- 12:09:47.767947
Done for 59000 rows----- 12:24:16.051351
Done for 60000 rows----- 12:42:38.662797
Done for 61000 rows----- 13:02:30.434092
Done for 62000 rows----- 13:14:52.922798
Done for 63000 rows----- 13:27:08.531276
Done for 64000 rows----- 13:39:11.569890
Done for 65000 rows----- 13:51:16.200141
Done for 66000 rows----- 14:03:19.868539
Done for 67000 rows----- 14:15:20.812690
Done for 68000 rows----- 14:27:43.969466
Done for 69000 rows----- 14:39:55.661146
Done for 70000 rows----- 14:52:29.932425
Done for 71000 rows----- 15:04:46.981082
Done for 72000 rows----- 15:17:06.042697
Done for 73000 rows----- 15:29:17.273240
Done for 74000 rows----- 15:41:17.789290
Done for 75000 rows----- 15:53:20.978416
Done for 76000 rows----- 16:05:22.451864
Done for 77000 rows----- 16:17:23.220427
Done for 78000 rows----- 16:29:26.854123

Done for 79000 rows----- 16:41:32.053351
Done for 80000 rows----- 16:53:36.683423
Done for 81000 rows----- 17:05:38.641272
Done for 82000 rows----- 17:17:40.989622
Done for 83000 rows----- 17:29:44.989306
Done for 84000 rows----- 17:41:46.511673
Done for 85000 rows----- 17:53:51.188719
Done for 86000 rows----- 18:05:52.067920
Done for 87000 rows----- 18:17:55.536395
Done for 88000 rows----- 18:29:58.916467
Done for 89000 rows----- 18:42:02.466849
Done for 90000 rows----- 18:54:10.110314
Done for 91000 rows----- 19:06:11.207587
Done for 92000 rows----- 19:18:14.453771
Done for 93000 rows----- 19:30:26.785359
Done for 94000 rows----- 19:42:47.078845
Done for 95000 rows----- 19:54:51.578216
Done for 96000 rows----- 20:07:02.462514
Done for 97000 rows----- 20:19:45.420329
Done for 98000 rows----- 20:32:13.898403
Done for 99000 rows----- 20:44:26.768651
Done for 100000 rows----- 20:56:49.705882
Done for 101000 rows----- 21:09:11.976479
Done for 102000 rows----- 21:21:20.939799
Done for 103000 rows----- 21:33:29.890246
Done for 104000 rows----- 21:45:43.317010
Done for 105000 rows----- 21:58:08.672618
Done for 106000 rows----- 22:10:26.801945
Done for 107000 rows----- 22:22:53.829969
Done for 108000 rows----- 22:38:38.912940
Done for 109000 rows----- 23:10:14.381908
Done for 110000 rows----- 23:27:59.320883
Done for 111000 rows----- 23:40:05.613654
Done for 112000 rows----- 23:52:19.547148
Done for 113000 rows----- 1 day, 0:04:35.789066
Done for 114000 rows----- 1 day, 0:16:46.709465
Done for 115000 rows----- 1 day, 0:29:10.186809
Done for 116000 rows----- 1 day, 0:42:24.475316
Done for 117000 rows----- 1 day, 0:55:05.888407
Done for 118000 rows----- 1 day, 1:07:29.245208

Done for 119000 rows----- 1 day, 1:19:46.067630
Done for 120000 rows----- 1 day, 1:32:08.220489
Done for 121000 rows----- 1 day, 1:44:26.599146
Done for 122000 rows----- 1 day, 1:58:42.831358
Done for 123000 rows----- 1 day, 2:11:10.041136
Done for 124000 rows----- 1 day, 2:23:03.803590
Done for 125000 rows----- 1 day, 2:35:30.102207
Done for 126000 rows----- 1 day, 2:47:24.984462
Done for 127000 rows----- 1 day, 2:59:37.147646
Done for 128000 rows----- 1 day, 3:11:27.375117
Done for 129000 rows----- 1 day, 3:23:21.371363
Done for 130000 rows----- 1 day, 3:35:31.839961
Done for 131000 rows----- 1 day, 3:47:46.006450
Done for 132000 rows----- 1 day, 3:59:50.164553
Done for 133000 rows----- 1 day, 4:12:06.119952
Done for 134000 rows----- 1 day, 4:24:06.921108
Done for 135000 rows----- 1 day, 4:36:00.397252
Done for 136000 rows----- 1 day, 4:47:51.647379
Done for 137000 rows----- 1 day, 4:59:42.988160
Done for 138000 rows----- 1 day, 5:11:34.508447
Done for 139000 rows----- 1 day, 5:23:29.527111
Done for 140000 rows----- 1 day, 5:36:07.064129
Done for 141000 rows----- 1 day, 5:48:03.229229
Done for 142000 rows----- 1 day, 6:01:16.474934
Done for 143000 rows----- 1 day, 6:18:37.308691
Done for 144000 rows----- 1 day, 6:37:05.863650
Done for 145000 rows----- 1 day, 6:55:32.586611
Done for 146000 rows----- 1 day, 7:14:25.859641
Done for 147000 rows----- 1 day, 7:26:48.561104
Done for 148000 rows----- 1 day, 7:39:08.616724
Done for 149000 rows----- 1 day, 7:50:58.617057
Done for 150000 rows----- 1 day, 8:02:47.780931
Done for 151000 rows----- 1 day, 8:14:35.722284
Done for 152000 rows----- 1 day, 8:26:56.229520
Done for 153000 rows----- 1 day, 8:39:17.297567
Done for 154000 rows----- 1 day, 8:51:38.905187
Done for 155000 rows----- 1 day, 9:03:26.258986
Done for 156000 rows----- 1 day, 9:15:15.497692
Done for 157000 rows----- 1 day, 9:27:07.263089
Done for 158000 rows----- 1 day, 9:38:54.684104


```

Done for 159000 rows----- 1 day, 9:50:42.043199
Done for 160000 rows----- 1 day, 10:02:30.449196
Done for 161000 rows----- 1 day, 10:14:51.143389
Done for 162000 rows----- 1 day, 10:26:50.702010
Done for 163000 rows----- 1 day, 10:39:02.685127
Done for 164000 rows----- 1 day, 10:50:51.308388
Done for 165000 rows----- 1 day, 11:02:40.607527
Done for 166000 rows----- 1 day, 11:14:51.167546
Done for 167000 rows----- 1 day, 11:26:58.586036
Done for 168000 rows----- 1 day, 11:39:01.515987
Done for 169000 rows----- 1 day, 11:50:49.141379
Done for 170000 rows----- 1 day, 12:02:35.543089
Done for 171000 rows----- 1 day, 12:14:23.110780
Done for 172000 rows----- 1 day, 12:26:15.192696
Done for 173000 rows----- 1 day, 12:38:05.473253
Done for 174000 rows----- 1 day, 12:49:51.449148
Done for 175000 rows----- 1 day, 13:01:40.016369
Done for 176000 rows----- 1 day, 13:13:27.561417
Done for 177000 rows----- 1 day, 13:25:16.990924
1 day, 13:33:53.778751

```

Reading from the file to make a Train_dataframe

```

In [25]: reg_train = pd.read_csv('reg_train.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5', 'smr1', 'smr2', 'smr3', 'smr4', 'smr5', 'UAvg', 'MAvg', 'rating'], header=None)
reg_train.head()

```

Out[25]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating
0	53406	33	3.550651	4.0	5.0	5.0	4.0	1.0	5.0	2.0	5.0	3.0	1.0	3.366667	4.172414	4
1	67390	33	3.550651	1.0	5.0	4.0	5.0	3.0	4.0	4.0	3.0	4.0	2.0	3.769231	4.172414	4
2	99540	33	3.550651	4.0	5.0	5.0	5.0	4.0	3.0	4.0	5.0	4.0	4.0	3.300000	4.172414	3
3	99865	33	3.550651	5.0	5.0	4.0	4.0	4.0	5.0	4.0	4.0	5.0	3.0	3.639344	4.172414	5
4	101620	33	3.550651	2.0	3.0	5.0	5.0	4.0	4.0	3.0	3.0	5.0	5.0	3.571429	4.172414	5

```
In [26]: reg_train.shape
```

```
Out[26]: (177717, 16)
```

-
- **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 [27]: # 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 [28]: sample_train_averages['global']
```

```
Out[28]: 3.5506507537264302
```

```

In [29]: start = datetime.now()

if os.path.isfile('reg_test.csv'):
    print("It is already created...")
else:

    print('preparing {} tuples for the dataset..\n'.format(len(sample_test_ratings)))
    with open('reg_test.csv', 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], sample_train_sparse_matrix).ravel()

                top_sim_users = user_sim.argsort()[::-1][1:] # we are ignoring 'The User' from its similar user

                # 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 similar 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 its similar use
rs.

    # 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_rating
s)))

    #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)
# 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)

```

preparing 11207 tuples for the dataset..

```

Done for 1000 rows----- 0:11:40.503526
Done for 2000 rows----- 0:23:20.362292
Done for 3000 rows----- 0:34:57.996210
Done for 4000 rows----- 0:46:35.723121
Done for 5000 rows----- 0:58:13.007265
Done for 6000 rows----- 1:09:48.416678
Done for 7000 rows----- 1:21:25.306104
Done for 8000 rows----- 1:33:02.188226
Done for 9000 rows----- 1:44:40.623882
Done for 10000 rows----- 1:56:23.158404
Done for 11000 rows----- 2:07:59.808997
2:10:23.638319

```

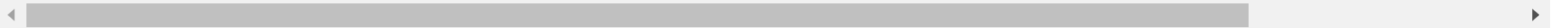
Reading from the file to make a test dataframe

```
In [30]: reg_test_df = pd.read_csv('reg_test.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
                    'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
                    'UAvg', 'MAvg', 'rating'], header=None)

reg_test_df.head(4)
```

Out[30]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5
0	808635	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651
1	941866	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651
2	1280761	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651
3	1737912	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651



-
- **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 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 [34]: !pip install scikit-surprise
```


Collecting scikit-surprise

Using cached <https://files.pythonhosted.org/packages/4d/fc/cd4210b247d1dca421c25994740cbbf03c5e980e31881f10eadf45fdab0/scikit-surprise-1.0.6.tar.gz>

Requirement already satisfied: joblib>=0.11 in c:\users\user\anaconda3\lib\site-packages (from scikit-surprise) (0.13.2)

Requirement already satisfied: numpy>=1.11.2 in c:\users\user\anaconda3\lib\site-packages (from scikit-surprise) (1.14.3)

Requirement already satisfied: scipy>=1.0.0 in c:\users\user\anaconda3\lib\site-packages (from scikit-surprise) (1.0.0)

Requirement already satisfied: six>=1.10.0 in c:\users\user\anaconda3\lib\site-packages (from scikit-surprise) (1.11.0)

Building wheels for collected packages: scikit-surprise

Running setup.py bdist_wheel for scikit-surprise: started

Running setup.py bdist_wheel for scikit-surprise: finished with status 'error'

Complete output from command c:\users\user\anaconda3\python.exe -u -c "import setuptools, tokenize;__file__='C:\\Users\\USER\\AppData\\Local\\Temp\\pip-install-aulpkruel\\scikit-surprise\\setup.py';f=getattr(tokenize, 'open', open)(__file__);code=f.read().replace('\r\n', '\n');f.close();exec(compile(code, __file__, 'exec'))" bdist_wheel -d C:\Users\USER\AppData\Local\Temp\pip-wheel-lxqcwfb --python-tag cp36:

running bdist_wheel

running build

running build_py

creating build

creating build\lib.win-amd64-3.6

creating build\lib.win-amd64-3.6\surprise

copying surprise\accuracy.py -> build\lib.win-amd64-3.6\surprise

copying surprise\builtin_datasets.py -> build\lib.win-amd64-3.6\surprise

copying surprise\dataset.py -> build\lib.win-amd64-3.6\surprise

copying surprise\dump.py -> build\lib.win-amd64-3.6\surprise

copying surprise\evaluate.py -> build\lib.win-amd64-3.6\surprise

copying surprise\reader.py -> build\lib.win-amd64-3.6\surprise

copying surprise\trainset.py -> build\lib.win-amd64-3.6\surprise

copying surprise\utils.py -> build\lib.win-amd64-3.6\surprise

copying surprise__init__.py -> build\lib.win-amd64-3.6\surprise

copying surprise__main__.py -> build\lib.win-amd64-3.6\surprise

creating build\lib.win-amd64-3.6\surprise\model_selection

copying surprise\model_selection\search.py -> build\lib.win-amd64-3.6\surprise\model_selection

copying surprise\model_selection\split.py -> build\lib.win-amd64-3.6\surprise\model_selection

copying surprise\model_selection\validation.py -> build\lib.win-amd64-3.6\surprise\model_selection

copying surprise\model_selection__init__.py -> build\lib.win-amd64-3.6\surprise\model_selection

creating build\lib.win-amd64-3.6\surprise\prediction_algorithms

```
copying surprise\prediction_algorithms\algo_base.py -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
copying surprise\prediction_algorithms\baseline_only.py -> build\lib.win-amd64-3.6\surprise\prediction_algorit
hms
copying surprise\prediction_algorithms\knns.py -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
copying surprise\prediction_algorithms\predictions.py -> build\lib.win-amd64-3.6\surprise\prediction_algorith
s
copying surprise\prediction_algorithms\random_pred.py -> build\lib.win-amd64-3.6\surprise\prediction_algorith
s
copying surprise\prediction_algorithms\__init__.py -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
running egg_info
writing scikit_surprise.egg-info\PKG-INFO
writing dependency_links to scikit_surprise.egg-info\dependency_links.txt
writing entry points to scikit_surprise.egg-info\entry_points.txt
writing requirements to scikit_surprise.egg-info\requires.txt
writing top-level names to scikit_surprise.egg-info\top_level.txt
reading manifest file 'scikit_surprise.egg-info\SOURCES.txt'
reading manifest template 'MANIFEST.in'
writing manifest file 'scikit_surprise.egg-info\SOURCES.txt'
copying surprise\similarities.c -> build\lib.win-amd64-3.6\surprise
copying surprise\similarities.pyx -> build\lib.win-amd64-3.6\surprise
copying surprise\prediction_algorithms\co_clustering.c -> build\lib.win-amd64-3.6\surprise\prediction_algorith
ms
copying surprise\prediction_algorithms\matrix_factorization.c -> build\lib.win-amd64-3.6\surprise\prediction_a
lgorithms
copying surprise\prediction_algorithms\optimize_baselines.c -> build\lib.win-amd64-3.6\surprise\prediction_alg
orithms
copying surprise\prediction_algorithms\slope_one.c -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
copying surprise\prediction_algorithms\co_clustering.pyx -> build\lib.win-amd64-3.6\surprise\prediction_algori
thms
copying surprise\prediction_algorithms\matrix_factorization.pyx -> build\lib.win-amd64-3.6\surprise\prediction
_algorithms
copying surprise\prediction_algorithms\optimize_baselines.pyx -> build\lib.win-amd64-3.6\surprise\prediction_a
lgorithms
copying surprise\prediction_algorithms\slope_one.pyx -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
running build_ext
building 'surprise.similarities' extension
error: Microsoft Visual C++ 14.0 is required. Get it with "Microsoft Visual C++ Build Tools": http://landinghub.visualstudio.com/visual-cpp-build-tools
```

```

Running setup.py clean for scikit-surprise
Failed to build scikit-surprise
Installing collected packages: scikit-surprise
Running setup.py install for scikit-surprise: started
Running setup.py install for scikit-surprise: finished with status 'error'
Complete output from command c:\users\user\anaconda3\python.exe -u -c "import setuptools, tokenize;__file__
='C:\\Users\\USER\\AppData\\Local\\Temp\\pip-install-aulpkrue\\scikit-surprise\\setup.py';f=getattr(tokenize, 'o
pen', open)(__file__);code=f.read().replace('\r\n', '\n');f.close();exec(compile(code, __file__, 'exec'))" insta
ll --record C:\Users\USER\AppData\Local\Temp\pip-record-8ymavvp\install-record.txt --single-version-externally-
managed --compile:
running install
running build
running build_py
creating build
creating build\lib.win-amd64-3.6
creating build\lib.win-amd64-3.6\surprise
copying surprise\accuracy.py -> build\lib.win-amd64-3.6\surprise
copying surprise\builtin_datasets.py -> build\lib.win-amd64-3.6\surprise
copying surprise\dataset.py -> build\lib.win-amd64-3.6\surprise
copying surprise\dump.py -> build\lib.win-amd64-3.6\surprise
copying surprise\evaluate.py -> build\lib.win-amd64-3.6\surprise
copying surprise\reader.py -> build\lib.win-amd64-3.6\surprise
copying surprise\trainset.py -> build\lib.win-amd64-3.6\surprise
copying surprise\utils.py -> build\lib.win-amd64-3.6\surprise
copying surprise\__init__.py -> build\lib.win-amd64-3.6\surprise
copying surprise\__main__.py -> build\lib.win-amd64-3.6\surprise
creating build\lib.win-amd64-3.6\surprise\model_selection
copying surprise\model_selection\search.py -> build\lib.win-amd64-3.6\surprise\model_selection
copying surprise\model_selection\split.py -> build\lib.win-amd64-3.6\surprise\model_selection
copying surprise\model_selection\validation.py -> build\lib.win-amd64-3.6\surprise\model_selection
copying surprise\model_selection\__init__.py -> build\lib.win-amd64-3.6\surprise\model_selection
creating build\lib.win-amd64-3.6\surprise\prediction_algorithms
copying surprise\prediction_algorithms\algo_base.py -> build\lib.win-amd64-3.6\surprise\prediction_algorith
s
copying surprise\prediction_algorithms\baseline_only.py -> build\lib.win-amd64-3.6\surprise\prediction_algor
ithms
copying surprise\prediction_algorithms\knns.py -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
copying surprise\prediction_algorithms\predictions.py -> build\lib.win-amd64-3.6\surprise\prediction_algorit
hms
copying surprise\prediction_algorithms\random_pred.py -> build\lib.win-amd64-3.6\surprise\prediction_algorit

```

hms

```

    copying surprise\prediction_algorithms\__init__.py -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
    running egg_info
    writing scikit_surprise.egg-info\PKG-INFO
    writing dependency_links to scikit_surprise.egg-info\dependency_links.txt
    writing entry points to scikit_surprise.egg-info\entry_points.txt
    writing requirements to scikit_surprise.egg-info\requires.txt
    writing top-level names to scikit_surprise.egg-info\top_level.txt
    reading manifest file 'scikit_surprise.egg-info\SOURCES.txt'
    reading manifest template 'MANIFEST.in'
    writing manifest file 'scikit_surprise.egg-info\SOURCES.txt'
    copying surprise\similarities.c -> build\lib.win-amd64-3.6\surprise
    copying surprise\similarities.pyx -> build\lib.win-amd64-3.6\surprise
    copying surprise\prediction_algorithms\co_clustering.c -> build\lib.win-amd64-3.6\surprise\prediction_algori
thms
    copying surprise\prediction_algorithms\matrix_factorization.c -> build\lib.win-amd64-3.6\surprise\prediction
_algorithms
    copying surprise\prediction_algorithms\optimize_baselines.c -> build\lib.win-amd64-3.6\surprise\prediction_a
lgorithms
    copying surprise\prediction_algorithms\slope_one.c -> build\lib.win-amd64-3.6\surprise\prediction_algorithms
    copying surprise\prediction_algorithms\co_clustering.pyx -> build\lib.win-amd64-3.6\surprise\prediction_algo
rithms
    copying surprise\prediction_algorithms\matrix_factorization.pyx -> build\lib.win-amd64-3.6\surprise\predicti
on_algorithms
    copying surprise\prediction_algorithms\optimize_baselines.pyx -> build\lib.win-amd64-3.6\surprise\prediction
_algorithms
    copying surprise\prediction_algorithms\slope_one.pyx -> build\lib.win-amd64-3.6\surprise\prediction_algorith
ms
    running build_ext
    building 'surprise.similarities' extension
    error: Microsoft Visual C++ 14.0 is required. Get it with "Microsoft Visual C++ Build Tools": http://landing
hub.visualstudio.com/visual-cpp-build-tools

```

Failed building wheel for scikit-surprise
 Command "c:\users\user\anaconda3\python.exe -u -c "import setuptools, tokenize;__file__='C:\\Users\\USER\\AppData\\Local\\Temp\\pip-install-aulpkrue\\scikit-surprise\\setup.py';f=getattr(tokenize, 'open', open)(__file__);code=f.read().replace('\\r\\n', '\\n');f.close();exec(compile(code, __file__, 'exec'))" install --record C:\\Users\\USER\\AppData\\Local\\Temp\\pip-record-8ymavvpp\\install-record.txt --single-version-externally-managed --compile" failed with error code 1 in C:\\Users\\USER\\AppData\\Local\\Temp\\pip-install-aulpkrue\\scikit-surprise\\
 You are using pip version 18.0, however version 19.0.3 is available.
 You should consider upgrading via the 'python -m pip install --upgrade pip' command.

In [35]: **from surprise import** Reader, Dataset

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 (http://surprise.readthedocs.io/en/stable/getting_started.html#load-dom-dataframe-py)

In [36]: *# 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 [38]: testset = list(zip(reg_test_df.user.values, reg_test_df.movie.values, reg_test_df.rating.values))
testset[:3]
```

```
Out[38]: [(808635, 71, 5), (941866, 71, 4), (1280761, 71, 1)]
```

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 [39]: models_evaluation_train = dict()
models_evaluation_test = dict()

models_evaluation_train, models_evaluation_test
```

```
Out[39]: ({}, {})
```

Utility functions for running regression models

```

In [40]: # 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 : {}'.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('MAPE : ', mape_test)

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

Utility functions for Surprise modes

```

In [41]: # it is just to makesure 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 objecs
#####
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 : {}\nMAPE : {}".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 [42]: import xgboost as xgb
```

```
In [71]: def eval_xgboost(algo, x_train, y_train, 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 : {}'.format(datetime.now()-start))
        print('Done \n')
        print(algo.best_estimator_)
        # 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}
        print(train_results)
```

```
In [59]: # prepare Train data
from sklearn.model_selection import RandomizedSearchCV
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']
params={"n_estimators":[5,10,50,100,500],
        "max_depth":[5,10,50,100,500]}
# initialize Our first XGBoost model...
#for i in param1:
#    for j in param2:
#print("n_estimators:",i)
#print("max_depth:",j)
clf= xgb.XGBRegressor(silent=False, random_state=15)#, n_estimators=i,max_depth=j)
first_xgb=RandomizedSearchCV(clf,params,n_jobs=-1)
eval_xgboost(first_xgb, x_train, y_train)

# 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()
```

[illegible]

[http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

[http://localhost:8888/nbconvert/html/AaIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AaIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

[illegible]

[http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

[illegible]

[http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

[illegible]

[illegible]

[illegible]

[illegible]

[http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

```

[03:05:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 60 extra nodes, 0 pruned nodes, max_depth=5
[03:05:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:05:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 56 extra nodes, 0 pruned nodes, max_depth=5
[03:05:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 50 extra nodes, 0 pruned nodes, max_depth=5
[03:05:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 50 extra nodes, 0 pruned nodes, max_depth=5
[03:05:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:05:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:05:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
Done. Time taken : 0:45:54.667063

```

Done

```

XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
             colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
             max_depth=5, min_child_weight=1, missing=None, n_estimators=500,
             n_jobs=1, nthread=None, objective='reg:linear', random_state=15,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
             silent=False, subsample=1)
Evaluating the model with TRAIN data...
{'rmse': 0.8102828012161456, 'mape': 24.01797621968059, 'predictions': array([3.9565158, 3.751518 , 4.0210133,
..., 3.3002386, 2.6105006,
2.6670866], dtype=float32)}

```

```
In [70]: first_xgb = xgb.XGBRegressor(silent=False, n_jobs=-1, n_estimators=500, max_depth=5), random_state=15, n_estimators=100)
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_xgb'] = train_results
models_evaluation_test['first_xgb'] = test_results

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

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix Movie.ipynb?download=false](http://localhost:8888/nbconvert/html/AAIC/netflix-prize-data/Netflix%20Movie.ipynb?download=false)

[illegible]

[illegible]


```

[03:29:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 60 extra nodes, 0 pruned nodes, max_depth=5
[03:29:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:29:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:29:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 56 extra nodes, 0 pruned nodes, max_depth=5
[03:29:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 50 extra nodes, 0 pruned nodes, max_depth=5
[03:29:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 50 extra nodes, 0 pruned nodes, max_depth=5
[03:29:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 58 extra nodes, 0 pruned nodes, max_depth=5
[03:29:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
[03:29:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 62 extra nodes, 0 pruned nodes, max_depth=5
Done. Time taken : 0:02:11.881643

```

Done

Evaluating the model with TRAIN data...

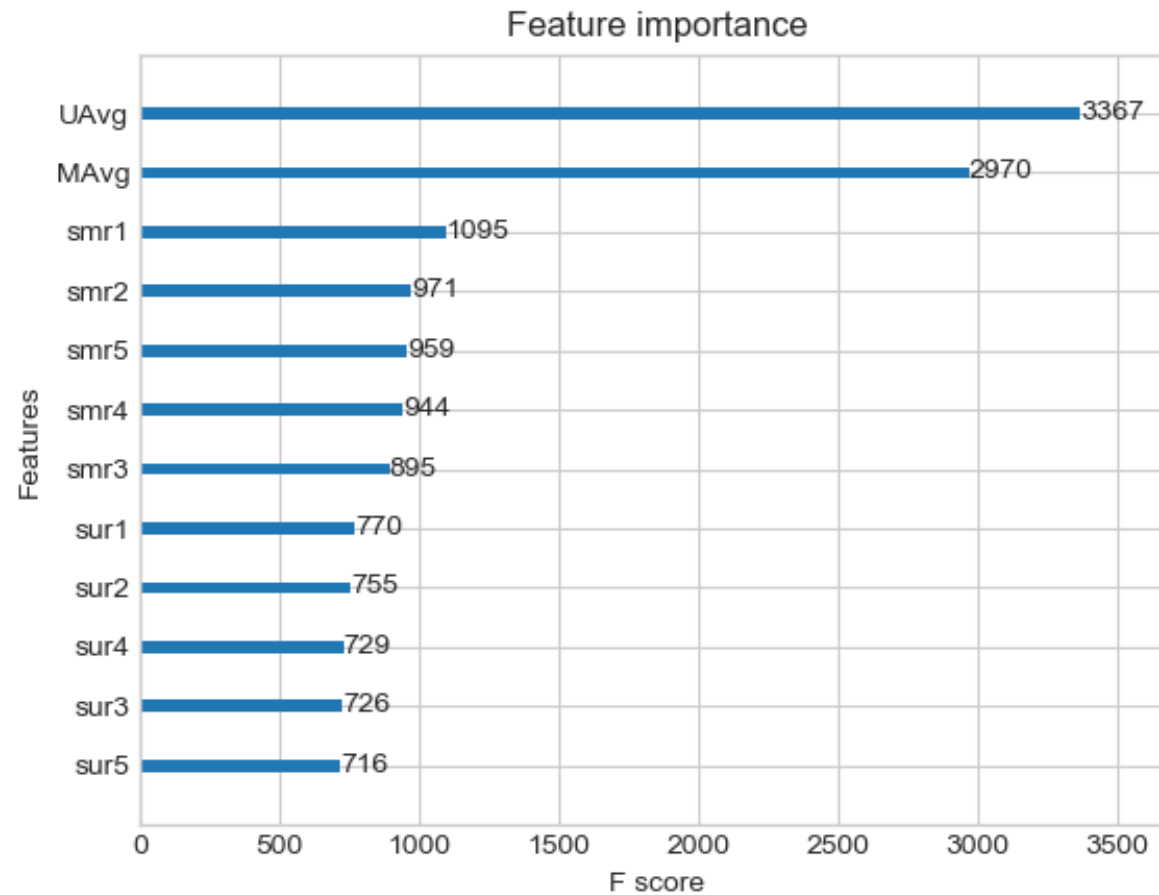
Evaluating Test data

TEST DATA

```

-----
RMSE : 1.1040796138786715
MAPE : 32.639338619055685

```



4.4.2 Surprise BaselineModel

```
In [62]: from surprise import BaselineOnly
```

Predicted_rating : (baseline prediction)

- http://surprise.readthedocs.io/en/stable/basic_algorithms.html#surprise.prediction_algorithms.baseline_only.BaselineOnly

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

- 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

$$\sum_{r_{ui} \in R_{train}} (r_{ui} - (\mu + b_u + b_i))^2 + \lambda (b_u^2 + b_i^2) . \text{ [mimimize } b_u, b_i]$$

```
In [65]: # 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 : 0:00:01.855652

Evaluating the model with train data..
time taken : 0:00:02.774529
-----
Train Data
-----
RMSE : 0.9339601465792874

MAPE : 29.580527712788847

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.290736
-----
Test Data
-----
RMSE : 1.0666751588756398

MAPE : 33.87684390285467

storing the test results in test dictionary...

-----
Total time taken to run this algorithm : 0:00:05.070996
```

4.4.3 XGBoost with initial 13 features + Surprise Baseline predictor

Updating Train Data

```
In [66]: # add our baseline_predicted value as our feature..
reg_train['bslpr'] = models_evaluation_train['bsl_algo']['predictions']
reg_train.head(2)
```

Out[66]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating	bslpr
0	53406	33	3.550651	4.0	5.0	5.0	4.0	1.0	5.0	2.0	5.0	3.0	1.0	3.366667	4.172414	4	3.994938
1	67390	33	3.550651	1.0	5.0	4.0	5.0	3.0	4.0	4.0	3.0	4.0	2.0	3.769231	4.172414	4	3.281518

Updating Test Data

```
In [67]: # 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[67]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5
0	808635	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651
1	941866	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651

```
In [74]: # 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']

params={"n_estimators": [5, 10, 50, 100, 500],
        "max_depth": [5, 10, 50, 100, 500]}
# initialize Our first XGBoost model...
xgb_bsl = xgb.XGBRegressor(silent=False, random_state=15)#, random_state=15, n_estimators=100)
xg = RandomizedSearchCV(xgb_bsl, params, n_jobs=-1)
#train_results, test_results =
eval_xgboost(xg, x_train, y_train)

# 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..

```
[05:01:29] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 7500 extra nodes, 0 pruned nodes, max_depth=24
[05:01:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 9930 extra nodes, 0 pruned nodes, max_depth=28
[05:01:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 12638 extra nodes, 0 pruned nodes, max_depth=28
[05:01:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 16056 extra nodes, 0 pruned nodes, max_depth=30
[05:01:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 21146 extra nodes, 0 pruned nodes, max_depth=33
[05:01:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 26656 extra nodes, 0 pruned nodes, max_depth=34
[05:01:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 32954 extra nodes, 0 pruned nodes, max_depth=33
[05:01:44] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 40344 extra nodes, 0 pruned nodes, max_depth=37
[05:01:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 49384 extra nodes, 0 pruned nodes, max_depth=45
[05:01:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 56536 extra nodes, 0 pruned nodes, max_depth=39
[05:01:55] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 66026 extra nodes, 0 pruned nodes, max_depth=44
[05:02:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 74014 extra nodes, 0 pruned nodes, max_depth=46
[05:02:04] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 83364 extra nodes, 0 pruned nodes, max_depth=42
[05:02:07] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 89708 extra nodes, 0 pruned nodes, max_depth=44
[05:02:11] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 95134 extra nodes, 0 pruned nodes, max_depth=50
[05:02:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98816 extra nodes, 0 pruned nodes, max_depth=46
[05:02:18] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 106716 extra nodes, 0 pruned nodes, max_depth=44
[05:02:22] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 109322 extra nodes, 0 pruned nodes, max_depth=48
[05:02:26] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 113544 extra nodes, 0 pruned nodes, max_depth=45
[05:02:29] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 119140 extra nodes, 0 pruned nodes, max_depth=45
```



```
h=50
[05:02:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 122024 extra nodes, 0 pruned nodes, max_dept
h=50
[05:02:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 123768 extra nodes, 0 pruned nodes, max_dept
h=45
[05:02:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 126950 extra nodes, 0 pruned nodes, max_dept
h=50
[05:02:45] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 128420 extra nodes, 0 pruned nodes, max_dept
h=50
[05:02:49] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 131858 extra nodes, 0 pruned nodes, max_dept
h=50
[05:02:53] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 134708 extra nodes, 0 pruned nodes, max_dept
h=50
[05:02:58] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 136882 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137734 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:06] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 140260 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 140018 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 142198 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 143590 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 145116 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:28] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 146862 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 147602 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 147630 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 150920 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:46] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151880 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:51] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151218 extra nodes, 0 pruned nodes, max_dept
h=50
[05:03:55] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 153160 extra nodes, 0 pruned nodes, max_dept
```

```
h=50
[05:03:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151920 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:04] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 154902 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:08] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 154326 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157036 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:17] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157608 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:22] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158070 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158034 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159436 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 155696 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160434 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:45] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162140 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162458 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:54] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 163300 extra nodes, 0 pruned nodes, max_dept
h=50
[05:04:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165874 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:03] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165822 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:08] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 164244 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159794 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:17] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166324 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:22] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162734 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 161348 extra nodes, 0 pruned nodes, max_dept
```

```
h=50
[05:05:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166552 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 168042 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165364 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:46] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166126 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160618 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:55] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 168452 extra nodes, 0 pruned nodes, max_dept
h=50
[05:05:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 167996 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:04] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 164606 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:09] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 163908 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162770 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:18] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159340 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162810 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:28] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159964 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157558 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160344 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:43] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159412 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158022 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 155496 extra nodes, 0 pruned nodes, max_dept
h=50
[05:06:57] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 153568 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 144842 extra nodes, 0 pruned nodes, max_dept
```

```
h=50
[05:07:08] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 141216 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:12] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 148056 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:17] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 142626 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:21] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137438 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:26] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137948 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:30] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 138366 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 134950 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 130700 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:44] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 128408 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:48] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 126262 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:53] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 123010 extra nodes, 0 pruned nodes, max_dept
h=50
[05:07:57] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 117430 extra nodes, 0 pruned nodes, max_dept
h=50
[05:08:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 106852 extra nodes, 0 pruned nodes, max_dept
h=50
[05:08:06] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 104384 extra nodes, 0 pruned nodes, max_dept
h=50
[05:08:11] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98224 extra nodes, 0 pruned nodes, max_dept
=50
[05:08:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 105336 extra nodes, 0 pruned nodes, max_dept
h=50
[05:08:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 100218 extra nodes, 0 pruned nodes, max_dept
h=50
[05:08:24] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98068 extra nodes, 0 pruned nodes, max_dept
=50
[05:08:29] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 94742 extra nodes, 0 pruned nodes, max_dept
=50
[05:08:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 91608 extra nodes, 0 pruned nodes, max_dept
```

```
=50
[05:08:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 87992 extra nodes, 0 pruned nodes, max_depth
=50
[05:08:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 82638 extra nodes, 0 pruned nodes, max_depth
=50
[05:08:46] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 80980 extra nodes, 0 pruned nodes, max_depth
=50
[05:08:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 75334 extra nodes, 0 pruned nodes, max_depth
=50
[05:08:54] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 74806 extra nodes, 0 pruned nodes, max_depth
=50
[05:08:58] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 70558 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 68324 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:07] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 60324 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:11] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 43102 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 59992 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 55100 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 53908 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 51770 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 48888 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:35] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 45850 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:39] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 42750 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:44] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 40352 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:48] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 37568 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 32538 extra nodes, 0 pruned nodes, max_depth
=50
[05:09:56] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 31888 extra nodes, 0 pruned nodes, max_depth
```

```
=50
[05:10:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 26632 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:03] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 25948 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:07] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 22834 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:11] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 19656 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 21522 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:18] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 17490 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:20] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 5900 extra nodes, 0 pruned nodes, max_depth=
38
[05:10:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 10264 extra nodes, 0 pruned nodes, max_depth
=50
[05:10:25] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 2698 extra nodes, 0 pruned nodes, max_depth=
31
[05:10:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 4356 extra nodes, 0 pruned nodes, max_depth=
40
[05:10:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 8082 extra nodes, 0 pruned nodes, max_depth=
50
[05:10:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 2788 extra nodes, 0 pruned nodes, max_depth=
36
[05:10:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 5244 extra nodes, 0 pruned nodes, max_depth=
47
[05:10:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 2932 extra nodes, 0 pruned nodes, max_depth=
34
[05:10:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 2862 extra nodes, 0 pruned nodes, max_depth=
38
[05:10:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 1692 extra nodes, 0 pruned nodes, max_depth=
32
[05:10:43] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 5200 extra nodes, 0 pruned nodes, max_depth=
38
[05:10:45] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 1172 extra nodes, 0 pruned nodes, max_depth=
27
[05:10:46] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 1820 extra nodes, 0 pruned nodes, max_depth=
37
[05:10:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 744 extra nodes, 0 pruned nodes, max_depth=2
```

```
6
[05:10:49] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 1644 extra nodes, 0 pruned nodes, max_depth=36
[05:10:49] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 216 extra nodes, 0 pruned nodes, max_depth=16
[05:10:51] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 836 extra nodes, 0 pruned nodes, max_depth=30
[05:10:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 1454 extra nodes, 0 pruned nodes, max_depth=29
[05:10:53] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 284 extra nodes, 0 pruned nodes, max_depth=23
[05:10:54] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 888 extra nodes, 0 pruned nodes, max_depth=26
[05:10:54] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 278 extra nodes, 0 pruned nodes, max_depth=17
[05:10:57] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 4324 extra nodes, 0 pruned nodes, max_depth=41
[05:10:58] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 696 extra nodes, 0 pruned nodes, max_depth=31
[05:10:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 616 extra nodes, 0 pruned nodes, max_depth=21
[05:10:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:10:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:10:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
```

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

```
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
[05:11:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 0 extra nodes, 0 pruned nodes, max_depth=0
Done. Time taken : 1:21:19.834434
```

Done

```
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
              max_depth=50, min_child_weight=1, missing=None, n_estimators=500,
              n_jobs=1, nthread=None, objective='reg:linear', random_state=15,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=False, subsample=1)
Evaluating the model with TRAIN data...
{'rmse': 0.0006400084066093161, 'mape': 0.014000415130326004, 'predictions': array([3.9999523, 3.9999816, 3.0001
533, ..., 3.999276 , 1.0001268,
1.0002437], dtype=float32)}
```



```
In [76]: xgb_bsl = xgb.XGBRegressor(silent=False, random_state=15, n_estimators=100,max_depth=50)
#xg = RandomizedSearchCV(xgb_bsl,params,n_jobs=-1)
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..

```
[05:15:49] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 7500 extra nodes, 0 pruned nodes, max_depth=24
[05:15:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 9930 extra nodes, 0 pruned nodes, max_depth=28
[05:15:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 12638 extra nodes, 0 pruned nodes, max_depth=28
[05:15:53] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 16056 extra nodes, 0 pruned nodes, max_depth=30
[05:15:55] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 21146 extra nodes, 0 pruned nodes, max_depth=33
[05:15:57] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 26656 extra nodes, 0 pruned nodes, max_depth=34
[05:15:59] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 32954 extra nodes, 0 pruned nodes, max_depth=33
[05:16:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 40344 extra nodes, 0 pruned nodes, max_depth=37
[05:16:05] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 49384 extra nodes, 0 pruned nodes, max_depth=45
[05:16:07] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 56536 extra nodes, 0 pruned nodes, max_depth=39
[05:16:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 66026 extra nodes, 0 pruned nodes, max_depth=44
[05:16:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 74014 extra nodes, 0 pruned nodes, max_depth=46
[05:16:17] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 83364 extra nodes, 0 pruned nodes, max_depth=42
[05:16:20] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 89708 extra nodes, 0 pruned nodes, max_depth=44
[05:16:24] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 95134 extra nodes, 0 pruned nodes, max_depth=50
[05:16:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98816 extra nodes, 0 pruned nodes, max_depth=46
[05:16:31] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 106716 extra nodes, 0 pruned nodes, max_depth=44
[05:16:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 109322 extra nodes, 0 pruned nodes, max_depth=48
[05:16:40] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 113544 extra nodes, 0 pruned nodes, max_depth=45
[05:16:44] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 119140 extra nodes, 0 pruned nodes, max_depth=45
```

h=50
[05:16:48] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 122024 extra nodes, 0 pruned nodes, max_dept
h=50
[05:16:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 123768 extra nodes, 0 pruned nodes, max_dept
h=45
[05:16:56] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 126950 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 128420 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:04] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 131858 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:08] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 134708 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:12] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 136882 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:17] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137734 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:21] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 140260 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:25] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 140018 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:30] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 142198 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:34] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 143590 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 145116 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:44] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 146862 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:48] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 147602 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:53] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 147630 extra nodes, 0 pruned nodes, max_dept
h=50
[05:17:57] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 150920 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:02] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151880 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:06] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151218 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 153160 extra nodes, 0 pruned nodes, max_dept

h=50
[05:18:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 151920 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 154902 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:24] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 154326 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:29] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157036 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157608 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158070 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158034 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159436 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:51] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 155696 extra nodes, 0 pruned nodes, max_dept
h=50
[05:18:56] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160434 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162140 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:05] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162458 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 163300 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165874 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165822 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 164244 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:28] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159794 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166324 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162734 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:42] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 161348 extra nodes, 0 pruned nodes, max_dept

h=50
[05:19:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166552 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 168042 extra nodes, 0 pruned nodes, max_dept
h=50
[05:19:56] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 165364 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 166126 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:05] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160618 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 168452 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:15] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 167996 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:22] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 164606 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:27] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 163908 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162770 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:36] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159340 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 162810 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:46] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159964 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:50] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 157558 extra nodes, 0 pruned nodes, max_dept
h=50
[05:20:55] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 160344 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:00] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 159412 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:04] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 158022 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:09] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 155496 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:13] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 153568 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:18] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 144842 extra nodes, 0 pruned nodes, max_dept

```
h=50
[05:21:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 141216 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:28] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 148056 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:33] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 142626 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:38] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137438 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:43] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 137948 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:47] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 138366 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:52] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 134950 extra nodes, 0 pruned nodes, max_dept
h=50
[05:21:56] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 130700 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:01] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 128408 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:05] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 126262 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:10] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 123010 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:14] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 117430 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:19] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 106852 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:23] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 104384 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:28] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98224 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:32] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 105336 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:37] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 100218 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:41] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 98068 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:45] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 94742 extra nodes, 0 pruned nodes, max_dept
h=50
[05:22:49] src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 91608 extra nodes, 0 pruned nodes, max_dept
```

=50

Done. Time taken : 0:07:02.908068

Done

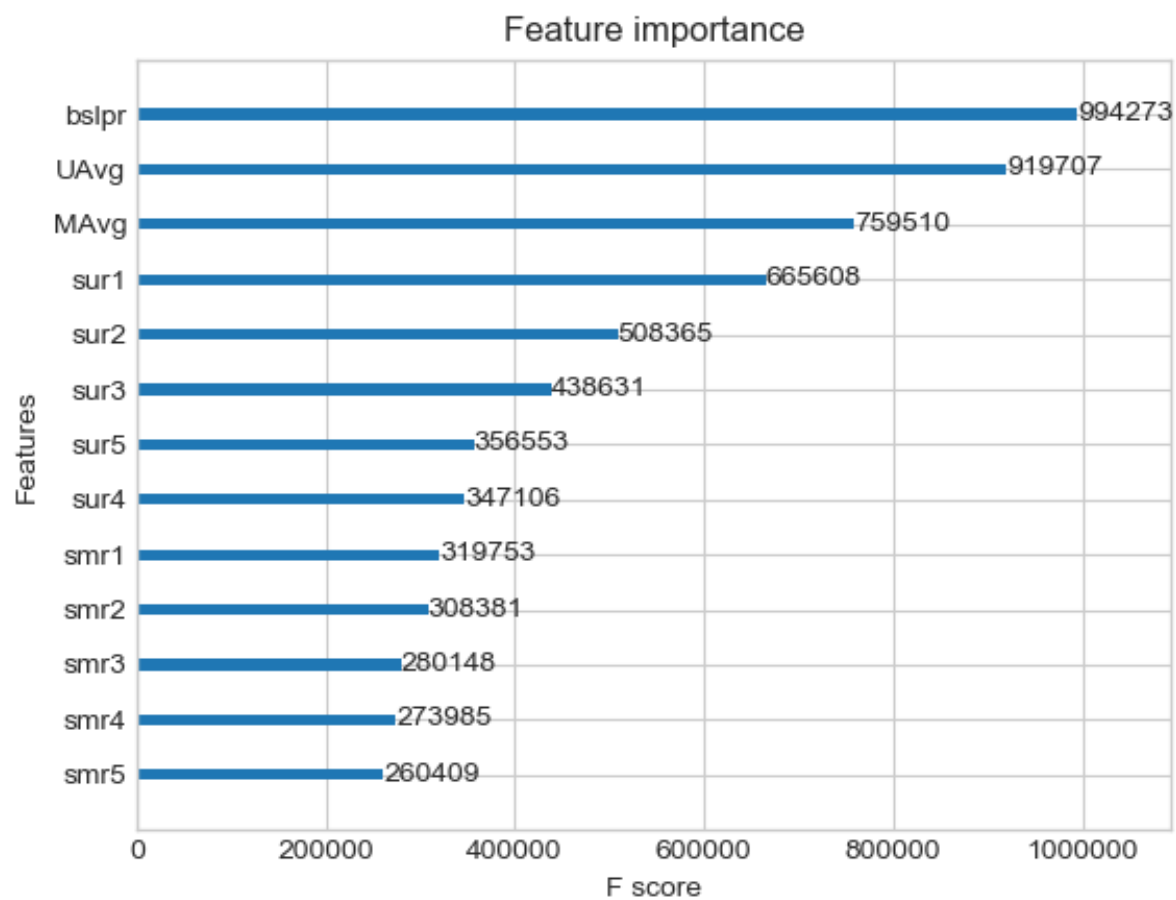
Evaluating the model with TRAIN data...

Evaluating Test data

TEST DATA

RMSE : 1.0718407691383005

MAPE : 33.629856818613206



4.4.4 Surprise KNNBaseline predictor

```
In [77]: from surprise import KNNBaseline
```


- KNN BASELINE

- http://surprise.readthedocs.io/en/stable/knn_inspired.html#surprise.prediction_algorithms.knns.KNNBaseline
(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
(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>
(<http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf>)

- predicted Rating : (*based on User-User similarity*)

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

- b_{ui} - Baseline prediction of (user, movie) rating
- $N_i^k(u)$ - 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):

$$\hat{r}_{ui} = b_{ui} + \frac{\sum_{j \in N_u^k(i)} \text{sim}(i, j) \cdot (r_{uj} - b_{uj})}{\sum_{j \in N_u^k(j)} \text{sim}(i, j)}$$

- *Notations follows same as above (user user based predicted rating)*

4.4.4.1 Surprise KNNBaseline with user user similarities

```
In [78]: # 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 : 0:02:56.141487

Evaluating the model with train data..
time taken : 0:06:10.659342
-----
Train Data
-----
RMSE : 0.35594237282596297

MAPE : 9.825864233671071

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.326127
-----
Test Data
-----
RMSE : 1.0660901365546087

MAPE : 33.91894303653171

storing the test results in test dictionary...

-----
Total time taken to run this algorithm : 0:09:07.296023
```

4.4.4.2 Surprise KNNBaseline with movie movie similarities

```
In [79]: # 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
```

```
Training the model...
Estimating biases using sgd...
Computing the pearson_baseline similarity matrix...
Done computing similarity matrix.
Done. time taken : 0:00:05.370094

Evaluating the model with train data..
time taken : 0:00:30.451290
-----
Train Data
-----
RMSE : 0.3592155923736147

MAPE : 9.54501751396511

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.307278
-----
Test Data
-----
RMSE : 1.066115023069518

MAPE : 33.92139376258056

storing the test results in test dictionary...

-----
Total time taken to run this algorithm : 0:00:36.269288
```

4.4.5 XGBoost with initial 13 features + Surprise Baseline predictor + KNNBaseline predictor

- First we will run XGBoost with predictions from both KNN's (that uses User_User and Item_Item similarities along with our previous features.
 - Then we will run XGBoost with just predictions from both knn models and predictions from our baseline model.

Preparing Train data

```
In [80]: # add the predicted values from both knns to this dataframe
reg_train['knn_bsl_u'] = models_evaluation_train['knn_bsl_u']['predictions']
reg_train['knn_bsl_m'] = models_evaluation_train['knn_bsl_m']['predictions']

reg_train.head(2)
```

Out[80]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating	bslpr	knn
0	53406	33	3.550651	4.0	5.0	5.0	4.0	1.0	5.0	2.0	5.0	3.0	1.0	3.366667	4.172414	4	3.994938	3.9
1	67390	33	3.550651	1.0	5.0	4.0	5.0	3.0	4.0	4.0	3.0	4.0	2.0	3.769231	4.172414	4	3.281518	3.0

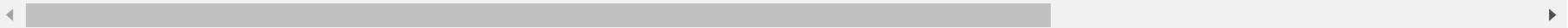
Preparing Test data

```
In [81]: reg_test_df['knn_bsl_u'] = models_evaluation_test['knn_bsl_u']['predictions']
reg_test_df['knn_bsl_m'] = models_evaluation_test['knn_bsl_m']['predictions']

reg_test_df.head(2)
```

Out[81]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5
0	808635	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651
1	941866	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651




```

In [83]: # 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']

params = {"n_estimators": [10, 50, 100, 300, 500],
          "max_depth": [5, 10, 15, 25, 50]}
# declare the model
clf = xgb.XGBRegressor(random_state=15)
xgb_knn_bsl = RandomizedSearchCV(clf, params, n_jobs=-1)
#train_results, test_results =
eval_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..

Done. Time taken : 1:09:00.445583

Done

```

XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
             colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
             max_depth=5, min_child_weight=1, missing=None, n_estimators=100,
             n_jobs=1, nthread=None, objective='reg:linear', random_state=15,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
             silent=True, subsample=1)

```

Evaluating the model with TRAIN data...

```

{'rmse': 0.8386894854749825, 'mape': 25.129227829799728, 'predictions': array([3.989349 , 3.9864907, 4.010354 ,
..., 3.3841488, 2.7137415,
2.9226072], dtype=float32)}

```

```
In [85]: xgb_knn_bsl = xgb.XGBRegressor(random_state=15,n_estimators=100,max_depth=5,n_jobs=-1)
        #= RandomizedSearchCV(clf,params,n_jobs=-1)
        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()
```

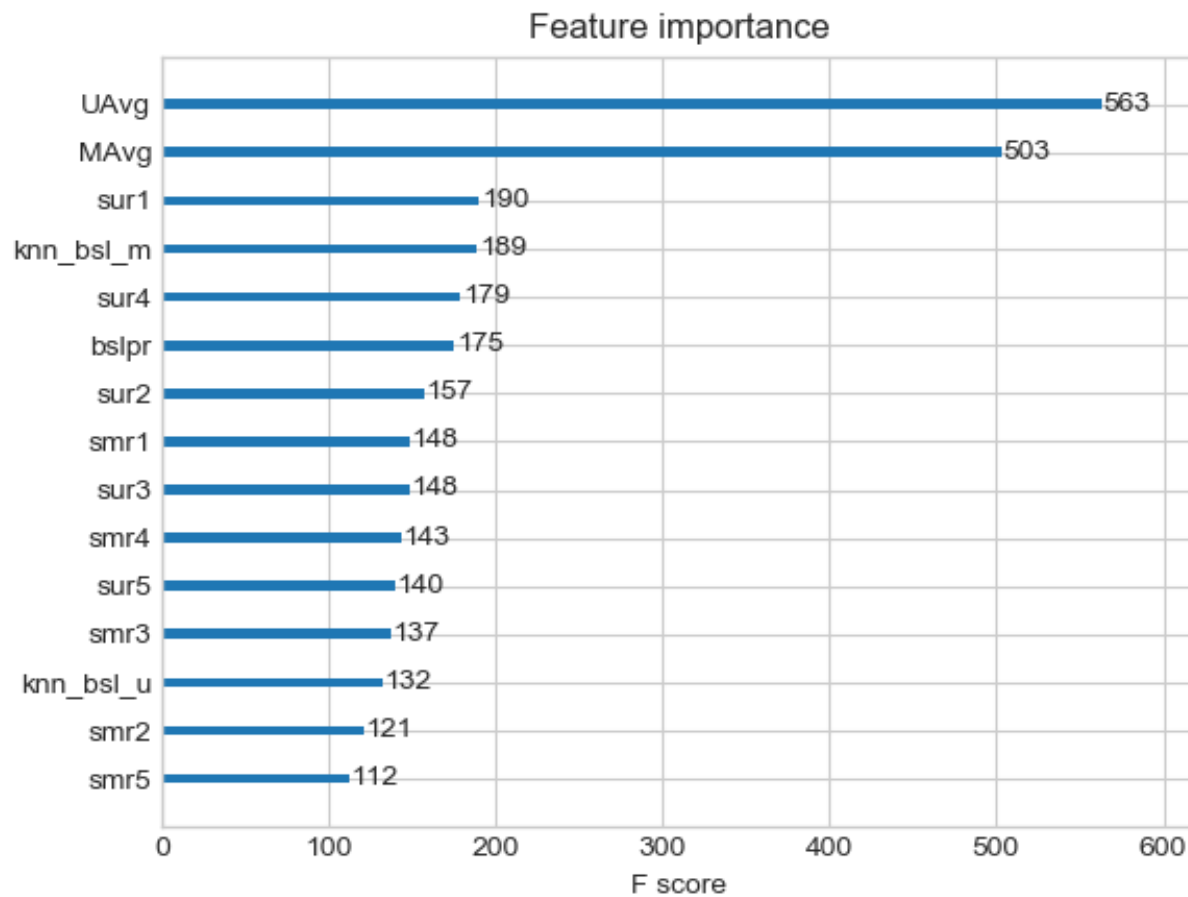
```
Training the model..  
Done. Time taken : 0:00:37.043709
```

```
Done
```

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

```
TEST DATA
```

```
-----  
RMSE : 1.0731568165186072  
MAPE : 33.39510588691888
```



4.4.6 Matrix Factorization Techniques

4.4.6.1 SVD Matrix Factorization User Movie intractions

```
In [86]: from surprise import SVD
```

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

- Predicted Rating :

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

- 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)
 ([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)

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

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

```
In [87]: # initiallize the model
svd = SVD(n_factors=100, biased=True, random_state=15, verbose=True)
svd_train_results, svd_test_results = run_surprise(svd, trainset, testset, verbose=True)

# 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...
Processing epoch 0
Processing epoch 1
Processing epoch 2
Processing epoch 3
Processing epoch 4
Processing epoch 5
Processing epoch 6
Processing epoch 7
Processing epoch 8
Processing epoch 9
Processing epoch 10
Processing epoch 11
Processing epoch 12
Processing epoch 13
Processing epoch 14
Processing epoch 15
Processing epoch 16
Processing epoch 17
Processing epoch 18
Processing epoch 19
Done. time taken : 0:00:28.329621

Evaluating the model with train data..
time taken : 0:00:03.457902
-----
Train Data
-----
RMSE : 0.6588520545653486

MAPE : 19.923644355504923

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.286719
-----
Test Data
-----
RMSE : 1.0663604556104023
```

MAPE : 33.853122640919224

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:00:32.212860

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

In [88]: `from surprise import SVDpp`

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

- Predicted Rating :

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

- 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)

$$- \sum_{r_{ui} \in R_{\text{train}}} \left(r_{ui} - \hat{r}_{ui} \right)^2 + \lambda \left(b_i^2 + b_u^2 + \|q_i\|^2 + \|p_u\|^2 + \|y_j\|^2 \right)$$

```
In [89]: # initialize the model
svdpp = SVDpp(n_factors=50, random_state=15, verbose=True)
svdpp_train_results, svdpp_test_results = run_surprise(svdpp, trainset, testset, verbose=True)

# 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...

processing epoch 0
processing epoch 1
processing epoch 2
processing epoch 3
processing epoch 4
processing epoch 5
processing epoch 6
processing epoch 7
processing epoch 8
processing epoch 9
processing epoch 10
processing epoch 11
processing epoch 12
processing epoch 13
processing epoch 14
processing epoch 15
processing epoch 16
processing epoch 17
processing epoch 18
processing epoch 19

Done. time taken : 0:07:51.414485

Evaluating the model with train data..

time taken : 0:00:21.679479

Train Data

RMSE : 0.6103784814253126

MAPE : 17.885904845659287

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.326160

Test Data

RMSE : 1.0667336314319098

MAPE : 33.77324600603565

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:08:13.570231

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

Preparing Train data

```
In [90]: # 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[90]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	...	smr4	smr5	UAvg	MAvg	rating	bslpr	knn_b
0	53406	33	3.550651	4.0	5.0	5.0	4.0	1.0	5.0	2.0	...	3.0	1.0	3.366667	4.172414	4	3.994938	3.9816
1	67390	33	3.550651	1.0	5.0	4.0	5.0	3.0	4.0	4.0	...	4.0	2.0	3.769231	4.172414	4	3.281518	3.0781

2 rows × 21 columns



Preparing Test data

```
In [91]: reg_test_df['svd'] = models_evaluation_test['svd']['predictions']
reg_test_df['svdpp'] = models_evaluation_test['svdpp']['predictions']

reg_test_df.head(2)
```

Out[91]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	...	smr4	smr5	UA
0	808635	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	...	3.550651	3.550651	3.5506
1	941866	71	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	3.550651	...	3.550651	3.550651	3.5506

2 rows × 21 columns



```
In [92]: # prepare x_train and y_train
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']

params = {"n_estimators": [10, 50, 100, 300, 500],
          "max_depth": [5, 10, 15, 25, 50]}

clf = xgb.XGBRegressor(random_state=15)
xgb_final = RandomizedSearchCV(clf, params, n_jobs=-1)
#train_results, test_results =
eval_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..

Done. Time taken : 1:28:50.423157

Done

```
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
              colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
              max_depth=5, min_child_weight=1, missing=None, n_estimators=300,
              n_jobs=1, nthread=None, objective='reg:linear', random_state=15,
              reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
              silent=True, subsample=1)
```

Evaluating the model with TRAIN data...

```
{'rmse': 0.8252417721787112, 'mape': 24.656976357438257, 'predictions': array([4.017031 , 3.9772937, 4.01476 ,
..., 3.3244882, 2.6878016,
      2.9244761], dtype=float32)}
```

```
In [94]: xgb_final = xgb.XGBRegressor(random_state=15,n_jobs=-1,n_estimators=300,max_depth=5)
# = RandomizedSearchCV(clf,params,n_jobs=-1)
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()
```



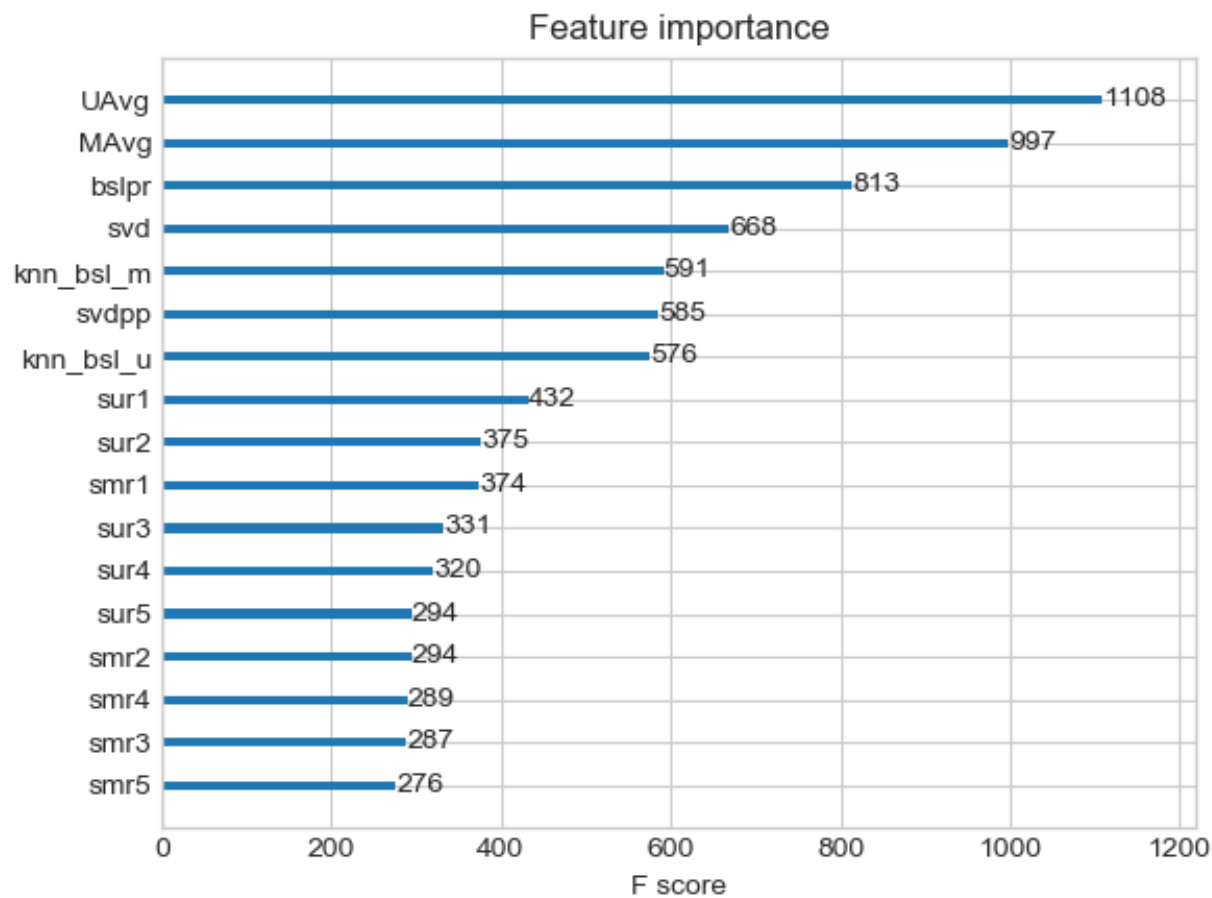
```
Training the model..  
Done. Time taken : 0:02:00.844933
```

```
Done
```

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

```
TEST DATA
```

```
-----  
RMSE : 1.0695683552716242  
MAPE : 33.600611413157374
```



4.4.8 XgBoost with Surprise Baseline + Surprise KNNbaseline + MF Techniques

```
In [95]: # prepare train data
x_train = reg_train[['knn_bsl_u', 'knn_bsl_m', 'svd', 'svdpp']]
y_train = reg_train['rating']

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

params = {"n_estimators": [10, 50, 100, 300, 500],
          "max_depth": [5, 10, 15, 25, 50]}
clf = xgb.XGBRegressor(random_state=15)
xgb_all_models = RandomizedSearchCV(clf, params, n_jobs=-1)
#train_results, test_results =
eval_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..

Done. Time taken : 0:50:43.344517

Done

```
XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
             colsample_bytrees=1, gamma=0, learning_rate=0.1, max_delta_step=0,
             max_depth=5, min_child_weight=1, missing=None, n_estimators=300,
             n_jobs=1, nthread=None, objective='reg:linear', random_state=15,
             reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
             silent=True, subsample=1)
```

Evaluating the model with TRAIN data...

```
{'rmse': 1.061079940008031, 'mape': 35.12051769544633, 'predictions': array([3.5358932, 3.5344315, 3.5580945,
..., 3.5085697, 3.5211134,
3.5425947], dtype=float32)}
```

```
In [98]: xgb_all_models = xgb.XGBRegressor(random_state=15,n_jobs=-1,n_estimators=300,max_depth=5)
# = RandomizedSearchCV(clf,params,n_jobs=-1)
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()
```

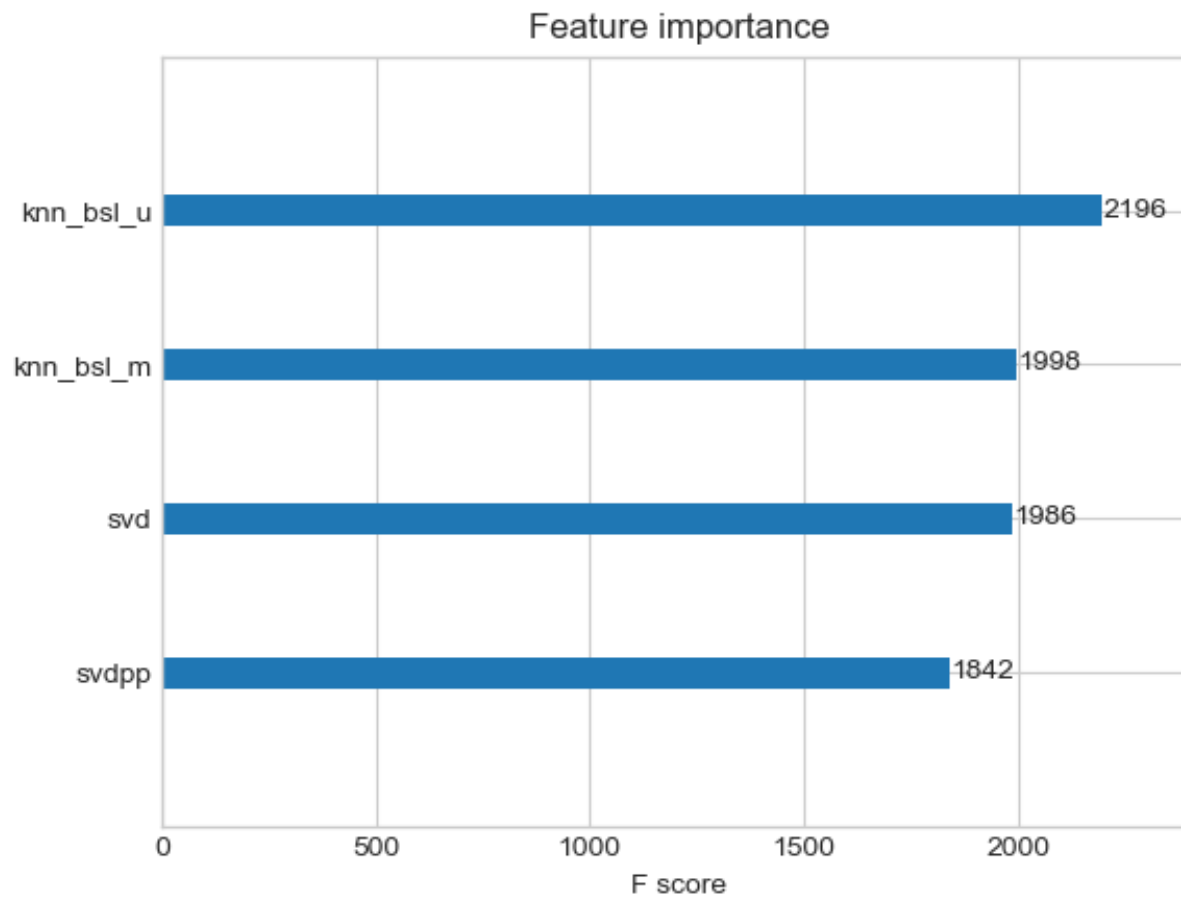
```
Training the model..  
Done. Time taken : 0:00:56.470461
```

```
Done
```

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

```
TEST DATA
```

```
-----  
RMSE : 1.0695163325135204  
MAPE : 34.37569586757193
```



4.5 Comparision between all models

```
In [100]: # Saving our TEST_RESULTS into a dataframe so that you don't have to run it again
pd.DataFrame(models_evaluation_test).to_csv('small_sample_results.csv')
models = pd.read_csv('small_sample_results.csv', index_col=0)
models.loc['rmse'].sort_values()
```

```
Out[100]: knn_bsl_u      1.0660901365546087
knn_bsl_m      1.066115023069518
svd      1.0663604556104023
bsl_algo      1.0666751588756398
svdpp      1.0667336314319098
xgb_all_models  1.0695163325135204
xgb_final      1.0695683552716242
xgb_bsl      1.0718407691383005
xgb_knn_bsl    1.0731568165186072
first_algo      1.102015760805355
first_xgb      1.1040796138786715
Name: rmse, dtype: object
```

In the above table we can see that knn baseline model with user-user similarity is tend to perform well compared to all the models.

SUMMARY:

- **Netflix movie recommendation is itself a very complex problem to solve. This implementation is based on the research paper Factor in the Neighbors: Scalable and Accurate Collaborative Filtering written by Yehuda Koren.**
- **The metric which is used in this problem is RMSE(Root Mean square error). This metric measures the differences between the values predicted by the model and actual values observed.**

step-by-step procedure taken to solve this problem:

CHAPTER-1: Data loading

1. The data itself comes with four txt files. so, write all the data into data.csv file and read the data into pandas dataframe\
2. Check whether any null values and duplicates are present in the data.
3. Split the data into 80:20 and write into train.csv and test.csv files.

CHAPTER-2: Exploratory data analysis

1. Analysis of rating is performed.
2. Analyzing no. of ratings given in a month.
3. Adding new column day of the week.
4. Analyze ratings given by user by plotting pdf and cdf.
5. Creating a sparse matrix for train and test data. where columns are movies and rows are users.
6. Getting the global averages of movie rating rated by each user and average rating rated per movie.

CHAPTER-3: Computing Similarity Matrix

1. Computing movie-movie similarity.

CHAPTER-4: Implementing Models:

1. This problem can be posed as both recommendation and regression problem.
2. So, we use XGBRegressor surprise library to solve this problem.
3. First, sample the data because solving this problem is a non-trivial task and featurize the data for both regression and surprise library.

Steps:

- A. Perform XGBRegression using featurized regression data.
- B. perform SurpriseBaseLine model on featurized surprised library data.
- C. ADD the output of the baseline model as a feature of regression data and perform regression data.

- D. perform SurpriseKNN for user-user and item-item.
- E. The output of the both user-user and item-item is add to the regression data and XGBRegressor is performed.
- F. Perform matrixfatorization of SVD and SVD++ add the output is added as a feature to the regression data.
- G. By dropping Suprisebaseline feature XgbRegressor is performed.
- H. Finally, by taking all the data XGBRegressor is performed.