# 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 ahs 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
  </ul>

```
    The first line of each file [combined_data_1.txt, combined_da
    ta_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 480
    189 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 b
e 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]:
```python
# this is just to know how much time will it take to run this entire ip
ython notebook
from datetime import datetime
# globalstart = datetime.now()
import pandas as pd
import numpy as np
```

```python
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

```python
In [0]: 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('dat
        a.csv')
            # We re reading from each of the four files and appendig each ratin
        g 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_dat
a_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 anoth
er 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)
```

```
Reading ratings from data_folder/combined_data_1.txt...
Done.

Reading ratings from data_folder/combined_data_2.txt...
Done.

Reading ratings from data_folder/combined_data_3.txt...
Done.

Reading ratings from data_folder/combined_data_4.txt...
Done.

Time taken : 0:05:03.705966
```

In [0]:
```
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 [0]: `df.head()`

Out[0]:

|  | 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 [0]: `df.describe()['rating']`

Out[0]:
```
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 [0]:
```
# just to make sure that all Nan containing rows are deleted..
print("No of Nan values in our dataframe : ", sum(df.isnull().any()))
```

No of Nan values in our dataframe :  0

### 3.1.3 Removing Duplicates

In [0]:
```
dup_bool = df.duplicated(['movie','user','rating'])
dups = sum(dup_bool) # by considering all columns..( including timestam
p)
print("There are {} duplicate rating entries in the data..".format(dups
))
```

There are 0 duplicate rating entries in the data..

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

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

Total data
--------------------------------------------------

Total no of ratings : 100480507
Total No of Users   : 480189
Total No of movies  : 17770

## 3.2 Spliting data into Train and Test(80:20)

In [0]:
```python
if not os.path.isfile('train.csv'):
    # create the dataframe and store it in the disk for offline purpose
s..
    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 purpose
s..
    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 [0]:
```python
# 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 [0]: print("Test data ")
        print("-"*50)
        print("\nTotal no of ratings :",test_df.shape[0])
        print("Total No of Users   :", len(np.unique(test_df.user)))
        print("Total No of movies  :", len(np.unique(test_df.movie)))
```

```
Test data
--------------------------------------------------

Total no of ratings : 20096102
Total No of Users   : 349312
Total No of movies  : 17757
```

## 3.3 Exploratory Data Analysis on Train data

```
In [0]: # method to make y-axis more readable
        def human(num, units = 'M'):
            units = units.lower()
            num = float(num)
            if units == 'k':
                return str(num/10**3) + " K"
            elif units == 'm':
                return str(num/10**6) + " M"
            elif units == 'b':
                return str(num/10**9) +  " B"
```

### 3.3.1 Distribution of ratings

```
In [0]: 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()
```

## Distribution of ratings over Training dataset



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

```
In [0]:  # It is used to skip the warning ''SettingWithCopyWarning''..
         pd.options.mode.chained_assignment = None  # default='warn'

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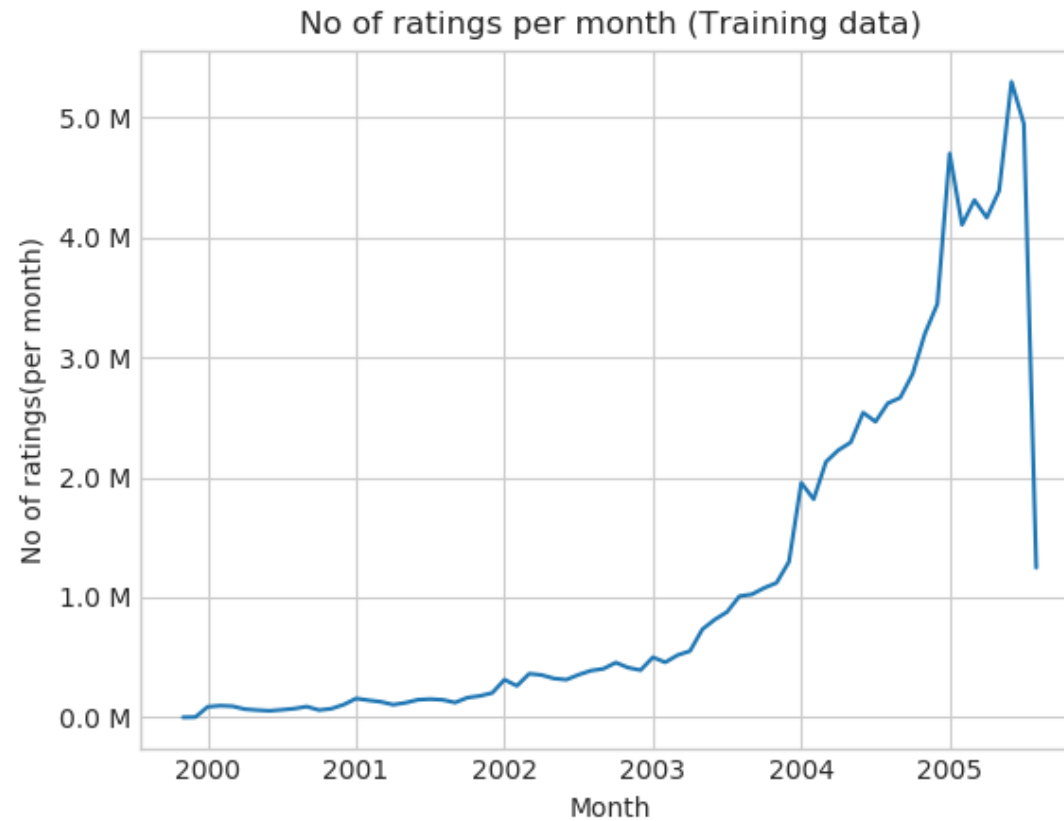
```
train_df.tail()
```

|  | 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 [0]:
```python
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()
```

## No of ratings per month (Training data)



### 3.3.3 Analysis on the Ratings given by user

```
In [0]: no_of_rated_movies_per_user = train_df.groupby(by='user')['rating'].cou
        nt().sort_values(ascending=False)

        no_of_rated_movies_per_user.head()
```

```
Out[0]: user
        305344     17112
```
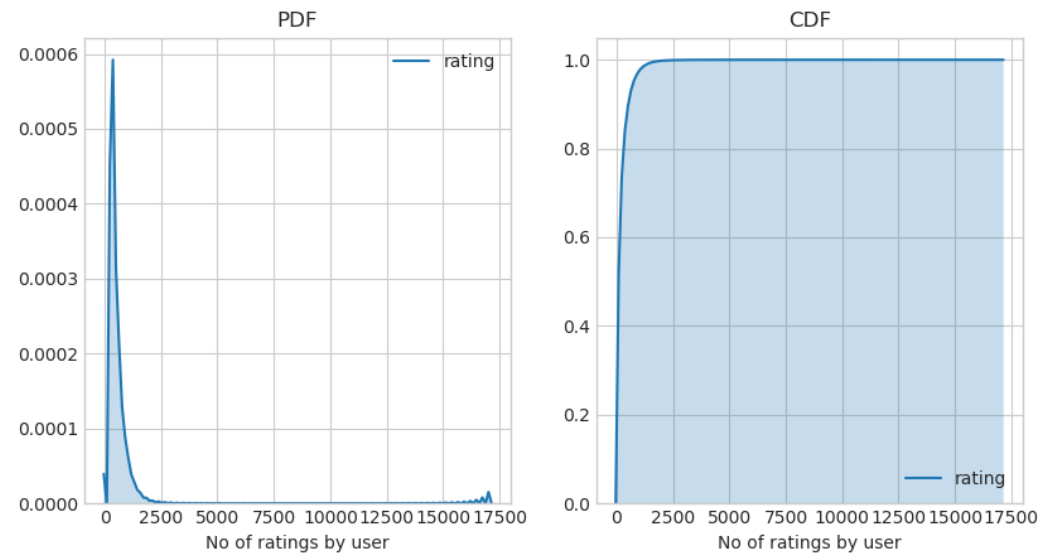
```
2439493     15896
387418      15402
1639792      9767
1461435      9447
Name: rating, dtype: int64
```

In [0]:
```python
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()
```

```python
In [0]: no_of_rated_movies_per_user.describe()
```

```
Out[0]: 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
```
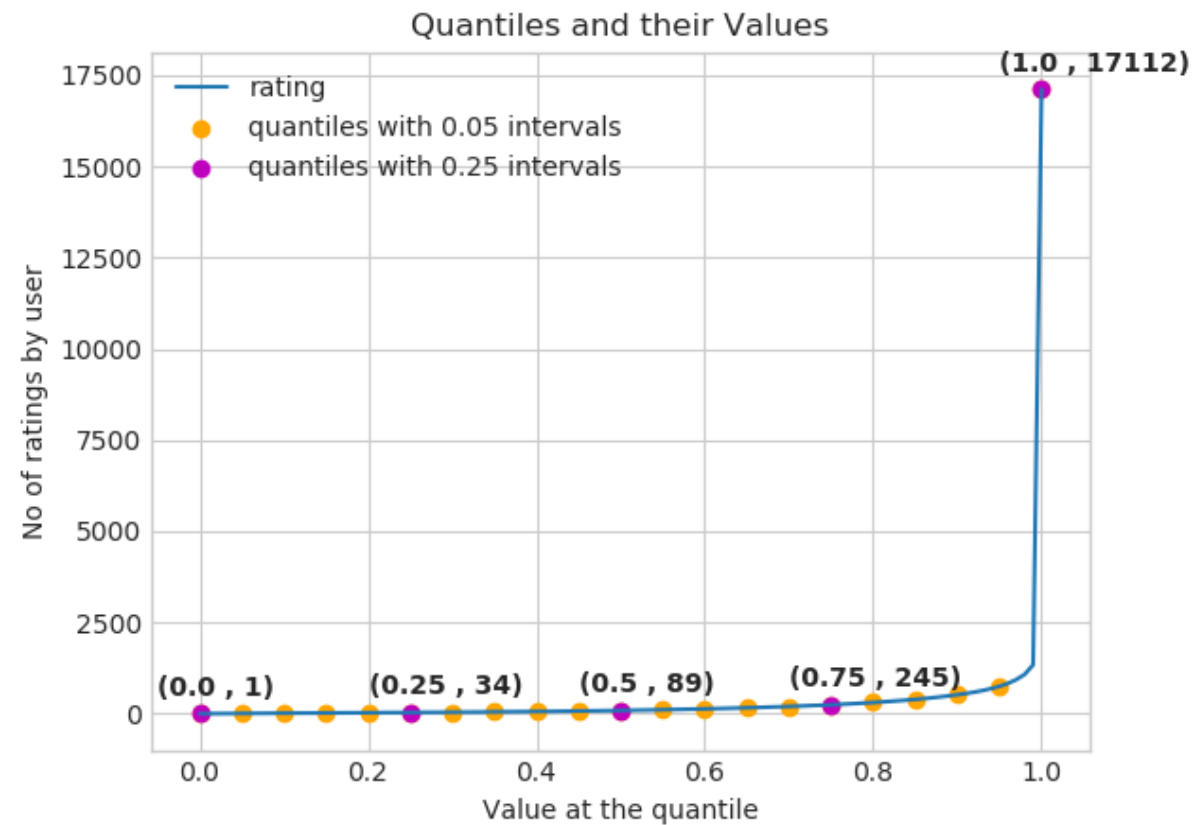
> There, is something interesting going on with the quantiles..

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

```python
In [0]: plt.title("Quantiles and their Values")
        quantiles.plot()
        # quantiles with 0.05 difference
        plt.scatter(x=quantiles.index[::5], y=quantiles.values[::5], c='orange'
        , label="quantiles with 0.05 intervals")
        # quantiles with 0.25 difference
        plt.scatter(x=quantiles.index[::25], y=quantiles.values[::25], c='m', l
        abel = "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()
```



Quantiles and their Values

```
In [0]:  quantiles[::5]

Out[0]:  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 [0]:
```python
print('\n No of ratings at last 5 percentile : {}\n'.format(sum(no_of_r
ated_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 [0]:
```python
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()
```

# RATINGS per Movie



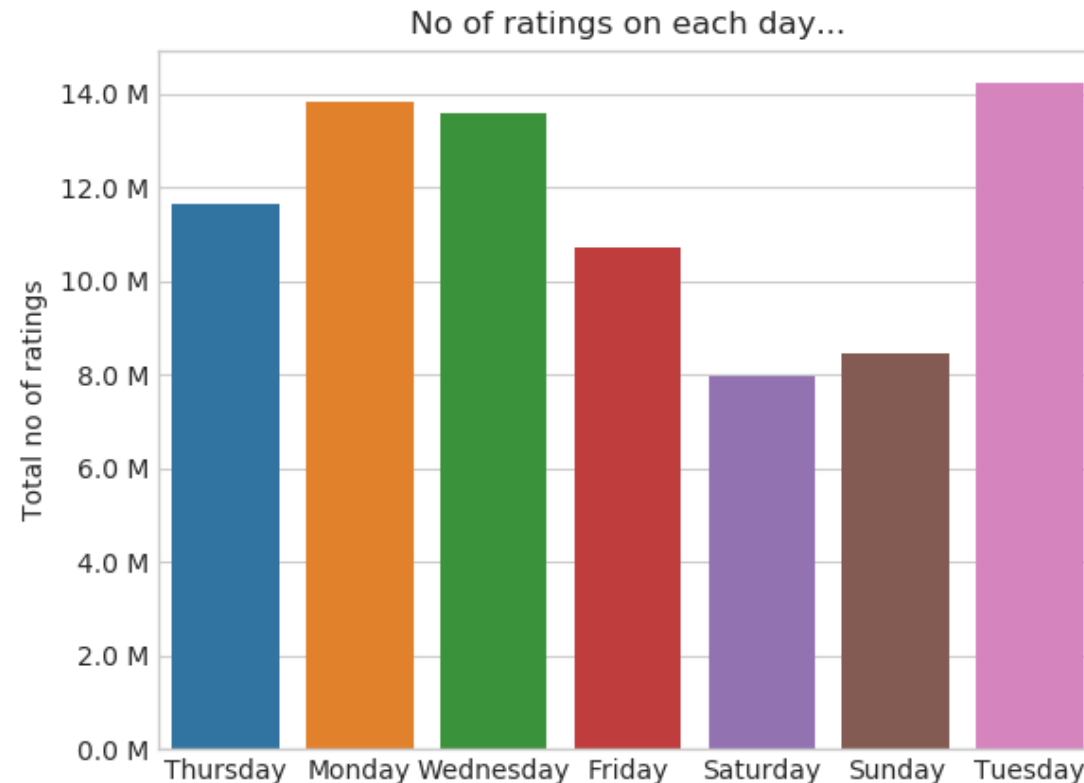- **It is very skewed.. just like nunmber of ratings given per user.**

```
- 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 hundereds of rating
s.
```

### 3.3.5 Number of ratings on each day of the week

```
In [0]:  fig, ax = plt.subplots()
         sns.countplot(x='day_of_week', data=train_df, ax=ax)
```

```
plt.title('No of ratings on each day...')
plt.ylabel('Total no of ratings')
plt.xlabel('')
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
plt.show()
```



No of ratings on each day...

In [0]:
```
start = datetime.now()
fig = plt.figure(figsize=plt.figaspect(.45))
sns.boxplot(y='rating', x='day_of_week', data=train_df)
plt.show()
print(datetime.now() - start)
```

```
0:01:10.003761
```

In [0]: 
```python
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 [0]:  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)
```

```
We are creating sparse_matrix from the dataframe..
Done. It's shape is : (user, movie) :  (2649430, 17771)
```

```
Saving it into disk for furthur usage..
Done..

0:01:13.804969
```

**The Sparsity of Train Sparse Matrix**

```
In [0]: us,mv = train_sparse_matrix.shape
        elem = train_sparse_matrix.count_nonzero()

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

```
Sparsity Of Train matrix : 99.8292709259195 %
```

### 3.3.6.2 Creating sparse matrix from test data frame

```
In [0]: 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, (tes
        t_df.user.values,
                                                    test_df.movie.values)))

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

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

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

```
We are creating sparse_matrix from the dataframe..
Done. It's shape is : (user, movie) :  (2649430, 17771)
Saving it into disk for furthur usage..
Done..

0:00:18.566120
```

**The Sparsity of Test data Matrix**

In [0]:
```
us,mv = test_sparse_matrix.shape
elem = test_sparse_matrix.count_nonzero()

print("Sparsity Of Test matrix : {} % ".format(  (1-(elem/(us*mv))) * 1
00) )
```

```
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 [0]:
```
# get the user averages in dictionary (key: user_id/movie_id, value: av
g 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 no
t)
    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
        # creae a dictonary of users and their average ratigns..
        average_ratings = { i : sum_of_ratings[i]/no_of_ratings[i]
                                        for i in range(u if of_users else m)
                                        if no_of_ratings[i] !=0}

        # return that dictionary of average ratings
        return average_ratings
```

### 3.3.7.1 finding global average of all movie ratings

```
In [0]: train_averages = dict()
        # get the global average of ratings in our train set.
        train_global_average = train_sparse_matrix.sum()/train_sparse_matrix.co
        unt_nonzero()
        train_averages['global'] = train_global_average
        train_averages
```

```
Out[0]: {'global': 3.582890686321557}
```

### 3.3.7.2 finding average rating per user

```
In [0]: train_averages['user'] = get_average_ratings(train_sparse_matrix, of_us
        ers=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 [0]: train_averages['movie'] =  get_average_ratings(train_sparse_matrix, of_
```

```
users=False)
print('\n AVerage rating of movie 15 :',train_averages['movie'][15])

 AVerage rating of movie 15 : 3.3038461538461537
```

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

In [0]:
```
start = datetime.now()
# draw pdfs for average rating per user and average
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(
.5))
fig.suptitle('Avg Ratings per User and per Movie', fontsize=15)

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

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

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

Avg Ratings per User and per Movie

```
0:00:35.003443
```

### 3.3.8 Cold Start problem

**3.3.8.1 Cold Start problem with Users**

```
In [0]: total_users = len(np.unique(df.user))
        users_train = len(train_averages['user'])
        new_users = total_users - users_train

        print('\nTotal number of Users  :', total_users)
        print('\nNumber of Users in Train data :', users_train)
        print("\nNo of Users that didn't appear in train data: {}({} %) \n ".fo
        rmat(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

```python
In [0]: 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 ".f
        ormat(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. usersbeing lare.

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

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

```
In [0]:
from sklearn.metrics.pairwise import cosine_similarity


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

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

    start = datetime.now()
    temp = 0
```

```python
    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()
```

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

In [0]:
```python
start = datetime.now()
u_u_sim_sparse, _ = compute_user_similarity(train_sparse_matrix, comput
e_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 : 0:03:20.300488  ]
computing done for 40 users [  time elapsed : 0:06:38.518391  ]
computing done for 60 users [  time elapsed : 0:09:53.143126  ]
computing done for 80 users [  time elapsed : 0:13:10.080447  ]
computing done for 100 users [  time elapsed : 0:16:24.711032  ]
Creating Sparse matrix from the computed similarities
```

```
-----------------------------------------------------------------------------
-------------------------------
Time taken : 0:16:33.618931
```

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

- We have **405,041 users** in out training set and computing similarities between them..( **17K dimensional vector..**) is time consuming..

- From above plot, It took roughly **8.88 sec** for computing simlilar users for **one user**

- We have **405,041 users** with us in training set.

- 
  $$405041 \times 8.88 = 3596764.08 \, \text{sec} = 59946.068 \, \text{min} = 999.101133333 \, \text{hours}$$
  $$= 41.629213889 \, \text{days} \ldots$$
    - 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 dimentsions using SVD, so that **it might** speed up the process...

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

start = datetime.now()

# initilaize the algorithm with some parameters..
# All of them are default except n_components. n_itr is for Randomized
 SVD solver.
netflix_svd = TruncatedSVD(n_components=500, algorithm='randomized', random_state=15)
trunc_svd = netflix_svd.fit_transform(train_sparse_matrix)

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

```
0:29:07.069783
```

Here,

- $\sum \longleftarrow$ (netflix_svd.**singular_values_** )

- $V^T \longleftarrow$ (netflix_svd.**components_)**

- $\bigcup$ 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 [0]:
```python
expl_var = np.cumsum(netflix_svd.explained_variance_ratio_)
```

In [0]:
```python
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)
# annote some (latentfactors, expl_var) to make it clear
ind = [1, 2,4,8,20, 60, 100, 200, 300, 400, 500]
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='bol
d')

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



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

plt.show()
```

```
In [0]: for i in ind:
            print("({}, {})".format(i, np.round(expl_var[i-1], 2)))
```

```
(1, 0.23)
(2, 0.26)
(4, 0.3)
(8, 0.34)
(20, 0.38)
(60, 0.44)
(100, 0.47)
(200, 0.53)
(300, 0.57)
(400, 0.61)
(500, 0.64)
```

> 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 [0]:
```python
# 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)
```

0:00:45.670265

In [0]:
```python
type(trunc_matrix), trunc_matrix.shape
```
Out[0]: (numpy.ndarray, (2649430, 500))

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

```
In [0]: if not os.path.isfile('trunc_sparse_matrix.npz'):
            # create that sparse 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 [0]: trunc_sparse_matrix.shape
```

```
Out[0]: (2649430, 500)
```

```
In [0]: start = datetime.now()
        trunc_u_u_sim_matrix, _ = compute_user_similarity(trunc_sparse_matrix,
        compute_for_few=True, top=50, verbose=True,
                                                  verb_for_n_rows=10)
        print("-"*50)
        print("time:",datetime.now()-start)
```

```
Computing top 50 similarities for each user..
computing done for 10 users [  time elapsed : 0:02:09.746324  ]
computing done for 20 users [  time elapsed : 0:04:16.017768  ]
computing done for 30 users [  time elapsed : 0:06:20.861163  ]
computing done for 40 users [  time elapsed : 0:08:24.933316  ]
computing done for 50 users [  time elapsed : 0:10:28.861485  ]
Creating Sparse matrix from the computed similarities
```

```
--------------------------------------------------
time: 0:10:52.658092
```

**: 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.

- 
  $405041 \times 12.18 ==== 4933399.38 \, \text{sec} ==== 82223.323 \, \text{min} ==== 137($
  $==== 57.099529861 \, \text{days...}$
  - Even we run on 4 cores parallelly (a typical system now a days), It will still take almost **(14 - 15)** days.

- **Why did this happen...??**

  - Just think about it. It's not that difficult.

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

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

-An alternative is to compute similar users for a particular user, whenenver 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 tha
    t information.
        - In production time, We might have to recompute similaritie
    s, if it is computed a long time ago. Because user preferences c
    hanges 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 Diction
aries**.
        -
        - **key     :** _userid_
        - __value__: _Again a dictionary_
            - __key__  : _Similar User_
            - __value__: _Similarity Value_

### 3.4.2 Computing Movie-Movie Similarity matrix

```
In [0]: start = datetime.now()
        if not os.path.isfile('m_m_sim_sparse.npz'):
            print("It seems you don't have that file. Computing movie_movie sim
        ilarity...")
            start = datetime.now()
            m_m_sim_sparse = cosine_similarity(X=train_sparse_matrix.T, dense_o
        utput=False)
            print("Done..")
            # store this sparse matrix in disk before using it. For future purp
        oses.
            print("Saving it to disk without the need of re-computing it agai
        n.. ")
            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
0:10:02.736054
```

In [0]: `m_m_sim_sparse.shape`

Out[0]: `(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 [0]:
```python
movie_ids = np.unique(m_m_sim_sparse.nonzero()[1])
```

In [0]:
```python
start = datetime.now()
similar_movies = dict()
for movie in movie_ids:
    # get the top similar movies and store them in the dictionary
    sim_movies = m_m_sim_sparse[movie].toarray().ravel().argsort()[::-1
][1:]
    similar_movies[movie] = sim_movies[:100]
print(datetime.now() - start)

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

```
0:00:33.411700
```

Out[0]:
```
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])
```

### 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....*

```python
In [0]:  # First Let's load the movie details into soe dataframe..
         # movie details are in 'netflix/movie_titles.csv'

         movie_titles = pd.read_csv("data_folder/movie_titles.csv", sep=',', hea
         der = None,
                                    names=['movie_id', 'year_of_release', 'titl
         e'], verbose=True,
                                    index_col = 'movie_id', encoding = "ISO-8859-1")

         movie_titles.head()
```

```
Tokenization took: 4.50 ms
Type conversion took: 165.72 ms
Parser memory cleanup took: 0.01 ms
```

Out[0]:

| movie_id | year_of_release | title |
|---|---|---|
| 1 | 2003.0 | Dinosaur Planet |

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

**Similar Movies for 'Vampire Journals'**

In [0]:
```python
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 on
ly 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 [0]:
```python
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 rever
se the array and ignore its similarity (ie.,1)
```

```
                                                          # and return its indices
(movie_ids)
```

In [0]:
```python
plt.plot(similarities[sim_indices], label='All the ratings')
plt.plot(similarities[sim_indices[:100]], label='top 100 similar movie
s')
plt.title("Similar Movies of {}(movie_id)".format(mv_id), fontsize=20)
plt.xlabel("Movies (Not Movie_Ids)", fontsize=15)
plt.ylabel("Cosine Similarity",fontsize=15)
plt.legend()
plt.show()
```

**Top 10 similar movies**

In [0]: `movie_titles.loc[sim_indices[:10]]`

Out[0]:

| movie_id | year_of_release | title |
|---|---|---|
| 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 [0]:  def get_sample_sparse_matrix(sparse_matrix, no_users, no_movies, path,
         verbose = True):
             """
                 It will get it from the ''path'' if it is present  or It will c
         reate
                 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(use
         rs), len(movies)))
             print("Original Matrix : Ratings -- {}\n".format(len(ratings)))

             # It just to make sure to get same sample everytime we run this pro
         gram..
             # 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_in
         ds..
             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[m
         ask], col_ind[mask])),
                                                      shape=(max(sample_users)+1
         , max(sample_movies)+1))

             if verbose:
```

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

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

    return sample_sparse_matrix
```

## 4.1 Sampling Data

### 4.1.1 Build sample train data from the train data

```
In [0]:  start = datetime.now()
         path = "sample/small/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 10k users and 1k movies from available data
             sample_train_sparse_matrix = get_sample_sparse_matrix(train_sparse_
         matrix, no_users=10000, no_movies=1000,
                                                           path = path)

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

```
It is present in your pwd, getting it from disk....
DONE..
0:00:00.035179
```

### 4.1.2 Build sample test data from the test data

```
In [0]: start = datetime.now()

path = "sample/small/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_ma
trix, no_users=5000, no_movies=500,
                                                        path = "sample/small/s
ample_test_sparse_matrix.npz")
print(datetime.now() - start)
```

```
It is present in your pwd, getting it from disk....
DONE..
0:00:00.028740
```

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

```
In [0]: sample_train_averages = dict()
```

### 4.2.1 Finding Global Average of all movie ratings

```
In [0]: # get the global average of ratings in our train set.
global_average = sample_train_sparse_matrix.sum()/sample_train_sparse_m
atrix.count_nonzero()
```

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

Out[0]: {'global': 3.581679377504138}

### 4.2.2 Finding Average rating per User

```
In [0]: sample_train_averages['user'] = get_average_ratings(sample_train_sparse
        _matrix, of_users=True)
        print('\nAverage rating of user 1515220 :',sample_train_averages['user'
        ][1515220])
```

Average rating of user 1515220 : 3.9655172413793105

### 4.2.3 Finding Average rating per Movie

```
In [0]: sample_train_averages['movie'] =  get_average_ratings(sample_train_spar
        se_matrix, of_users=False)
        print('\n AVerage rating of movie 15153 :',sample_train_averages['movi
        e'][15153])
```

AVerage rating of movie 15153 : 2.6458333333333335

## 4.3 Featurizing data

```
In [0]: print('\n No of ratings in Our Sampled train matrix is : {}\n'.format(s
        ample_train_sparse_matrix.count_nonzero()))
        print('\n No of ratings in Our Sampled test  matrix is : {}\n'.format(s
        ample_test_sparse_matrix.count_nonzero()))
```

No of ratings in Our Sampled train matrix is : 129286

No of ratings in Our Sampled test  matrix is : 7333

## 4.3.1 Featurizing data for regression problem

### 4.3.1.1 Featurizing train data

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

In [0]:

```python
############################################################
# It took me almost 10 hours to prepare this train dataset.#
############################################################
start = datetime.now()
if os.path.isfile('sample/small/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('sample/small/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()
```

```python
            # we will make it's length "5" by adding movie averages to
.
            top_sim_users_ratings = list(top_ratings[top_ratings != 0]
[:5])
            top_sim_users_ratings.extend([sample_train_averages['movie'
][movie]]*(5 - len(top_sim_users_ratings)))
        #     print(top_sim_users_ratings, end=" ")


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


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

```
                row.append(sample_train_averages['user'][user])
                # Avg_movie rating
                row.append(sample_train_averages['movie'][movie])

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

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


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

```
preparing 129286 tuples for the dataset..

Done for 10000 rows----- 0:53:13.974716
Done for 20000 rows----- 1:47:58.228942
Done for 30000 rows----- 2:42:46.963119
Done for 40000 rows----- 3:36:44.807894
Done for 50000 rows----- 4:28:55.311500
Done for 60000 rows----- 5:24:18.493104
Done for 70000 rows----- 6:17:39.669922
Done for 80000 rows----- 7:11:23.970879
Done for 90000 rows----- 8:05:33.787770
Done for 100000 rows----- 9:00:25.463562
Done for 110000 rows----- 9:51:28.530010
Done for 120000 rows----- 10:42:05.382141
11:30:13.699183
```

**Reading from the file to make a Train_dataframe**

In [5]:
```
reg_train = pd.read_csv('reg_train.csv', names = ['user', 'movie', 'GAv
```

```
g', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5','smr1', 'smr2', 'smr3', 'smr
4', 'smr5', 'UAvg', 'MAvg', 'rating'], header=None)
reg_train.head()
```

Out[5]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr2 | smr3 | smr4 | smr5 | U |
|---|--------|-------|----------|------|------|------|------|------|------|------|------|------|------|-------|
| 0 | 53406 | 33 | 3.581679 | 4.0 | 5.0 | 5.0 | 4.0 | 1.0 | 5.0 | 2.0 | 5.0 | 3.0 | 1.0 | 3.37( |
| 1 | 99540 | 33 | 3.581679 | 5.0 | 5.0 | 5.0 | 4.0 | 5.0 | 3.0 | 4.0 | 4.0 | 3.0 | 5.0 | 3.55! |
| 2 | 99865 | 33 | 3.581679 | 5.0 | 5.0 | 4.0 | 5.0 | 3.0 | 5.0 | 4.0 | 4.0 | 5.0 | 4.0 | 3.71< |
| 3 | 101620 | 33 | 3.581679 | 2.0 | 3.0 | 5.0 | 5.0 | 4.0 | 4.0 | 3.0 | 3.0 | 4.0 | 5.0 | 3.58< |
| 4 | 112974 | 33 | 3.581679 | 5.0 | 5.0 | 5.0 | 5.0 | 5.0 | 3.0 | 5.0 | 5.0 | 5.0 | 3.0 | 3.75( |

- **GAvg** : Average rating of all the ratings

- **Similar users rating of this movie**:
    - sur1, sur2, sur3, sur4, sur5 ( top 5 similar users who rated that movie.. )

- **Similar movies rated by this user**:
    - smr1, smr2, smr3, smr4, smr5 ( top 5 similar movies rated by this movie.. )

- **UAvg** : User's Average rating

- **MAvg** : Average rating of this movie

- **rating** : Rating of this movie by this user.

**4.3.1.2 Featurizing test data**

```
In [0]:  # get users, movies and ratings from the Sampled Test
         sample_test_users, sample_test_movies, sample_test_ratings = sparse.fin
         d(sample_test_sparse_matrix)
```

```
In [0]:  sample_train_averages['global']
```

```
Out[0]:  3.581679377504138
```

```
In [0]:  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_t
         est_ratings)))
             with open('sample/small/reg_test.csv', mode='w') as reg_data_file:
                 count = 0
                 for (user, movie, rating)  in zip(sample_test_users, sample_tes
         t_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 i
         gnoring '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['mo
```

```python
vie'][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['gl
obal']]*(5 - len(top_sim_users_ratings)))
            #print(top_sim_users_ratings)
        except:
            print(user, movie)
            # we just want KeyErrors to be resolved. Not every Exce
ption...
            raise


        #-------------------- Ratings by "user"  to similar movies
 of "movie" --------------------
        try:
            # compute the similar movies of the "movie"
            movie_sim = cosine_similarity(sample_train_sparse_matri
x[:,movie].T, sample_train_sparse_matrix.T).ravel()
            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 u
ser..
            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['u
ser'][user]]*(5-len(top_sim_movies_ratings)))
            #print(top_sim_movies_ratings)
        except (IndexError, KeyError):
            #print(top_sim_movies_ratings, end=" : -- ")
```

```python
                    top_sim_movies_ratings.extend([sample_train_averages['g
lobal']]*(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 7333 tuples for the dataset..

Done for 1000 rows----- 0:04:29.293783
Done for 2000 rows----- 0:08:57.208002
Done for 3000 rows----- 0:13:30.333223
Done for 4000 rows----- 0:18:04.050813
Done for 5000 rows----- 0:22:38.671673
Done for 6000 rows----- 0:27:09.697009
Done for 7000 rows----- 0:31:41.933568
 0:33:12.529731

**Reading from the file to make a test dataframe**

In [4]:
```
reg_test_df = pd.read_csv('reg_test.csv', names = ['user', 'movie', 'GA
vg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
                                                    'smr1', 'smr
2', 'smr3', 'smr4', 'smr5',
                                                    'UAvg', 'MAv
g', 'rating'], header=None)
reg_test_df.head(4)
```

Out[4]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | sm |
|---|--------|-------|----------|----------|----------|----------|----------|----------|----------|--------|
| **0** | 808635 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.5816 |
| **1** | 941866 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.5816 |

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | sm |
|---|---|---|---|---|---|---|---|---|---|---|
| **2** | 1737912 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.5816 |
| **3** | 1849204 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.5816 |

- **GAvg** : Average rating of all the ratings

- **Similar users rating of this movie**:
    - sur1, sur2, sur3, sur4, sur5 ( top 5 simiular users who rated that movie.. )

- **Similar movies rated by this user**:
    - smr1, smr2, smr3, smr4, smr5 ( top 5 simiular movies rated by this movie.. )

- **UAvg** : User AVerage rating

- **MAvg** : Average rating of this movie

- **rating** : Rating of this movie by this user.

### 4.3.2 Transforming data for Surprise models

In [7]: `!pip3 install Surprise`

```
Collecting Surprise
  Downloading https://files.pythonhosted.org/packages/61/de/e5cba868220
1fcf9c3719a6fdda95693468ed061945493dea2dd37c5618b/surprise-0.1-py2.py3-
none-any.whl
Collecting scikit-surprise (from Surprise)
  Downloading https://files.pythonhosted.org/packages/4d/fc/cd4210b247d
```

```
1dca421c25994740cbbf03c5e980e31881f10eaddf45fdab0/scikit-surprise-1.0.
6.tar.gz (3.3MB)
    100% |████████████████████████████████| 3.3MB 429kB/s eta 0:00:01
Collecting joblib>=0.11 (from scikit-surprise->Surprise)
  Downloading https://files.pythonhosted.org/packages/cd/c1/50a758e8247
561e58cb87305b1e90b171b8c767b15b12a1734001f41d356/joblib-0.13.2-py2.py3
-none-any.whl (278kB)
    100% |████████████████████████████████| 286kB 5.1MB/s eta 0:00:01
Requirement already satisfied: numpy>=1.11.2 in /usr/local/lib/python3.
5/site-packages (from scikit-surprise->Surprise)
Requirement already satisfied: scipy>=1.0.0 in /usr/local/lib/python3.
5/site-packages (from scikit-surprise->Surprise)
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.5/
dist-packages (from scikit-surprise->Surprise)
Requirement already satisfied: mkl-random in /usr/local/lib/python3.5/d
ist-packages (from numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: mkl-fft in /usr/local/lib/python3.5/dist
-packages (from numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: mkl in /usr/local/lib/python3.5/dist-pac
kages (from numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: icc-rt in /usr/local/lib/python3.5/dist-
packages (from numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: tbb4py in /usr/local/lib/python3.5/dist-
packages (from numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: intel-numpy in /usr/local/lib/python3.5/
dist-packages (from scipy>=1.0.0->scikit-surprise->Surprise)
Requirement already satisfied: intel-openmp in /usr/local/lib/python3.
5/dist-packages (from mkl->numpy>=1.11.2->scikit-surprise->Surprise)
Requirement already satisfied: tbb==2019.* in /usr/local/lib/python3.5/
dist-packages (from tbb4py->numpy>=1.11.2->scikit-surprise->Surprise)
Building wheels for collected packages: scikit-surprise
  Running setup.py bdist_wheel for scikit-surprise ... done
  Stored in directory: /home/rahulgupta743/.cache/pip/wheels/ec/c0/55/3
a28eab06b53c220015063ebbdb81213cd3dcbb72c088251ec
Successfully built scikit-surprise
Installing collected packages: joblib, scikit-surprise, Surprise
Successfully installed Surprise-0.1 joblib-0.13.2 scikit-surprise-1.0.6
```

```
In [6]: from surprise import Reader, Dataset
```

In [ ]:

#### 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 saperate format for TRAIN and TEST data, which will be useful for training the models like SVD, KNNBaseLineOnly....etc..,in Surprise.

- We can form the trainset from a file, or from a Pandas DataFrame.
  http://surprise.readthedocs.io/en/stable/getting_started.html#load-dom-dataframe-py

In [7]:
```python
# 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 sur
prise 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 impotant)

In [8]:
```python
testset = list(zip(reg_test_df.user.values, reg_test_df.movie.values, r
eg_test_df.rating.values))
testset[:3]
```

Out[8]: [(808635, 71, 5), (941866, 71, 4), (1737912, 71, 3)]

# 4.4 Applying Machine Learning models

- Global dictionary that stores rmse and mape for all the models....
    - It stores the metrics in a dictionary of dictionaries

> **keys** : model names(string)
>
> **value**: dict(**key** : metric, **value** : value )

```
In [9]:  models_evaluation_train = dict()
         models_evaluation_test = dict()

         models_evaluation_train, models_evaluation_test
```

Out[9]: ({}, {})

> **Utility functions for running regression models**

```
In [10]:  # 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

          ###################################################################
          ###################################################################
```

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

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


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

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

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


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

```python
                            '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 [11]:
```python
# 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):
```

```python
    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 tes
t 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 Surpri
se)
    train_preds = algo.test(trainset.build_testset())
    # get predicted ratings from the train predictions..
```

```python
    train_actual_ratings, train_pred_ratings = get_ratings(train_preds)
    # get ''rmse'' and ''mape'' from the train predictions.
    train_rmse, train_mape = get_errors(train_preds)
    print('time taken : {}'.format(datetime.now()-st))

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

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

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

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

```
        test['mape'] = test_mape
        test['predictions'] = test_pred_ratings

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

        # return two dictionaries train and test
        return train, test
```

### 4.4.1 XGBoost with initial 13 features

In [12]:
```
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']
```

In [18]:
```
from matplotlib import warnings
warnings.simplefilter('ignore')
```

In [19]:
```
from sklearn.model_selection import GridSearchCV
parameters = {"learning_rate":[0.05,0.10,0.15,0.20],"max_depth":[3,4,5,
6,8,10],"n_estimators":[100,500,1000,1500]}
xgb1 = xgb.XGBRegressor()
gcv = GridSearchCV(xgb1,parameters,verbose=10,n_jobs=-1)
gcv.fit(x_train,y_train)
print(gcv.cv_results_.keys())
```

```
Fitting 3 folds for each of 96 candidates, totalling 288 fits
[CV] n_estimators=100, learning_rate=0.05, max_depth=3 ...............
[CV] n_estimators=100, learning_rate=0.05, max_depth=3 ..............
[CV] n_estimators=100, learning_rate=0.05, max_depth=3 ..............
[CV] n_estimators=500, learning_rate=0.05, max_depth=3 ..............
```

```
[CV]  n_estimators=100, learning_rate=0.05, max_depth=3, score=0.333353
22802151734, total=   6.3s
[CV] n_estimators=500, learning_rate=0.05, max_depth=3 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=3, score=0.405567
92038444295, total=   6.4s
[CV] n_estimators=500, learning_rate=0.05, max_depth=3 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=3, score=0.354814
0906249332, total=   6.4s
[CV] n_estimators=1000, learning_rate=0.05, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=3, score=0.344718
84296840943, total=  33.7s
[CV] n_estimators=1000, learning_rate=0.05, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=3, score=0.413589
7134890881, total=  34.1s
[CV] n_estimators=1000, learning_rate=0.05, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=3, score=0.360156
88028524795, total=  34.2s
[CV] n_estimators=1500, learning_rate=0.05, max_depth=3 ..............
```

`[Parallel(n_jobs=-1)]: Done     5 tasks        | elapsed:    42.8s`

```
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=3, score=0.34533
26182836068, total= 1.2min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=3, score=0.41381
062154053105, total= 1.2min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=3, score=0.36131
192958343494, total= 1.3min
[CV] n_estimators=100, learning_rate=0.05, max_depth=4 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=4, score=0.338844
232679811, total=  10.7s
[CV] n_estimators=100, learning_rate=0.05, max_depth=4 ...............
```

`[Parallel(n_jobs=-1)]: Done  10 tasks        | elapsed:  2.2min`

```
[CV]  n_estimators=100, learning_rate=0.05, max_depth=4, score=0.410925
0442193061, total=  10.8s
[CV] n_estimators=100, learning_rate=0.05, max_depth=4 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=4, score=0.356877
```

```
16218578247, total=  10.6s
[CV] n_estimators=500, learning_rate=0.05, max_depth=4 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=3, score=0.34610
69668452934, total= 1.9min
[CV] n_estimators=500, learning_rate=0.05, max_depth=4 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=3, score=0.41422
682771972685, total= 1.9min
[CV] n_estimators=500, learning_rate=0.05, max_depth=4 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=4, score=0.346011
9810458636, total=  46.7s
[CV] n_estimators=1000, learning_rate=0.05, max_depth=4 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=4, score=0.414182
48571603367, total=  45.5s
[CV] n_estimators=1000, learning_rate=0.05, max_depth=4 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=3, score=0.36291
425110612857, total= 1.8min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=4 ..............
```

[Parallel(n_jobs=-1)]: Done  17 tasks       | elapsed:  3.8min

```
[CV]  n_estimators=500, learning_rate=0.05, max_depth=4, score=0.360041
4654163966, total=  45.1s
[CV] n_estimators=1500, learning_rate=0.05, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=4, score=0.34699
20430637169, total= 1.6min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=4, score=0.41426
30420680181, total= 1.6min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=4, score=0.36160
754676465146, total= 1.6min
[CV] n_estimators=100, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=5, score=0.341511
19744120784, total=  13.6s
[CV] n_estimators=100, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=5, score=0.413019
0481042858, total=  13.5s
[CV] n_estimators=100, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=5, score=0.357993
```

```
77911263736, total=  13.5s
[CV] n_estimators=500, learning_rate=0.05, max_depth=5 ...............
```

[Parallel(n_jobs=-1)]: Done  24 tasks       | elapsed:  6.3min

```
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=4, score=0.34696
284831231766, total= 2.5min
[CV] n_estimators=500, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=5, score=0.346450
334232057, total= 1.1min
[CV] n_estimators=500, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=4, score=0.41461
1742945127, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=4, score=0.36231
55940411314, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=5, score=0.414644
4939121347, total= 1.0min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=5, score=0.361603
868203877, total=  59.2s
[CV] n_estimators=1500, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=5, score=0.34671
533056899007, total= 2.1min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=5, score=0.36353
033058463363, total= 2.1min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=5, score=0.41641
043614803885, total= 2.1min
[CV] n_estimators=100, learning_rate=0.05, max_depth=6 ...............
```

[Parallel(n_jobs=-1)]: Done  33 tasks       | elapsed: 10.2min

```
[CV]  n_estimators=100, learning_rate=0.05, max_depth=6, score=0.343715
4558587961, total=  17.1s
[CV] n_estimators=100, learning_rate=0.05, max_depth=6 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=6, score=0.413675
92415289867, total=  17.1s
```

```
[CV] n_estimators=100, learning_rate=0.05, max_depth=6 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=6, score=0.359314
73351304793, total=  17.0s
[CV] n_estimators=500, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=5, score=0.34628
110023115066, total= 3.2min
[CV] n_estimators=500, learning_rate=0.05, max_depth=6 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=6, score=0.345466
1664214626, total= 1.3min
[CV] n_estimators=500, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=6, score=0.412869
5408516433, total= 1.3min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=5, score=0.41438
087923205985, total= 3.2min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=5, score=0.36214
134943546067, total= 3.3min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=6, score=0.365084
5989474299, total= 1.4min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=6 ..............
```

```
[Parallel(n_jobs=-1)]: Done  42 tasks      | elapsed: 14.0min
```

```
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=6, score=0.34658
45739779165, total= 2.5min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=6, score=0.41033
49688511494, total= 2.5min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=6 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=6, score=0.36461
97343010006, total= 2.5min
[CV] n_estimators=100, learning_rate=0.05, max_depth=8 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=8, score=0.343761
0899258416, total=  22.2s
[CV] n_estimators=100, learning_rate=0.05, max_depth=8 ...............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=8, score=0.412180
0054906912, total=  22.5s
[CV] n_estimators=100, learning_rate=0.05, max_depth=8 ..............
```

```
[CV]  n_estimators=100, learning_rate=0.05, max_depth=8, score=0.360041
08943986496, total=  25.2s
[CV] n_estimators=500, learning_rate=0.05, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=6, score=0.34432
94907322802, total= 3.7min
[CV] n_estimators=500, learning_rate=0.05, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=8, score=0.340624
8963697257, total= 1.9min
[CV] n_estimators=500, learning_rate=0.05, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=8, score=0.412366
4722084526, total= 2.0min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=8 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=6, score=0.40541
69432611414, total= 3.9min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=8 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=6, score=0.36166
88947043427, total= 4.0min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=8 ..............
```

[Parallel(n_jobs=-1)]: Done  53 tasks       | elapsed: 20.5min

```
[CV]  n_estimators=500, learning_rate=0.05, max_depth=8, score=0.363632
3134914382, total= 1.9min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=8, score=0.33770
450327947554, total= 3.8min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=8, score=0.40425
21432181291, total= 3.9min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=8, score=0.34916
076569309124, total= 3.8min
[CV] n_estimators=100, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=10, score=0.34281
28706251101, total=  32.1s
[CV] n_estimators=100, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=10, score=0.40696
18398435473, total=  33.0s
[CV] n_estimators=100, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.05, max_depth=10, score=0.35847
```

```
27725557177, total=   32.4s
[CV] n_estimators=500, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=8, score=0.32917
65701971641, total= 5.7min
[CV] n_estimators=500, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=10, score=0.33824
52074271849, total= 2.6min
[CV] n_estimators=500, learning_rate=0.05, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=8, score=0.39473
92506478245, total= 5.7min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=10 .............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=8, score=0.33687
05185308334, total= 5.8min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=10 .............
```

[Parallel(n_jobs=-1)]: Done  64 tasks       | elapsed: 30.8min

```
[CV]  n_estimators=500, learning_rate=0.05, max_depth=10, score=0.40057
691131594764, total= 2.7min
[CV] n_estimators=1000, learning_rate=0.05, max_depth=10 .............
[CV]  n_estimators=500, learning_rate=0.05, max_depth=10, score=0.34878
96449174335, total= 2.6min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=10, score=0.3211
836140202253, total= 5.1min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=10, score=0.3811
178110221838, total= 5.1min
[CV] n_estimators=1500, learning_rate=0.05, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.05, max_depth=10, score=0.3322
7366405980574, total= 5.2min
[CV] n_estimators=100, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=3, score=0.3419347
1193121566, total=   7.5s
[CV] n_estimators=100, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=3, score=0.4122896
5509922877, total=   7.1s
[CV] n_estimators=100, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=3, score=0.3583924
8349235514, total=   7.1s
```

```
[CV] n_estimators=500, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=3, score=0.3465708
5424365186, total=  34.2s
[CV] n_estimators=500, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=3, score=0.4140064
5631252536, total=  35.3s
[CV] n_estimators=500, learning_rate=0.1, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=3, score=0.3628629
698375066, total=  37.2s
[CV] n_estimators=1000, learning_rate=0.1, max_depth=3 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=3, score=0.347963
8547108983, total= 1.2min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=3 ...............
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=10, score=0.3089
0150487901324, total= 7.5min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=3 ...............

[Parallel(n_jobs=-1)]: Done  77 tasks      | elapsed: 40.2min

[CV]  n_estimators=1000, learning_rate=0.1, max_depth=3, score=0.414298
03360739925, total= 1.2min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=3 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=3, score=0.363627
01669969144, total= 1.2min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=3 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=3, score=0.349379
49733961104, total= 1.8min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=3 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=3, score=0.415385
22974120984, total= 1.9min
[CV] n_estimators=100, learning_rate=0.1, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=4, score=0.3444399
95008401, total=  10.7s
[CV] n_estimators=100, learning_rate=0.1, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=4, score=0.4134117
548702042, total=  10.5s
[CV] n_estimators=100, learning_rate=0.1, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=4, score=0.3589456
642000472, total=  11.0s
[CV] n_estimators=500, learning_rate=0.1, max_depth=4 ................
```

```
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=10, score=0.3693
356255615923, total= 7.6min
[CV] n_estimators=500, learning_rate=0.1, max_depth=4 ................
[CV]  n_estimators=1500, learning_rate=0.05, max_depth=10, score=0.3183
387472999252, total= 7.5min
[CV] n_estimators=500, learning_rate=0.1, max_depth=4 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=4, score=0.3478100
2141737727, total=  51.9s
[CV] n_estimators=1000, learning_rate=0.1, max_depth=4 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=4, score=0.4145986
955081724, total=  49.3s
[CV] n_estimators=1000, learning_rate=0.1, max_depth=4 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=3, score=0.364237
4004822251, total= 1.9min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=4 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=4, score=0.3627483
920747545, total=  49.6s
[CV] n_estimators=1500, learning_rate=0.1, max_depth=4 ...............

[Parallel(n_jobs=-1)]: Done  90 tasks      | elapsed: 45.6min

[CV]  n_estimators=1000, learning_rate=0.1, max_depth=4, score=0.348743
54795902063, total= 1.7min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=4 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=4, score=0.415867
3707367043, total= 1.6min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=4 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=4, score=0.363868
3401161712, total= 1.6min
[CV] n_estimators=100, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=5, score=0.3452675
368846406, total=  13.1s
[CV] n_estimators=100, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=5, score=0.4141746
28437734, total=  13.0s
[CV] n_estimators=100, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=5, score=0.3597030
7762526177, total=  13.1s
[CV] n_estimators=500, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=4, score=0.349521
0674624896, total= 2.5min
```

```
0074024090, total= 2.5min
[CV] n_estimators=500, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=5, score=0.3506374
996592181, total= 1.1min
[CV] n_estimators=500, learning_rate=0.1, max_depth=5 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=5, score=0.4146430
313820815, total= 1.1min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=4, score=0.412870
206538561, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=4, score=0.359980
8915170155, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=5, score=0.3658289
125927263, total= 1.1min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=5, score=0.347884
66800963486, total= 2.1min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=5, score=0.409117
15182256, total= 2.1min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=5, score=0.360721
8934669131, total= 2.1min
[CV] n_estimators=100, learning_rate=0.1, max_depth=6 ................

[Parallel(n_jobs=-1)]: Done 105 tasks       | elapsed: 52.0min

[CV]  n_estimators=100, learning_rate=0.1, max_depth=6, score=0.3444913
0974463427, total=  15.2s
[CV] n_estimators=100, learning_rate=0.1, max_depth=6 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=6, score=0.4133796
44650533, total=  14.7s
[CV] n_estimators=100, learning_rate=0.1, max_depth=6 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=6, score=0.3614108
872084496, total=  14.8s
[CV] n_estimators=500, learning_rate=0.1, max_depth=6 ................
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=5, score=0.345634
3307743631, total= 3.1min
[CV] n_estimators=500, learning_rate=0.1, max_depth=6 ................
```

```
[CV]  n_estimators=500, learning_rate=0.1, max_depth=6 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=6, score=0.3463826
348621443, total= 1.4min
[CV] n_estimators=500, learning_rate=0.1, max_depth=6 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=6, score=0.4110498
253671475, total= 1.4min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=5, score=0.354934
7307324418, total= 3.1min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=5, score=0.402904
76851693446, total= 3.1min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=6, score=0.3613433
3256864454, total= 1.4min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=6, score=0.340449
8443151232, total= 2.5min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=6, score=0.350573
39483195327, total= 2.5min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=6, score=0.400417
4944174106, total= 2.5min
[CV] n_estimators=100, learning_rate=0.1, max_depth=8 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=8, score=0.3422171
205027527, total=  24.2s
[CV] n_estimators=100, learning_rate=0.1, max_depth=8 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=8, score=0.4112657
2673861805, total=  24.1s
[CV] n_estimators=100, learning_rate=0.1, max_depth=8 ................
[CV]  n_estimators=100, learning_rate=0.1, max_depth=8, score=0.3617429
6411749707, total=  24.0s
[CV] n_estimators=500, learning_rate=0.1, max_depth=8 ................

[Parallel(n_jobs=-1)]: Done 120 tasks      | elapsed: 59.0min

[CV]  n_estimators=1500, learning_rate=0.1, max_depth=6, score=0.331759
16364285307, total= 3.9min
[CV] n_estimators=500, learning_rate=0.1, max_depth=8 ................
[CV]  n_estimators=500, learning_rate=0.1, max_depth=8, score=0.3359058
```

```
                    3738196855, total= 1.9min
[CV] n_estimators=500, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=6, score=0.388241
86208395123, total= 4.1min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=8, score=0.3975330
334946584, total= 2.0min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=6, score=0.337747
38322887066, total= 4.1min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=8, score=0.3491818
7560738717, total= 2.0min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=8, score=0.316322
2227752325, total= 3.7min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=8, score=0.377262
79450187766, total= 3.7min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=8 ...............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=8, score=0.323260
8101631389, total= 3.8min
[CV] n_estimators=100, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.1, max_depth=10, score=0.336307
70051724046, total=  31.6s
[CV] n_estimators=100, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.1, max_depth=10, score=0.402368
1367638654, total=  31.6s
[CV] n_estimators=100, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.1, max_depth=10, score=0.352966
87493830814, total=  31.6s
[CV] n_estimators=500, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=8, score=0.299189
4930268939, total= 5.7min
[CV] n_estimators=500, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=10, score=0.316356
90235715386, total= 2.6min
[CV] n_estimators=500, learning_rate=0.1, max_depth=10 ...............
[CV]  n_estimators=500, learning_rate=0.1, max_depth=10, score=0.376489
```

```
76414835256, total= 2.6min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=8, score=0.303327
8170070792, total= 5.8min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=8, score=0.359567
66343075164, total= 5.9min
[CV] n_estimators=1000, learning_rate=0.1, max_depth=10 ..............
```

[Parallel(n_jobs=-1)]: Done 137 tasks        | elapsed: 72.4min

```
[CV]  n_estimators=500, learning_rate=0.1, max_depth=10, score=0.324626
6634481111, total= 2.6min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=10, score=0.29222
675684464205, total= 5.0min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=10, score=0.35047
36379984983, total= 5.0min
[CV] n_estimators=1500, learning_rate=0.1, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.1, max_depth=10, score=0.29633
208948972667, total= 5.0min
[CV] n_estimators=100, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=3, score=0.343944
91665406746, total=   7.8s
[CV] n_estimators=100, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=3, score=0.411526
4305753069, total=   7.9s
[CV] n_estimators=100, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=3, score=0.357840
862708642, total=   7.8s
[CV] n_estimators=500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=3, score=0.347681
8771264861, total=  37.8s
[CV] n_estimators=500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=3, score=0.414536
2166910299, total=  37.9s
[CV] n_estimators=500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=3, score=0.361501
6300402472, total=  38.2s
```

```
63004053472, total=  38.2s
[CV] n_estimators=1000, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=3, score=0.34988
813989406375, total= 1.3min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=10, score=0.27759
522056860564, total= 7.7min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=3, score=0.41458
263679874985, total= 1.2min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=3, score=0.36273
579076175255, total= 1.3min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=3, score=0.35151
55449927043, total= 1.7min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=3 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=3, score=0.41321
546572996004, total= 1.7min
[CV] n_estimators=100, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=4, score=0.343960
8004362013, total=  10.7s
[CV] n_estimators=100, learning_rate=0.15, max_depth=4 ..............
```

[Parallel(n_jobs=-1)]: Done 154 tasks      | elapsed: 85.1min

```
[CV]  n_estimators=100, learning_rate=0.15, max_depth=4, score=0.412642
6254019465, total=  11.4s
[CV] n_estimators=100, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=4, score=0.358833
7338258737, total=  11.0s
[CV] n_estimators=500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=10, score=0.28253
64436715787, total= 7.5min
[CV] n_estimators=500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1500, learning_rate=0.1, max_depth=10, score=0.33819
046447397766, total= 7.7min
[CV] n_estimators=500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=4, score=0.348098
8718798198, total=  51.3s
```

```
[CV] n_estimators=1000, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=3, score=0.36379
420695380604, total= 2.0min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=4, score=0.415055
9044386498, total=  49.6s
[CV] n_estimators=1000, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=4, score=0.361073
96600215186, total=  50.9s
[CV] n_estimators=1500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=4, score=0.34881
06726269264, total= 1.6min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=4, score=0.41226
49026972588, total= 1.7min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=4 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=4, score=0.36145
60817988568, total= 1.7min
[CV] n_estimators=100, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=5, score=0.343698
80908067385, total=  13.1s
[CV] n_estimators=100, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=5, score=0.412142
79672842186, total=  13.3s
[CV] n_estimators=100, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=5, score=0.360009
98351828056, total=  13.1s
[CV] n_estimators=500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=4, score=0.34603
86968850413, total= 2.5min
[CV] n_estimators=500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=5, score=0.348764
27132518784, total= 1.1min
[CV] n_estimators=500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=4, score=0.40625
318015153655, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=5, score=0.412293
0121352256, total= 1.1min
```

```
[CV] n_estimators=1000, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=4, score=0.35582
87505545471, total= 2.5min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=5 ..............

[Parallel(n_jobs=-1)]: Done 173 tasks       | elapsed: 91.3min

[CV]  n_estimators=500, learning_rate=0.15, max_depth=5, score=0.361717
1902301561, total= 1.1min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=5, score=0.33898
58445133438, total= 2.2min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=5, score=0.40172
950229639337, total= 2.2min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=5 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=5, score=0.35189
34210475866, total= 2.2min
[CV] n_estimators=100, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=6, score=0.346432
5240794279, total=  16.6s
[CV] n_estimators=100, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=6, score=0.412711
91745002067, total=  16.4s
[CV] n_estimators=100, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=6, score=0.361419
7325039229, total=  16.4s
[CV] n_estimators=500, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=5, score=0.32957
55269180717, total= 3.3min
[CV] n_estimators=500, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=6, score=0.347071
09671119235, total= 1.4min
[CV] n_estimators=500, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=6, score=0.404380
55975865306, total= 1.4min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=6 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=5, score=0.38906
84632624111, total= 3.2min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=6 ..............
```

```
[CV]   n_estimators=1500, learning_rate=0.15, max_depth=5, score=0.33482
32176644037, total= 3.3min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=6 ..............
[CV]   n_estimators=500, learning_rate=0.15, max_depth=6, score=0.357543
39538004315, total= 1.4min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=6 ..............
[CV]   n_estimators=1000, learning_rate=0.15, max_depth=6, score=0.33135
127867609215, total= 2.7min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=6 ..............
[CV]   n_estimators=1000, learning_rate=0.15, max_depth=6, score=0.38602
705261848275, total= 2.7min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=6 ..............
[CV]   n_estimators=1000, learning_rate=0.15, max_depth=6, score=0.33628
698291610504, total= 2.7min
[CV] n_estimators=100, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=100, learning_rate=0.15, max_depth=8, score=0.341572
44344036286, total=  24.7s
[CV] n_estimators=100, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=100, learning_rate=0.15, max_depth=8, score=0.407699
11497097355, total=  24.6s
[CV] n_estimators=100, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=100, learning_rate=0.15, max_depth=8, score=0.359947
2253560434, total=  25.0s
[CV] n_estimators=500, learning_rate=0.15, max_depth=8 ..............

[Parallel(n_jobs=-1)]: Done 192 tasks      | elapsed: 101.1min

[CV]   n_estimators=1500, learning_rate=0.15, max_depth=6, score=0.31507
30868581575, total= 4.2min
[CV] n_estimators=500, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=500, learning_rate=0.15, max_depth=8, score=0.322437
40390712516, total= 2.0min
[CV] n_estimators=500, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=500, learning_rate=0.15, max_depth=8, score=0.377611
37394739847, total= 2.0min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=8 ..............
[CV]   n_estimators=1500, learning_rate=0.15, max_depth=6, score=0.37031
235712574384, total= 4.3min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=8 ..............

[CV]   n_estimators=1500, learning_rate=0.15, max_depth=6, score=0.31536
```

```
61904130539, total= 4.3min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=8 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=8, score=0.329937
3507525467, total= 2.0min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=8, score=0.29032
31809228234, total= 3.8min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=8, score=0.34932
88652915657, total= 3.9min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=8 ..............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=8, score=0.29746
07044416976, total= 3.9min
[CV] n_estimators=100, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=10, score=0.32889
983838857495, total=  32.2s
[CV] n_estimators=100, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=10, score=0.39642
239960089987, total=  32.2s
[CV] n_estimators=100, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=100, learning_rate=0.15, max_depth=10, score=0.34885
23494764998, total=  32.1s
[CV] n_estimators=500, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=8, score=0.27167
967587877784, total= 5.9min
[CV] n_estimators=500, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=10, score=0.29267
667135916775, total= 2.7min
[CV] n_estimators=500, learning_rate=0.15, max_depth=10 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=8, score=0.33105
800140523417, total= 6.1min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=10 .............
[CV]  n_estimators=500, learning_rate=0.15, max_depth=10, score=0.35111
249320404686, total= 2.7min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=10 .............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=8, score=0.27700
113899149825, total= 6.1min
[CV] n_estimators=1000, learning_rate=0.15, max_depth=10 .............

[CV]  n_estimators=500, learning_rate=0.15, max_depth=10, score=0.30224
```

```
736778057326, total= 2.7min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=10, score=0.2688
1257951300164, total= 5.2min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=10, score=0.3284
4970722941763, total= 5.2min
[CV] n_estimators=1500, learning_rate=0.15, max_depth=10 .............
[CV]  n_estimators=1000, learning_rate=0.15, max_depth=10, score=0.2765
153269100896, total= 5.1min
[CV] n_estimators=100, learning_rate=0.2, max_depth=3 ................
```

[Parallel(n_jobs=-1)]: Done 213 tasks       | elapsed: 120.5min

```
[CV]  n_estimators=100, learning_rate=0.2, max_depth=3, score=0.3434643
331671196, total=   8.0s
[CV] n_estimators=100, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=3, score=0.4131715
476682637, total=   8.1s
[CV] n_estimators=100, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=3, score=0.3580435
239397376, total=   8.0s
[CV] n_estimators=500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=3, score=0.3465935
7267954327, total=   39.7s
[CV] n_estimators=500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=3, score=0.4156307
6994522985, total=   39.4s
[CV] n_estimators=500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=3, score=0.3618298
2036446665, total=   39.3s
[CV] n_estimators=1000, learning_rate=0.2, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=3, score=0.349616
6195742141, total= 1.4min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=3 ..............
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=10, score=0.2572
204471540882, total= 7.8min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=3 ..............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=3, score=0.414911
57874813456, total= 1.3min
```

```
5787401545U, total= 1.5min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=3, score=0.362410
9431606237, total= 1.3min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=3, score=0.350453
06441068064, total= 2.0min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=3 ................
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=3, score=0.413093
2473480048, total= 2.0min
[CV] n_estimators=100, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=4, score=0.3440374
077014339, total=  10.6s
[CV] n_estimators=100, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=4, score=0.4133492
489739212, total=  11.0s
[CV] n_estimators=100, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=4, score=0.3589101
107135173, total=  10.9s
[CV] n_estimators=500, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=10, score=0.3193
1257616781716, total= 8.1min
[CV] n_estimators=500, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=4, score=0.3473140
0061032336, total=  47.0s
[CV] n_estimators=500, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=1500, learning_rate=0.15, max_depth=10, score=0.2645
272611911631, total= 8.2min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=3, score=0.362145
5787798973, total= 1.8min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=4, score=0.4135120
435079032, total=  48.7s
[CV] n_estimators=1000, learning_rate=0.2, max_depth=4 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=4, score=0.3614920
423630763, total=  48.8s
[CV] n_estimators=1500, learning_rate=0.2, max_depth=4 ................

[Parallel(n_jobs=-1)]: Done 234 tasks      | elapsed: 130.3min
```

```
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=4, score=0.346636
429284927, total= 1.7min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=4 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=4, score=0.407212
8136176816, total= 1.7min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=4 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=4, score=0.357760
4982748249, total= 1.8min
[CV] n_estimators=100, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=5, score=0.3458464
767108342, total=  13.6s
[CV] n_estimators=100, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=5, score=0.4116466
256396878, total=  13.6s
[CV] n_estimators=100, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=5, score=0.3583902
505172691, total=  13.8s
[CV] n_estimators=500, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=4, score=0.340328
21325837237, total= 2.6min
[CV] n_estimators=500, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=5, score=0.3449482
8134062616, total= 1.1min
[CV] n_estimators=500, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=4, score=0.395498
4894256317, total= 2.6min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=5, score=0.4056392
853430516, total= 1.1min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=4, score=0.351251
44111204876, total= 2.7min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=5, score=0.3533698
409245997, total= 1.1min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=5, score=0.327341
3618026658, total= 2.2min

[CV] n_estimators=1500, learning_rate=0.2, max_depth=5 ...............
```

```
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=5, score=0.390616
3396812037, total= 2.2min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=5 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=5, score=0.339225
96928025006, total= 2.3min
[CV] n_estimators=100, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=6, score=0.3432332
490414478, total=  17.8s
[CV] n_estimators=100, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=6, score=0.4097018
82492607, total=  17.2s
[CV] n_estimators=100, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=100, learning_rate=0.2, max_depth=6, score=0.3631166
643914978, total=  17.1s
[CV] n_estimators=500, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=5, score=0.313730
6083085256, total= 3.4min
[CV] n_estimators=500, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=500, learning_rate=0.2, max_depth=6, score=0.3295526
9453472535, total= 1.3min
[CV] n_estimators=500, learning_rate=0.2, max_depth=6 ................
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=5, score=0.373818
29933369326, total= 3.4min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=6 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=5, score=0.308648
92436233704, total= 3.3min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=6 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=6, score=0.3948248
212949502, total= 1.4min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=6 ...............
```

[Parallel(n_jobs=-1)]: Done 257 tasks      | elapsed: 140.2min

```
[CV]  n_estimators=500, learning_rate=0.2, max_depth=6, score=0.3422534
8407968104, total= 1.4min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=6, score=0.302711
0499878972, total= 2.7min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=6, score=0.314234
```

```
84752702424, total= 2.6min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=6 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=6, score=0.367892
53685968304, total= 2.7min
[CV] n_estimators=100, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=8, score=0.3388341
157699469, total=  23.4s
[CV] n_estimators=100, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=8, score=0.4005832
0857697466, total=  23.3s
[CV] n_estimators=100, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=8, score=0.3558593
2518251684, total=  23.4s
[CV] n_estimators=500, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=6, score=0.283380
43798330304, total= 4.0min
[CV] n_estimators=500, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=8, score=0.2994033
73745741, total= 2.0min
[CV] n_estimators=500, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=8, score=0.3523937
7742859973, total= 1.9min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=6, score=0.346334
40426500883, total= 4.0min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=6, score=0.292601
90478788384, total= 4.0min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=8, score=0.3126777
831843853, total= 2.0min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=8, score=0.265231
05288897386, total= 3.9min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=8 ...............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=8, score=0.320697
1201343519, total= 3.8min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=8 ...............

[CV]  n_estimators=1000, learning_rate=0.2, max_depth=8, score=0.271537
```

```
6422044933, total= 3.9min
[CV] n_estimators=100, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=10, score=0.321218
7999376255, total=  31.2s
[CV] n_estimators=100, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=10, score=0.380580
1418731807, total=  30.8s
[CV] n_estimators=100, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=100, learning_rate=0.2, max_depth=10, score=0.335072
6998119158, total=  30.4s
[CV] n_estimators=500, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=8, score=0.243584
53798369228, total= 5.8min
[CV] n_estimators=500, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=10, score=0.267839
0224108287, total= 2.6min
[CV] n_estimators=500, learning_rate=0.2, max_depth=10 ...............
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=8, score=0.300903
3383383126, total= 5.7min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=10 ..............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=10, score=0.328077
1040305551, total= 2.6min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=10 ..............
```

[Parallel(n_jobs=-1)]: Done 280 tasks       | elapsed: 157.5min

```
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=8, score=0.247856
3731146498, total= 5.7min
[CV] n_estimators=1000, learning_rate=0.2, max_depth=10 ..............
[CV]  n_estimators=500, learning_rate=0.2, max_depth=10, score=0.276051
82965959474, total= 2.4min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=10, score=0.24425
880252324217, total= 4.8min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=10, score=0.30710
14505879489, total= 4.8min
[CV] n_estimators=1500, learning_rate=0.2, max_depth=10 ..............
[CV]  n_estimators=1000, learning_rate=0.2, max_depth=10, score=0.25143
3183822643, total= 4.9min
[CV]  n_estimators=1500, learning_rate=0.2, max_depth=10, score=0.33577
```

```
[CV]   n_estimators=1500, learning_rate=0.2, max_depth=10, score=0.23577
859699474235, total= 7.0min
[CV]   n_estimators=1500, learning_rate=0.2, max_depth=10, score=0.30060
131651399424, total= 5.5min
[CV]   n_estimators=1500, learning_rate=0.2, max_depth=10, score=0.24336
511159415464, total= 5.7min
```

```
[Parallel(n_jobs=-1)]: Done 288 out of 288 | elapsed: 169.1min finished
```

```
dict_keys(['param_learning_rate', 'split0_test_score', 'params', 'split
1_train_score', 'split1_test_score', 'param_max_depth', 'param_n_estima
tors', 'mean_test_score', 'split2_test_score', 'mean_fit_time', 'mean_s
core_time', 'mean_train_score', 'split2_train_score', 'split0_train_sco
re', 'rank_test_score', 'std_fit_time', 'std_test_score', 'std_score_ti
me', 'std_train_score'])
```

In [14]:
```python
import xgboost as xgb
```

In [24]:
```python
print(gcv.best_params_)
```

```
{'n_estimators': 500, 'learning_rate': 0.1, 'max_depth': 5}
```

In [25]:
```python
# prepare Train data
x_train = reg_train.drop(['user','movie','rating'], axis=1)
y_train = reg_train['rating']

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

# initialize Our first XGBoost model...
first_xgb = xgb.XGBRegressor(silent=False, n_jobs=13, random_state=15,
n_estimators=500,learning_rate=0.1,max_depth=5)
train_results, test_results = run_xgboost(first_xgb, x_train, y_train,
x_test, y_test)

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

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

```
Training the model..
Done. Time taken : 0:01:06.333358

Done

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

TEST DATA
------------------------------
RMSE :   1.1017397109523788
MAPE :   33.355770244225134
```

### 4.4.2 Suprise BaselineModel

```
In [26]:   from surprise import BaselineOnly
```

**Predicted_rating : ( baseline prediction )**

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

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

- $\boldsymbol{\mu}$ : Average of all trainings in training data.
- $\boldsymbol{b}_u$ : User bias
- $\boldsymbol{b}_i$ : Item bias (movie biases)

**Optimization function ( Least Squares Problem )**

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

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

In [28]:
```python
# options are to specify.., how to compute those user and item biases
bsl_options = {'method': 'sgd',
               'learning_rate': .001
              }
my_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(my_bsl_algo, trainse
t, testset, verbose=True)


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

```
Training the model...
Estimating biases using sgd...
Done. time taken : 0:00:00.943969

Evaluating the model with train data..
time taken : 0:00:01.170131
---------------
Train Data
---------------
RMSE : 0.9347153928678286

MAPE : 29.389572652358183

adding train results in the dictionary..
```

```
Evaluating for test data...
time taken : 0:00:00.082818
---------------
Test Data
---------------
RMSE : 1.0730330260516174

MAPE : 35.04995544572911

storing the test results in test dictionary...

----------------------------------------------
Total time taken to run this algorithm : 0:00:02.197859
```

### 4.4.3 XGBoost with initial 13 features + Surprise Baseline predictor

**Updating Train Data**

In [29]:
```python
# add our baseline_predicted value as our feature..
reg_train['bslpr'] = models_evaluation_train['bsl_algo']['predictions']
reg_train.head(2)
```

Out[29]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr2 | smr3 | smr4 | smr5 | UA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 53406 | 33 | 3.581679 | 4.0 | 5.0 | 5.0 | 4.0 | 1.0 | 5.0 | 2.0 | 5.0 | 3.0 | 1.0 | 3.370: |
| **1** | 99540 | 33 | 3.581679 | 5.0 | 5.0 | 5.0 | 4.0 | 5.0 | 3.0 | 4.0 | 4.0 | 3.0 | 5.0 | 3.555! |

**Updating Test Data**

In [30]:
```python
# add that baseline predicted ratings with Surprise to the test data as
 well
reg_test_df['bslpr']  = models_evaluation_test['bsl_algo']['prediction
```

```
                            s']

                            reg_test_df.head(2)
```

Out[30]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr |
|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 808635 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |
| **1** | 941866 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |

In [31]:
```python
# prepare train data
x_train = reg_train.drop(['user', 'movie','rating'], axis=1)
y_train = reg_train['rating']

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

# initialize Our first XGBoost model...
xgb_bsl = xgb.XGBRegressor(silent=False, n_jobs=13, random_state=15,n_e
stimators=500,learning_rate=0.1,max_depth=5)
train_results, test_results = run_xgboost(xgb_bsl, x_train, y_train, x_
test, y_test)

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

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

```
Training the model..
Done. Time taken : 0:01:30.942400

Done

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

```
TEST DATA
-------------------------------
RMSE :  1.0831711919892506
MAPE :  34.0777117877425
```

### 4.4.4 Surprise KNNBaseline predictor

In [32]:
```python
from surprise import KNNBaseline
```

- KNN BASELINE
  - http://surprise.readthedocs.io/en/stable/knn_inspired.html#surprise.prediction_algorithms.

- PEARSON_BASELINE SIMILARITY
  - http://surprise.readthedocs.io/en/stable/similarities.html#surprise.similarities.pearson_bas

- SHRINKAGE
  - *2.2 Neighborhood Models* in http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf

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

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

- $b_{ui}$ - *Baseline prediction* of (user,movie) rating
- $N_i^k(u)$ - Set of **K similar** users (neighbours) of **user (u)** who rated **movie(i)**
- *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 ):

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

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

**4.4.4.1 Surprise KNNBaseline with user user similarities**

In [33]:
```python
# we specify , how to compute similarities and what to consider with si
m_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_r
ate 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)
```

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

```
Training the model...
Estimating biases using sgd...
Computing the pearson_baseline similarity matrix...
Done computing similarity matrix.
Done. time taken : 0:00:37.370540

Evaluating the model with train data..
time taken : 0:01:36.083696
---------------
Train Data
---------------
RMSE : 0.33642097416508826

MAPE : 9.145093375416348

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.088073
---------------
Test Data
---------------
RMSE : 1.0726493739667242

MAPE : 35.02094499698424

storing the test results in test dictionary...

------------------------------------------------
Total time taken to run this algorithm : 0:02:13.543819
```

**4.4.4.2 Surprise KNNBaseline with movie movie similarities**

```python
# we specify , how to compute similarities and what to consider with si
m_options to our algorithm

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

sim_options = {'user_based' : False,
               'name': 'pearson_baseline',
               'shrinkage': 100,
               'min_support': 2
              }
# we keep other parameters like regularization parameter and learning_r
ate 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 : 0:00:01.416301

Evaluating the model with train data..
time taken : 0:00:08.674137
--------------
Train Data
--------------
RMSE : 0.32584796251610554

MAPE : 8.447062581998374
```

```
adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.085992
---------------
Test Data
---------------
RMSE : 1.072758832653683

MAPE : 35.02269653015042

storing the test results in test dictionary...

-----------------------------------------------
Total time taken to run this algorithm : 0:00:10.178321
```

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

- - o First we will run XGBoost with predictions from both KNN's ( that uses User_User and Item_Item similarities along with our previous features.

- - o Then we will run XGBoost with just predictions form both knn models and preditions from our baseline model.

**Preparing Train data**

```
In [35]:  # add the predicted values from both knns to this dataframe
          reg_train['knn_bsl_u'] = models_evaluation_train['knn_bsl_u']['predicti
          ons']
          reg_train['knn_bsl_m'] = models_evaluation_train['knn_bsl_m']['predicti
          ons']
```

```
reg_train.head(2)
```

Out[35]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr2 | smr3 | smr4 | smr5 | U/ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 53406 | 33 | 3.581679 | 4.0 | 5.0 | 5.0 | 4.0 | 1.0 | 5.0 | 2.0 | 5.0 | 3.0 | 1.0 | 3.370: |
| 1 | 99540 | 33 | 3.581679 | 5.0 | 5.0 | 5.0 | 4.0 | 5.0 | 3.0 | 4.0 | 4.0 | 3.0 | 5.0 | 3.5555 |

◄ ▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬▬ ►

**Preparing Test data**

In [36]:
```
reg_test_df['knn_bsl_u'] = models_evaluation_test['knn_bsl_u']['predict
ions']
reg_test_df['knn_bsl_m'] = models_evaluation_test['knn_bsl_m']['predict
ions']

reg_test_df.head(2)
```

Out[36]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 808635 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |
| 1 | 941866 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |

◄ ▬▬▬▬▬▬▬▬▬▬▬ ►

In [37]:
```
# prepare the train data....
x_train = reg_train.drop(['user', 'movie', 'rating'], axis=1)
y_train = reg_train['rating']

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

# declare the model
xgb_knn_bsl = xgb.XGBRegressor(n_jobs=10, random_state=15,n_estimators=
500,learning_rate=0.1,max_depth=5)
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 : 0:01:38.495779

Done

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

TEST DATA
------------------------------
RMSE :  1.0801399036424078
MAPE :  34.243392274024146
```

### 4.4.6 Matrix Factorization Techniques

**4.4.6.1 SVD Matrix Factorization User Movie intractions**

In [38]:
```python
from surprise import SVD
```

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

## - Predicted Rating :

- $ \large  \hat r_{ui} = \mu + b_u + b_i + q_i^Tp_u $

    - $\pmb q_i$ - Representation of item(movie) in latent facto
r space

    - $\pmb p_u$ - Representation of user in new latent factor s
pace


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


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

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

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

```
In [39]: # 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, tests
         et, 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 : 0:00:08.112851

Evaluating the model with train data..
time taken : 0:00:01.540855
---------------
Train Data
---------------
RMSE : 0.6574721240954099

MAPE : 19.704901088660474

adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.080680
---------------
Test Data
---------------
RMSE : 1.0726046873826458

MAPE : 35.01953535988152

storing the test results in test dictionary...
```

```
-----------------------------------------------
Total time taken to run this algorithm : 0:00:09.736683
```

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

In [40]:
```python
from surprise import SVDpp
```

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

# - Predicted Rating :

- $ \large \hat{r}_{ui} = \mu + b_u + b_i + q_i^T\left(p_u + |I_u|^{-\frac{1}{2}} \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)

- $ \large \sum_{r_{ui} \in R_{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 [41]:
```python
# initiallize the model
```

```python
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 : 0:02:10.605341

Evaluating the model with train data..
time taken : 0:00:07.993197
---------------
Train Data
---------------
RMSE : 0.6032438403305899

MAPE : 17.49285063490268
```

```
adding train results in the dictionary..

Evaluating for test data...
time taken : 0:00:00.081978
---------------
Test Data
---------------
RMSE : 1.0728491944183447

MAPE : 35.03817913919887

storing the test results in test dictionary...

-----------------------------------------------
Total time taken to run this algorithm : 0:02:18.682182
```

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

**Preparing Train data**

In [42]:
```python
# 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[42]:

| | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr2 | ... | smr4 | smr5 | UAvg |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 53406 | 33 | 3.581679 | 4.0 | 5.0 | 5.0 | 4.0 | 1.0 | 5.0 | 2.0 | ... | 3.0 | 1.0 | 3.370370 |
| **1** | 99540 | 33 | 3.581679 | 5.0 | 5.0 | 5.0 | 4.0 | 5.0 | 3.0 | 4.0 | ... | 3.0 | 5.0 | 3.555556 |

2 rows × 21 columns

**Preparing Test data**

```
In [43]:   reg_test_df['svd'] = models_evaluation_test['svd']['predictions']
           reg_test_df['svdpp'] = models_evaluation_test['svdpp']['predictions']

           reg_test_df.head(2)
```

Out[43]:

|   | user | movie | GAvg | sur1 | sur2 | sur3 | sur4 | sur5 | smr1 | smr |
|---|------|-------|------|------|------|------|------|------|------|-----|
| **0** | 808635 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |
| **1** | 941866 | 71 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.581679 | 3.58167 |

**2 rows × 21 columns**

```
In [44]:   # 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']


           xgb_final = xgb.XGBRegressor(n_jobs=10, random_state=15,n_estimators=50
           0,learning_rate=0.1,max_depth=5)
           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 : 0:01:39.095277

Done

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

TEST DATA
------------------------------
RMSE :  1.1087132923749305
MAPE :  33.15708167163116
```

### 4.4.8 XgBoost with Surprise Baseline + Surprise KNNbaseline + MF Techniques

In [45]:
```python
# 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']


xgb_all_models = xgb.XGBRegressor(n_jobs=10, random_state=15,n_estimato
rs=500,learning_rate=0.1,max_depth=5)
train_results, test_results = run_xgboost(xgb_all_models, x_train, y_tr
ain, 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 : 0:01:06.078385

Done

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

TEST DATA
------------------------------
RMSE :  1.075851597942998
MAPE :  34.94487253615702
```

## 4.5 Comparision between all models

```
In [47]:    # Saving our TEST_RESULTS into a dataframe so that you don't have to ru
            n it again
            pd.DataFrame(models_evaluation_test).to_csv('small_sample_results1.csv'
            )
            models = pd.read_csv('small_sample_results1.csv', index_col=0)
            models.loc['rmse'].sort_values()
```

```
Out[47]:    svd                 1.0726046873826458
            knn_bsl_u           1.0726493739667242
            knn_bsl_m            1.072758832653683
            svdpp               1.0728491944183447
            bsl_algo            1.0730330260516174
            xgb_all_models       1.075851597942998
            xgb_knn_bsl         1.0801399036424078
            xgb_bsl             1.0831711919892506
            first_algo          1.1017397109523788
```

```
xgb_final       1.1087132923749305
Name: rmse, dtype: object
```