



1. Business Problem

1.1 Problem Description

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

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

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

1.2 Problem Statement

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

1.3 Sources

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

1.4 Real world/Business Objectives and constraints

Objectives:

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

Constraints:

1. Some form of interpretability.

2. Machine Learning Problem

2.1 Data

2.1.1 Data Overview

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

Data files :

- combined_data_1.txt
- combined_data_2.txt
- combined_data_3.txt

- combined_data_4.txt
- movie_titles.csv

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

CustomerID,Rating,Date

MovieIDs range from 1 to 17770 sequentially.

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

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

Dates have the format YYYY-MM-DD.

2.1.2 Example Data point

1:

```
1488844,3,2005-09-06
822109,5,2005-05-13
885013,4,2005-10-19
30878,4,2005-12-26
823519,3,2004-05-03
893988,3,2005-11-17
124105,4,2004-08-05
1248029,3,2004-04-22
1842128,4,2004-05-09
2238063,3,2005-05-11
1503895,4,2005-05-19
2207774,5,2005-06-06
2590061,3,2004-08-12
2442,3,2004-04-14
543865,4,2004-05-28
1209119,4,2004-03-23
804919,4,2004-06-10
1086807,3,2004-12-28
1711859,4,2005-05-08
372233,5,2005-11-23
1080361,3,2005-03-28
1245640,3,2005-12-19
```

558634,4,2004-12-14
2165002,4,2004-04-06
1181550,3,2004-02-01
1227322,4,2004-02-06
427928,4,2004-02-26
814701,5,2005-09-29
808731,4,2005-10-31
662870,5,2005-08-24
337541,5,2005-03-23
786312,3,2004-11-16
1133214,4,2004-03-07
1537427,4,2004-03-29
1209954,5,2005-05-09
2381599,3,2005-09-12
525356,2,2004-07-11
1910569,4,2004-04-12
2263586,4,2004-08-20
2421815,2,2004-02-26
1009622,1,2005-01-19
1481961,2,2005-05-24
401047,4,2005-06-03
2179073,3,2004-08-29
1434636,3,2004-05-01
93986,5,2005-10-06
1308744,5,2005-10-29
2647871,4,2005-12-30
1905581,5,2005-08-16
2508819,3,2004-05-18
1578279,1,2005-05-19
1159695,4,2005-02-15
2588432,3,2005-03-31
2423091,3,2005-09-12
470232,4,2004-04-08
2148699,2,2004-06-05
1342007,3,2004-07-16
466135,4,2004-07-13
2472440,3,2005-08-13
1283744,3,2004-04-17
1927580,4,2004-11-08
716874,5,2005-05-06
4326,4,2005-10-29

2.2 Mapping the real world problem to a Machine

Learning Problem

2.2.1 Type of Machine Learning Problem

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

The given problem is a Recommendation problem

It can also be seen as a Regression problem

2.2.2 Performance metric

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

2.2.3 Machine Learning Objective and Constraints

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

In [1]:

```
# this is just to know how much time will it take to run this entire ipython
from datetime import datetime
# globalstart = datetime.now()
import pandas as pd
import numpy as np
import matplotlib
matplotlib.use('nbagg')

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

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

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

import plotly
import plotly.offline as offline
import plotly.graph_objs as go
offline.init_notebook_mode()
```

3. Exploratory Data Analysis

3.1 Preprocessing

3.1.1 Converting / Merging whole data to required format: u_i , m_j , r_{ij}

In [0]:

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

    row = list()
    files=['data_folder/combined_data_1.txt', 'data_folder/combined_data_2.txt',
           'data_folder/combined_data_3.txt', 'data_folder/combined_data_4.txt']
    for file in files:
        print("Reading ratings from {}".format(file))
        with open(file) as f:
            for line in f:
                del row[:] # you don't have to do this.
                line = line.strip()
                if line.endswith(':'):
                    # All below are ratings for this movie, until another movie
                    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[14]:

	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[7]:

```
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 timestamp)
print("There are {} duplicate rating entries in the data..".format(dups))
```

There are 0 duplicate rating entries in the data..

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

In [0]:

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

Total data

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

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

In [0]:

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

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

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

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

In [0]:

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

Training data

```
Total no of ratings : 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"
```

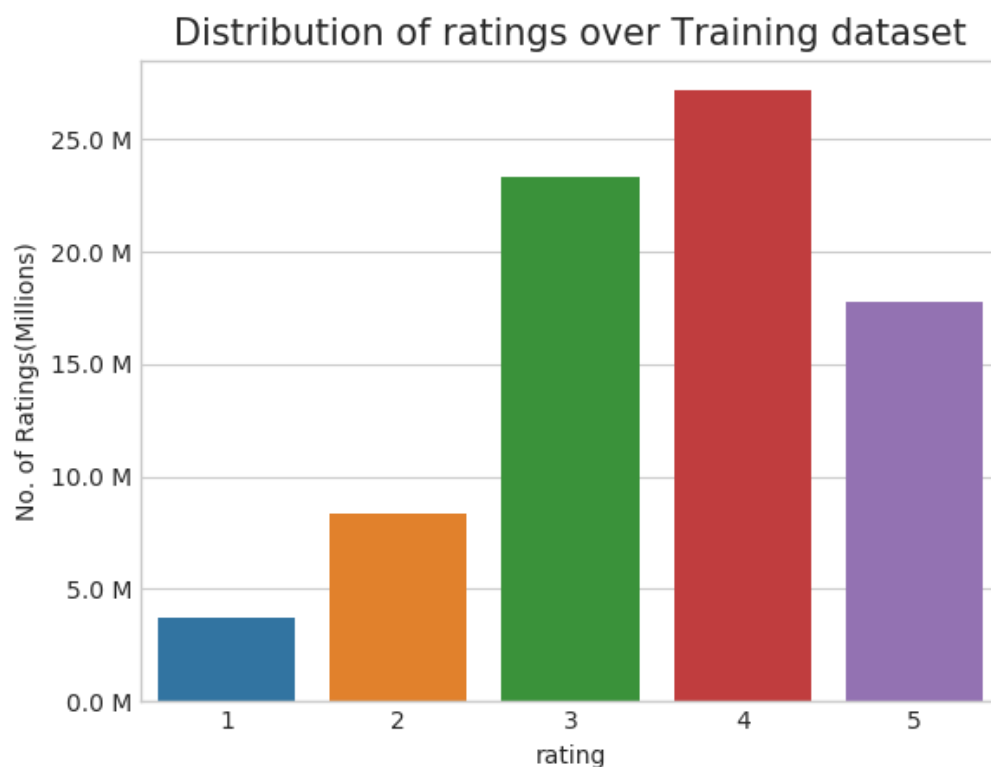
3.3.1 Distribution of ratings

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()
```

<IPython.core.display.Javascript object>



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

In [0]:

```
# It is used to skip the warning 'SettingWithCopyWarning'..  
pd.options.mode.chained_assignment = None # default='warn'  
  
train_df['day_of_week'] = train_df.date.dt.weekday_name  
  
train_df.tail()
```

Out[17]:

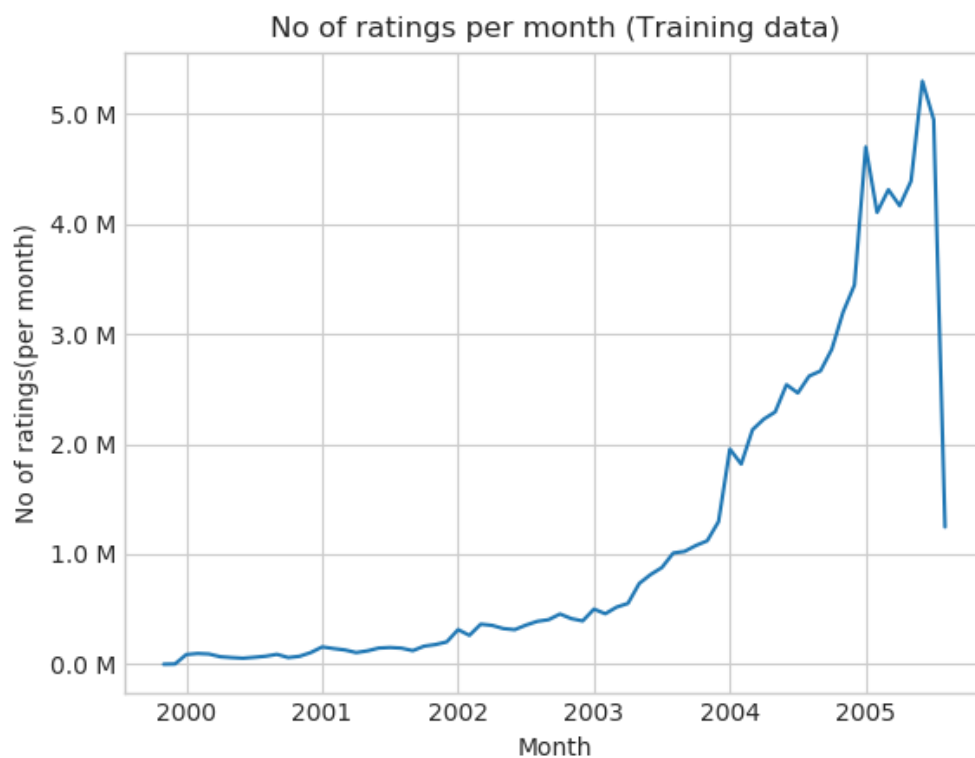
	movie	user	rating	date	day_of_week
80384400	12074	2033618	4	2005-08-08	Monday
80384401	862	1797061	3	2005-08-08	Monday
80384402	10986	1498715	5	2005-08-08	Monday
80384403	14861	500016	4	2005-08-08	Monday
80384404	5926	1044015	5	2005-08-08	Monday

3.3.2 Number of Ratings per a month

In [0]:

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

<IPython.core.display.Javascript object>



3.3.3 Analysis on the Ratings given by user

In [0]:

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

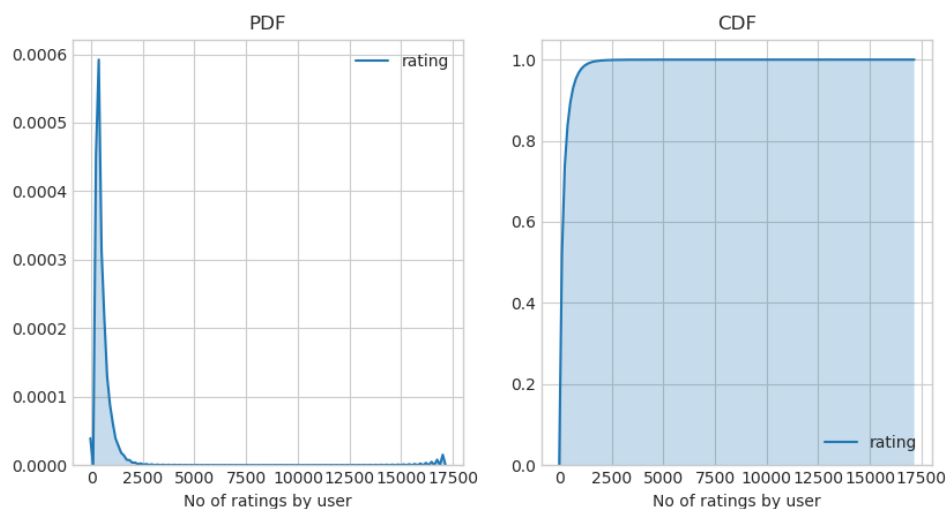
Out[20]:

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

In [0]:

```
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()
```

<IPython.core.display.Javascript object>



In [0]:

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

Out[22]:

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

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

In [0]:

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

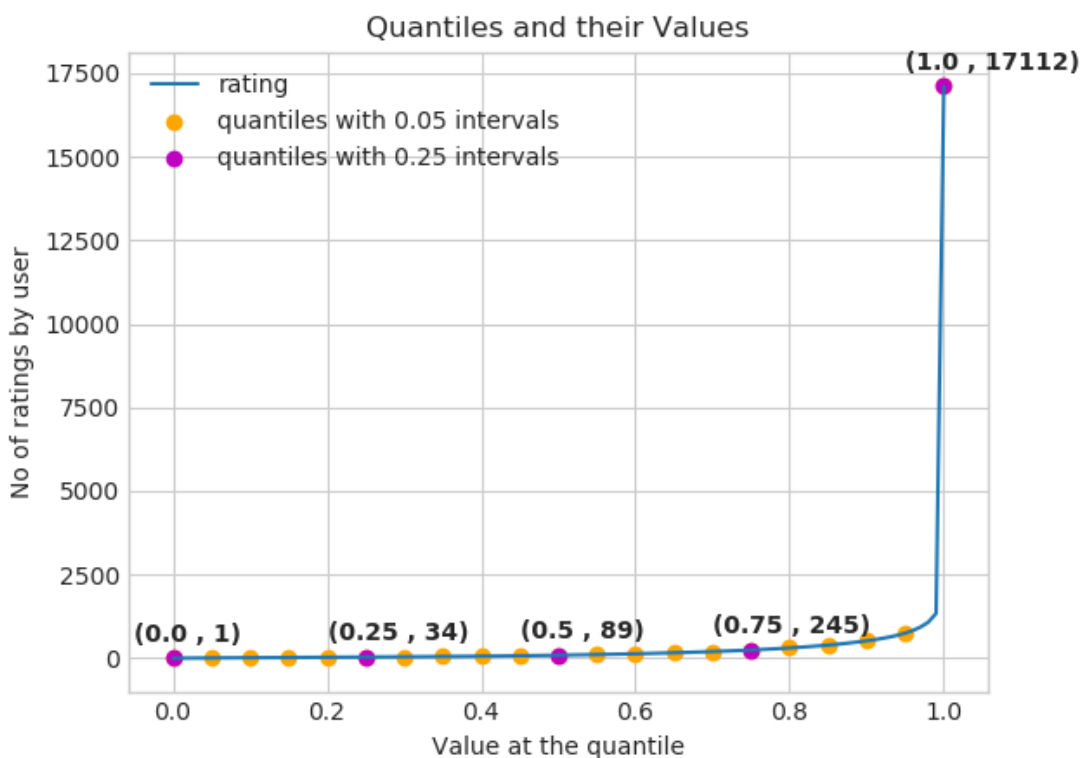

In [0]:

```
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 = '0.05 intervals')
# quantiles with 0.25 difference
plt.scatter(x=quantiles.index[::25], y=quantiles.values[::25], c='m', label = '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()
```

<IPython.core.display.Javascript object>



In [0]:

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

Out[25]:

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]:

```
print('\n No of ratings at last 5 percentile : {}'.format(sum(no_of Rated_r
```

No of ratings at last 5 percentile : 20305

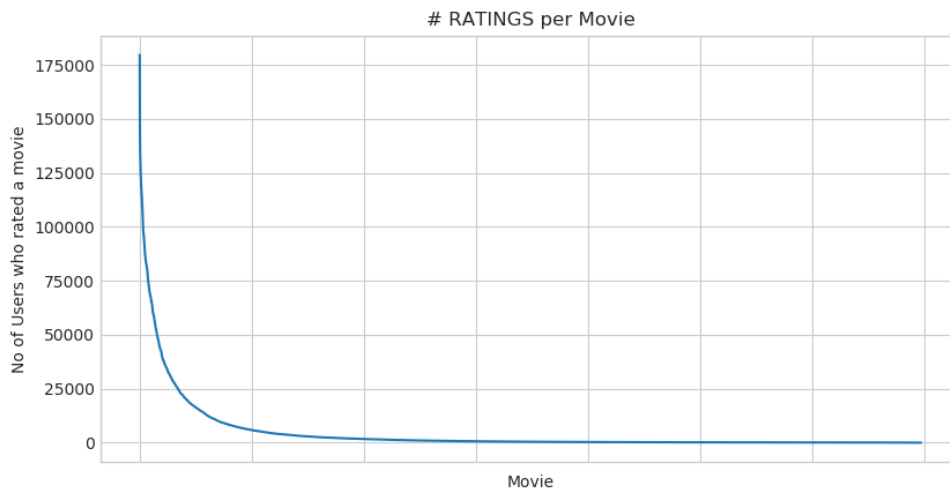
3.3.4 Analysis of ratings of a movie given by a user

In [0]:

```
no_of_ratings_per_movie = train_df.groupby(by='movie')['rating'].count().sort
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()
```

<IPython.core.display.Javascript object>



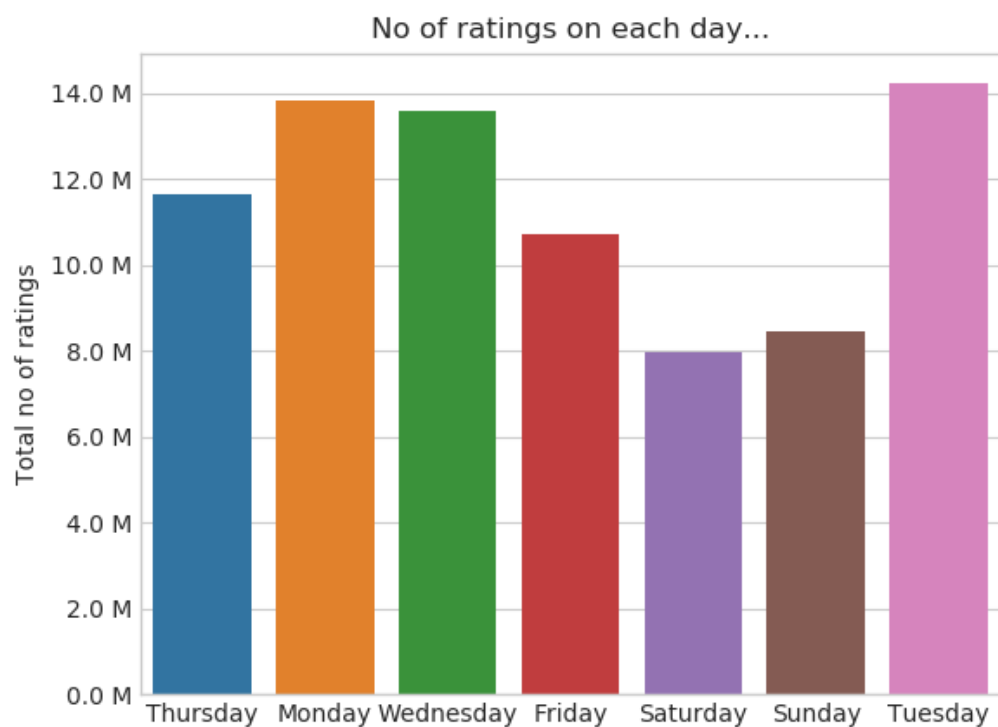
- **It is very skewed.. just like nunmber of ratings given per user.**
 - There are some movies (which are very popular) which are rated by huge number of users.
 - But most of the movies(like 90%) got some hundereds of ratings.

3.3.5 Number of ratings on each day of the week

In [0]:

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

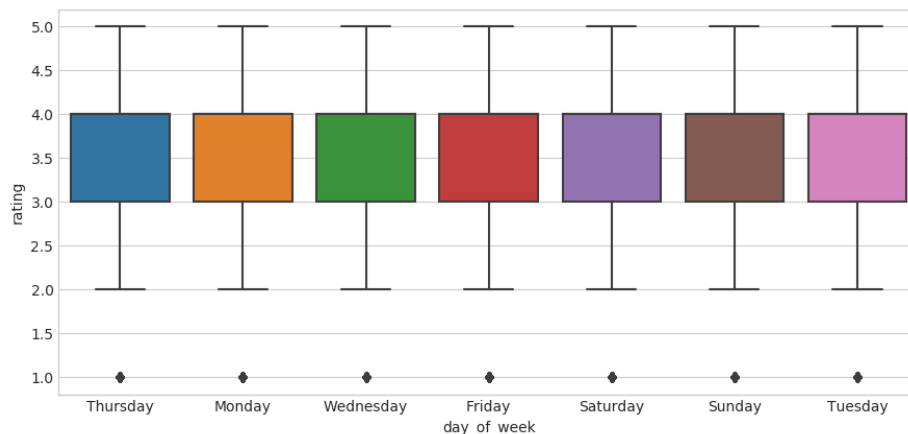
<IPython.core.display.Javascript object>



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)
```

<IPython.core.display.Javascript object>



0:01:10.003761

In [0]:

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

Average ratings

```
-----
day_of_week
Friday      3.585274
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 [2]:

```
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_id.values,
                                                                    train_df.movie.values)),)

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

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

```
It is present in your pwd, getting it from disk....
DONE..
0:00:03.009200
```

The Sparsity of Train Sparse Matrix

In [3]:

```
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 [4]:

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

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

DONE..

0:00:00.779965

The Sparsity of Test data Matrix

In [5]:

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

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

Sparsity Of Test matrix : 99.95731772988694 %

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

In [2]:

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

def get_average_ratings(sparse_matrix, of_users):

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

    # ".A1" is for converting Column_Matrix to 1-D numpy array
    sum_of_ratings = sparse_matrix.sum(axis=ax).A1
    # Boolean matrix of ratings ( whether a user rated that movie or not)
    is_rated = sparse_matrix!=0
    # no of ratings that each user OR movie..
    no_of_ratings = is_rated.sum(axis=ax).A1

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

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

3.3.7.1 finding global average of all movie ratings

In [0]:

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

Out[36]:

```
{'global': 3.582890686321557}
```

3.3.7.2 finding average rating per user

In [0]:

```
train_averages['user'] = get_average_ratings(train_sparse_matrix, of_users=Tr  
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=Tr  
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]:

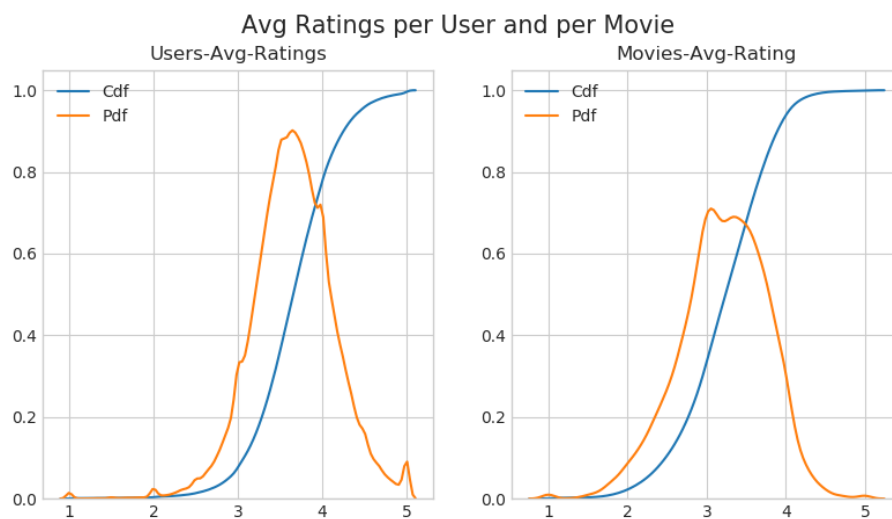
```
start = datetime.now()
# draw pdfs for average rating per user and average
fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))
fig.suptitle('Avg Ratings per User and per Movie', fontsize=15)

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)
```

<IPython.core.display.Javascript object>



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 ".format(new_users, new_users/total_users*100))
```

Total number of Users : 480189

Number of Users in Train data : 405041

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

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

3.3.8.2 Cold Start problem with Movies

In [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 ".format(new_movies, new_movies/total_movies*100))
```

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

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

    start = datetime.now()
    temp = 0

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

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

        # add them to our rows, cols and data
        rows.extend([row]*top)
        cols.extend(top_sim_ind)
        data.extend(top_sim_val)
        time_taken.append(datetime.now().timestamp() - prev.timestamp())
        if verbose:
            if temp%verb_for_n_rows == 0:
                print("computing done for {} users [ time elapsed : {} ]".format(temp, datetime.now()-start))

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

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

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

In [0]:

```
start = datetime.now()
u_u_sim_sparse, _ = compute_user_similarity(train_sparse_matrix, compute_for_
                                             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

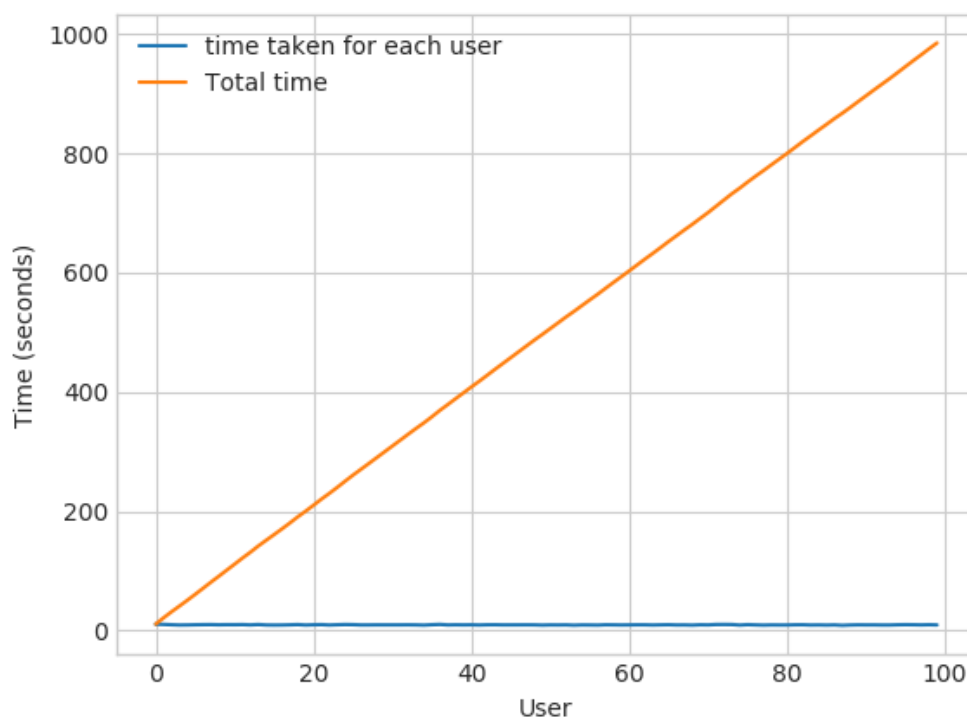
]

computing done for 100 users [time elapsed : 0:16:24.711032

]

Creating Sparse matrix from the computed similarities

<IPython.core.display.Javascript object>



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 our training set and computing similarities between them..(**17K dimensional vector**..) is time consuming..
- From above plot, It took roughly **8.88 sec** for computing similar users for **one user**
- We have **405,041 users** with us in training set.
- $405041 \times 8.88 = 3596764.08 \text{ sec} = 59946.068 \text{ min} = 999.101133333 \text{ hours} = 41.6$
 - Even if we run on 4 cores parallelly (a typical system now a days), It will still take almost **10 and 1/2 days**.

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

In [0]:

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

start = datetime.now()

# initialize the algorithm with some parameters..
# All of them are default except n_components. n_iter is for Randomized SVD so
netflix_svd = TruncatedSVD(n_components=500, algorithm='randomized', random_s
trunc_svd = netflix_svd.fit_transform(train_sparse_matrix)

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

0:29:07.069783

Here,

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

In [0]:

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

In [0]:

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

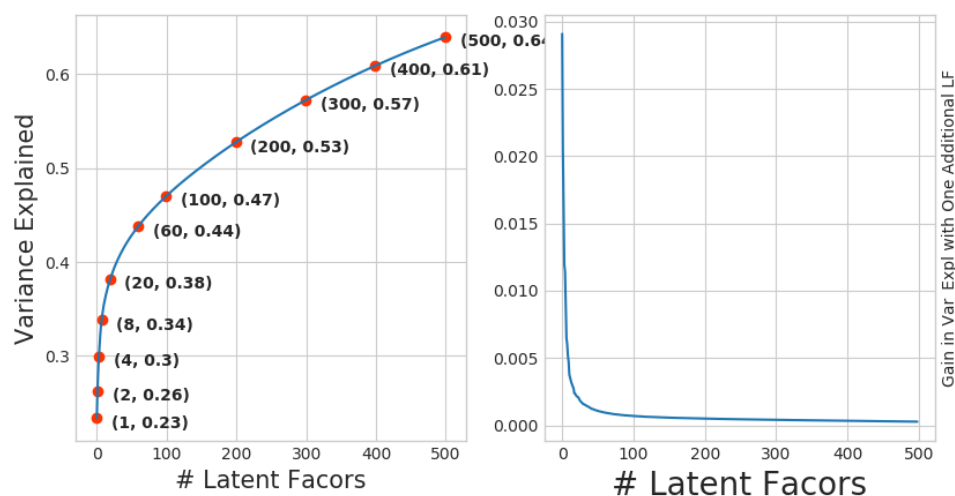
ax1.set_ylabel("Variance Explained", fontsize=15)
ax1.set_xlabel("# Latent Facors", fontsize=15)
ax1.plot(expl_var)
# annotate some (latentfactors, expl_var) to make it clear
ind = [1, 2, 4, 8, 20, 60, 100, 200, 300, 400, 500]
ax1.scatter(x = [i-1 for i in ind], y = expl_var[[i-1 for i in ind]], c='#ff3333')
for i in ind:
    ax1.annotate(s = "({}, {})".format(i, np.round(expl_var[i-1], 2)), xy=(i-1, expl_var[i-1]),
                xytext = ( i+20, expl_var[i-1] - 0.01), fontweight='bold')

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

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

plt.show()
```

<IPython.core.display.Javascript object>



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]:

```
# 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]:

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

Out[53]:

(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 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[55]:

(2649430, 500)

In [0]:

```
start = datetime.now()
trunc_u_u_sim_matrix, _ = compute_user_similarity(trunc_sparse_matrix, compute_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

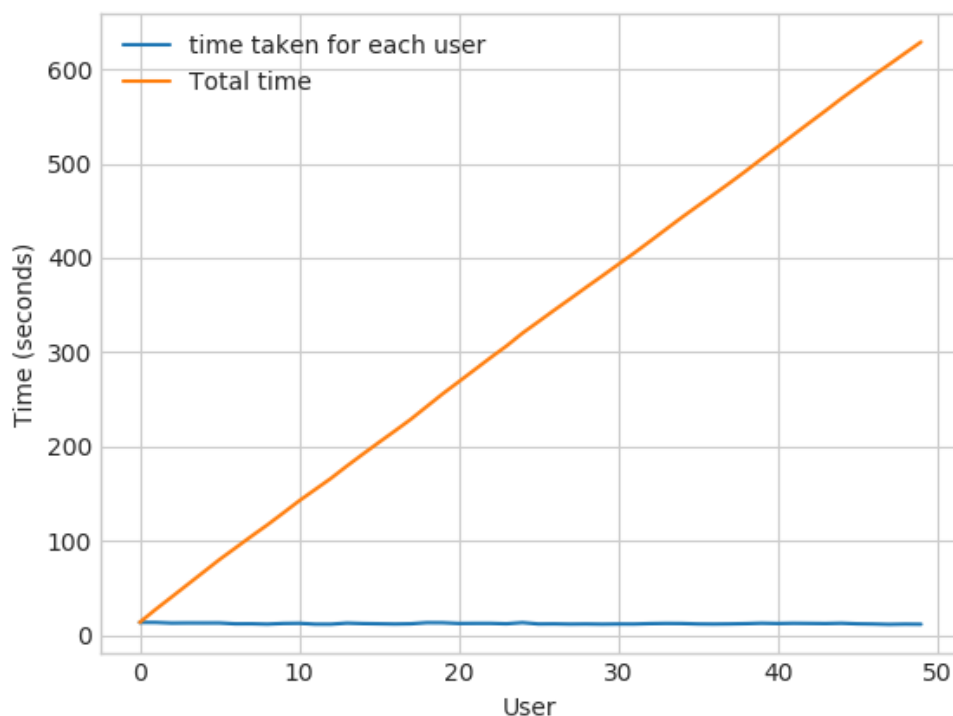
]

computing done for 50 users [time elapsed : 0:10:28.861485

]

Creating Sparse matrix from the computed similarities

<IPython.core.display.Javascript object>



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 similar users for **one user**
- We have **405041 users** with us in training set.
- $405041 \times 12.18 == 4933399.38 \text{ sec} == 82223.323 \text{ min} == 1370.38871 \text{ days}$
 - Even we run on 4 cores parallelly (a typical system now a days), It will still take almost (14 - 15) days.

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

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

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

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

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

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

3.4.2 Computing Movie-Movie Similarity matrix

In [0]:

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

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

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

It seems you don't have that file. Computing movie_movie similarity...

Done..

Saving it to disk without the need of re-computing it again..

Done..

It's a (17771, 17771) dimensional matrix

0:10:02.736054

In [0]:

```
m_m_sim_sparse.shape
```

Out[59]:

(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]:

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

In [0]:

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

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

0:00:33.411700

Out[62]:

```
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....

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=',', header =  
                           names=['movie_id', 'year_of_release', 'title'], ve  
                           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[64]:

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

Similar Movies for 'Vampire Journals'

In [0]:

```
mv_id = 67

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

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

print("\nWe have {} movies which are similar to this and we will get only top 10")
```

Movie -----> Vampire Journals

It has 270 Ratings from users.

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

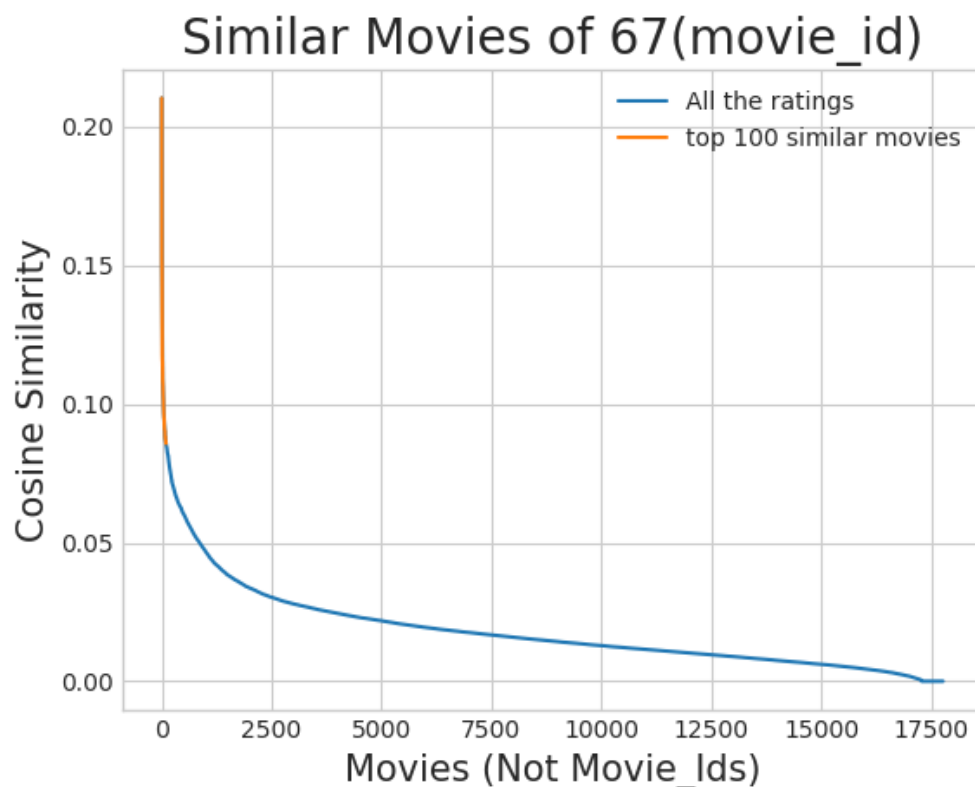
In [0]:

[illegible]

In [0]:

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

<IPython.core.display.Javascript object>



Top 10 similar movies

In [0]:

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

Out[68]:

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

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

4. Machine Learning Models



In [3]:

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

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

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

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

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

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

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

    return sample_sparse_matrix
```

4.1 Sampling Data

4.1.1 Build sample train data from the train data

In [8]:

```
start = datetime.now()
path = "sample_train_sparse_matrix.npz"
if os.path.isfile(path):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    sample_train_sparse_matrix = sparse.load_npz(path)
    print("DONE..")
else:
    # get 15k users and 1.5k movies from available data
    sample_train_sparse_matrix = get_sample_sparse_matrix(train_sparse_matrix,
                                                            path = path)

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

Original Matrix : (users, movies) -- (405041 17424)

Original Matrix : Ratings -- 80384405

Sampled Matrix : (users, movies) -- (15000 1500)

Sampled Matrix : Ratings -- 276144

Saving it into disk for furthur usage..

Done..

0:00:28.305985

4.1.2 Build sample test data from the test data

In [9]:

```
start = datetime.now()

path = "sample_test_sparse_matrix.npz"
if os.path.isfile(path):
    print("It is present in your pwd, getting it from disk....")
    # just get it from the disk instead of computing it
    sample_test_sparse_matrix = sparse.load_npz(path)
    print("DONE..")
else:
    # get 7.5k users and 750 movies from available data
    sample_test_sparse_matrix = get_sample_sparse_matrix(test_sparse_matrix,
                                                         path = "sample_test_sparse_r
print(datetime.now() - start)
```

Original Matrix : (users, movies) -- (349312 17757)

Original Matrix : Ratings -- 20096102

Sampled Matrix : (users, movies) -- (7500 750)

Sampled Matrix : Ratings -- 19273

Saving it into disk for furthur usage..

Done..

0:00:07.031471

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

In [10]:

```
sample_train_averages = dict()
```

4.2.1 Finding Global Average of all movie ratings

In [11]:

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

Out[11]:

```
{'global': 3.5405513065646907}
```

4.2.2 Finding Average rating per User

In [12]:

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

Average rating of user 1515220 : 3.953488372093023

4.2.3 Finding Average rating per Movie

In [13]:

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

AVerage rating of movie 15153 : 2.6184210526315788

4.3 Featurizing data

In [14]:

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

No of ratings in Our Sampled train matrix is : 276144

No of ratings in Our Sampled test matrix is : 19273

4.3.1 Featurizing data for regression problem

4.3.1.1 Featurizing train data

In [15]:

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


In [18]:

```
#####
# It took me almost 10 hours to prepare this train dataset.#
#####
start = datetime.now()
if os.path.isfile('reg_train.csv'):
    print("File already exists you don't have to prepare again..." )
else:
    print('preparing {} tuples for the dataset..\n'.format(len(sample_train_r
with open('reg_train.csv', mode='w') as reg_data_file:
    count = 0
    for (user, movie, rating) in zip(sample_train_users, sample_train_mo
        st = datetime.now()
        #     print(user, movie)
        #----- Ratings of "movie" by similar users of "us
        # compute the similar Users of the "user"
        user_sim = cosine_similarity(sample_train_sparse_matrix[user], sa
        top_sim_users = user_sim.argsort()[::-1][1:] # we are ignoring '
        # get the ratings of most similar users for this movie
        top_ratings = sample_train_sparse_matrix[top_sim_users, movie].to
        # 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'][movi
        #     print(top_sim_users_ratings, end=" ")

        #----- Ratings by "user" to similar movies of "n
        # compute the similar movies of the "movie"
        movie_sim = cosine_similarity(sample_train_sparse_matrix[:,movie]
        top_sim_movies = movie_sim.argsort()[::-1][1:] # we are ignoring
        # get the ratings of most similar movie rated by this user..
        top_ratings = sample_train_sparse_matrix[user, top_sim_movies].to
        # 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
        #     print(top_sim_movies_ratings, end=" : -- ")

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

```

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

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

print(datetime.now() - start)

```

preparing 276144 tuples for the dataset..

```

Done for 1000 rows----- 0:04:43.563130
Done for 2000 rows----- 0:09:30.963902
Done for 3000 rows----- 0:14:14.426430
Done for 4000 rows----- 0:18:59.703637
Done for 5000 rows----- 0:23:42.246337
Done for 6000 rows----- 0:28:25.823029
Done for 7000 rows----- 0:33:09.586466
Done for 8000 rows----- 0:38:00.338071
Done for 9000 rows----- 0:42:51.666675
Done for 10000 rows----- 0:47:42.910121
Done for 11000 rows----- 0:52:31.074167
Done for 12000 rows----- 0:57:16.607306
Done for 13000 rows----- 1:02:05.925336
Done for 14000 rows----- 1:06:55.038330
Done for 15000 rows----- 1:11:44.651786
Done for 16000 rows----- 1:16:30.790844
Done for 17000 rows----- 1:21:12.281023
Done for 18000 rows----- 1:25:54.056424

```

Reading from the file to make a Train_dataframe

In [2]:

```
reg_train = pd.read_csv('reg_train.csv', names = ['user', 'movie', 'GAvg', 's
reg_train.head()
```

Out[2]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3
0	174683	10	3.540551	4.0	5.0	4.0	3.0	4.0	3.0	3.0	3.0
1	233949	10	3.540551	4.0	4.0	5.0	5.0	5.0	2.0	3.0	4.0
2	767518	10	3.540551	4.0	5.0	4.0	3.0	4.0	5.0	3.0	3.0
3	894393	10	3.540551	5.0	5.0	5.0	4.0	3.0	4.0	4.0	4.0
4	951907	10	3.540551	4.0	5.0	3.0	4.0	4.0	4.0	5.0	4.0

-
- **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 [20]:

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

In [21]:

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

Out[21]:

3.5405513065646907

In [36]:

```
start = datetime.now()

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

    print('preparing {} tuples for the dataset..\n'.format(len(sample_test_ratings)))
    with open('reg_test.csv', mode='w') as reg_data_file:
        count = 0
        for (user, movie, rating) in zip(sample_test_users, sample_test_movies, sample_test_ratings):
            st = datetime.now()

            #----- Ratings of "movie" by similar users of "user"
            #print(user, movie)
            try:
                # compute the similar Users of the "user"
                user_sim = cosine_similarity(sample_train_sparse_matrix[user, :].T, sample_train_sparse_matrix)
                top_sim_users = user_sim.argsort()[::-1][1:] # we are ignoring the user itself
                # get the ratings of most similar users for this movie
                top_ratings = sample_train_sparse_matrix[top_sim_users, movie]
                # we will make it's length "5" by adding movie averages to it
                top_sim_users_ratings = list(top_ratings[top_ratings != 0][:5])
                top_sim_users_ratings.extend([sample_train_averages['movie', movie]] * (5 - len(top_sim_users_ratings)))
                # print(top_sim_users_ratings, end="--")

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

            #----- Ratings by "user" to similar movies of "movie"
            try:
                # compute the similar movies of the "movie"
                movie_sim = cosine_similarity(sample_train_sparse_matrix[:, movie], sample_train_sparse_matrix)
                top_sim_movies = movie_sim.argsort()[::-1][1:] # we are ignoring the movie itself
                # get the ratings of most similar movie rated by this user..
                top_ratings = sample_train_sparse_matrix[user, top_sim_movies]
                # we will make it's length "5" by adding user averages to it
                top_sim_movies_ratings = list(top_ratings[top_ratings != 0][:5])
                top_sim_movies_ratings.extend([sample_train_averages['user', user]] * (5 - len(top_sim_movies_ratings)))
                #print(top_sim_movies_ratings)
            except (IndexError, KeyError):
                #print(top_sim_movies_ratings, end=" : -- ")
```

```

        top_sim_movies_ratings.extend([sample_train_averages['global']
        #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()))
print("",datetime.now() - start)

```

preparing 19273 tuples for the dataset..

Done for 1000 rows----- 0:04:36.552196

Done for 2000 rows----- 0:09:13.377324

```

Done for 3000 rows----- 0:13:53.273002
Done for 4000 rows----- 0:18:30.980331
Done for 5000 rows----- 0:23:08.347486
Done for 6000 rows----- 0:27:49.417384
Done for 7000 rows----- 0:32:37.135531
Done for 8000 rows----- 0:37:23.502513
Done for 9000 rows----- 0:42:01.777354
Done for 10000 rows----- 0:46:38.994316
Done for 11000 rows----- 0:51:16.032628
Done for 12000 rows----- 0:55:53.470153
Done for 13000 rows----- 1:00:31.110496
Done for 14000 rows----- 1:05:08.322624
Done for 15000 rows----- 1:09:45.398019
Done for 16000 rows----- 1:14:22.515516
Done for 17000 rows----- 1:18:59.469327
Done for 18000 rows----- 1:23:39.638616
Done for 19000 rows----- 1:28:24.485672
1:29:41.994036

```

__Reading from the file to make a test dataframe __

In [3]:

```

reg_test_df = pd.read_csv('reg_test.csv', names = ['user', 'movie', 'GAvg',
                                                    'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
                                                    'UAvg', 'MAvg', 'RAvg'])
reg_test_df.head(4)

```

Out[3]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5
0	808635	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551
1	898730	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551
2	941866	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551
3	1280761	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551

- **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 [4]:

```
from surprise import Reader, Dataset
```

4.3.2.1 Transforming train data

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

In [5]:

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

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

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

4.3.2.2 Transforming test data

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

In [6]:

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

Out[6]:

```
[(808635, 71, 5), (898730, 71, 3), (941866, 71, 4)]
```

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 [7]:

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

models_evaluation_train, models_evaluation_test
```

Out[7]:

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

Utility functions for running regression models

In [8]:

```
# 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_
    mape = np.mean(np.abs( (y_true - y_pred)/y_true )) * 100
    return rmse, mape

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

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

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

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

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

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

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

```
    print('MAPE : ', mape_test)

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

Utility functions for Surprise modes

In [9]:

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

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

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

    return actual, pred

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

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

    return rmse, mape*100

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

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

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

    # ----- Evaluating train data-----#
```

```

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

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

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

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

if verbose:
    print('-'*15)
    print('Test Data')
    print('-'*15)
    print("RMSE : {}\nMAPE : {}".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 [10]:

```
import xgboost as xgb
```

In [11]:

```
# 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']
```

In [12]:

```
from sklearn.metrics import mean_squared_error
from sklearn.metrics import make_scorer

def rmse(y_true, y_pred):
    rmse = mean_squared_error(y_true, y_pred, squared=False)
    return rmse

scorer = make_scorer(rmse)
```

Hyper parmameter tuning xgboost

In [13]:

```
from sklearn.model_selection import GridSearchCV

parameters = {
    'max_depth' : [1, 2, 3, 4, 5, 6, 7],
    'n_estimators' : [10, 50, 100, 200, 300, 500]
}
xgbdt = xgb.XGBRegressor(n_jobs=18)
clf = GridSearchCV(xgbdt, parameters, cv=2, scoring = scorer, return_train_score=True)
clf.fit(x_train, y_train)
```

Out[13]:

```
GridSearchCV(cv=2, error_score=nan,
             estimator=XGBRegressor(base_score=None, booster=None,
                                     colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None, gamma=None,
                                     gpu_id=None, importance_type=None,
                                     interaction_constraints=None,
                                     learning_rate=None, max_depth=None,
                                     max_features=None, min_child_weight=None,
                                     missing=nan, monotone_constraints=None,
                                     n_estimators=None, objective='reg:squarederror',
                                     random_state=None, reg_alpha=None,
                                     reg_lambda=None, scale_pos_weight=None,
                                     subsample=None, tree_method=None,
                                     validate_parameters=False,
                                     verbosity=None),
             iid='deprecated', n_jobs=-1,
             param_grid={'max_depth': [1, 2, 3, 4, 5, 6, 7],
                          'n_estimators': [10, 50, 100, 200, 300, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring=make_scorer(rmse), verbose=0)
```

In [14]:

```
results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['mean_test_score'])
results.head(10)
```

Out[14]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max
11	54.807801	0.277757	0.222904	0.015459	
10	37.050697	0.033909	0.174547	0.002008	
16	45.279694	0.046874	0.335104	0.004987	
9	26.617979	0.105219	0.177525	0.005984	
17	67.563874	0.060585	0.251328	0.012965	
21	38.544831	0.174288	0.189992	0.005485	
22	54.360793	0.233370	0.300205	0.052851	
15	33.012030	0.006981	0.197971	0.010472	
20	21.417908	0.121675	0.156082	0.006482	
14	18.306701	0.054854	0.178024	0.013464	

In [15]:

```
# initialize Our first XGBoost model...
first_xgb = xgb.XGBRegressor(silent=False, random_state=15, max_depth = 2, n_
train_results, test_results = run_xgboost(first_xgb, x_train, y_train, x_test

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

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

Training the model..

Done. Time taken : 0:00:12.585355

Done

Evaluating the model with TRAIN data...

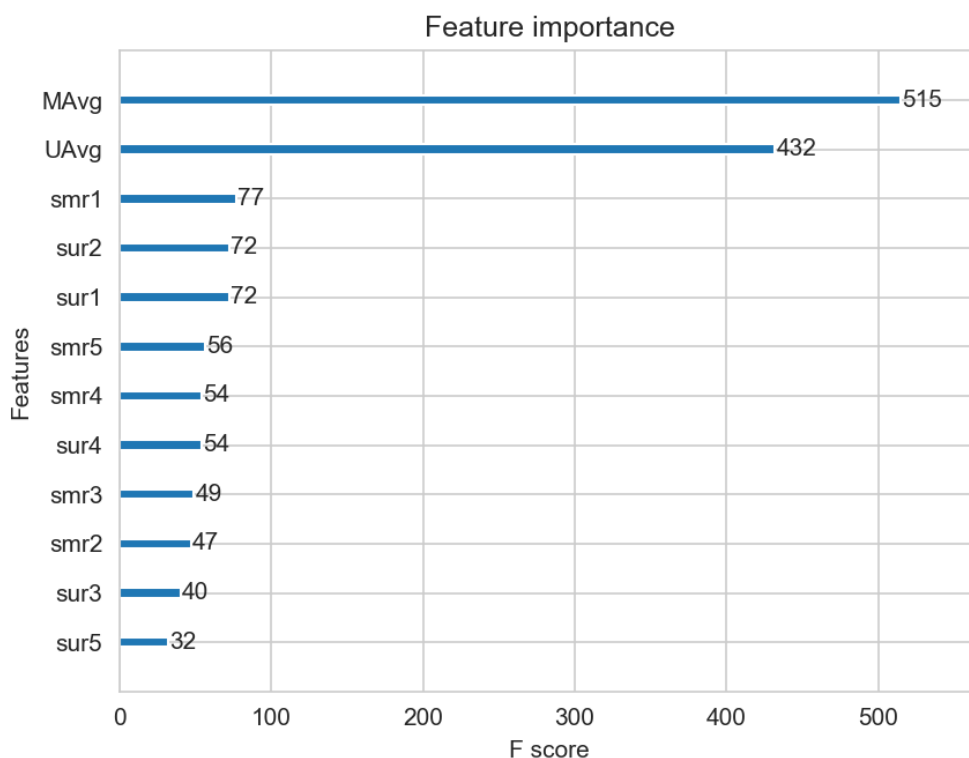
Evaluating Test data

TEST DATA

RMSE : 1.0795164032626583

MAPE : 34.89138465404024

<IPython.core.display.Javascript object>



4.4.2 Surprise BaselineModel

In [16]:

```
from surprise import BaselineOnly
```

__Predicted_rating : (baseline prediction) __

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

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

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

__Optimization function (Least Squares Problem) __

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

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

In [17]:

```
# options are to specify.., how to compute those user and item biases
bsl_options = {'method': 'sgd',
               'learning_rate': 0.1
              }
bsl_algo = BaselineOnly(bsl_options=bsl_options)
# run this algorithm.., It will return the train and test results..
bsl_train_results, bsl_test_results = run_surprise(bsl_algo, trainset, testset)

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

Training the model...

Estimating biases using sgd...

Done. time taken : 0:00:01.281574

Evaluating the model with train data..

time taken : 0:00:01.237689

Train Data

RMSE : 0.9190730925196674

MAPE : 27.973152796690297

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.212432

Test Data

RMSE : 1.0818422475962859

MAPE : 35.22518865420994

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:00:02.732692

4.4.3 XGBoost with initial 13 features + Surprise Baseline predictor

Updating Train Data

In [18]:

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

Out[18]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3
0	174683	10	3.540551	4.0	5.0	4.0	3.0	4.0	3.0	3.0	3.0
1	233949	10	3.540551	4.0	4.0	5.0	5.0	5.0	2.0	3.0	4.0

Updating Test Data

In [19]:

```
# add that baseline predicted ratings with Surprise to the test data as well
reg_test_df['bslpr'] = models_evaluation_test['bsl_algo']['predictions']
reg_test_df.head(2)
```

Out[19]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5
0	808635	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551
1	898730	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551

In [20]:

```
# 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']
```

In [21]:

```
parameters = {
    'max_depth' : [1, 2, 3, 4, 5, 6, 7],
    'n_estimators' : [10, 50, 100, 200, 300, 500]
}
xgbdt = xgb.XGBRegressor(n_jobs=18)
clf = GridSearchCV(xgbdt, parameters, cv=2, scoring = scorer, return_train_score=True)
clf.fit(x_train, y_train)
```

Out[21]:

```
GridSearchCV(cv=2, error_score=nan,
             estimator=XGBRegressor(base_score=None, booster=None,
                                     colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None, gamma=None,
                                     max_depth=None, min_child_weight=None,
                                     monotone_constraints=None, n_estimators=None,
                                     n_jobs=None, num_parallel_tree=None,
                                     objective='reg:squarederror',
                                     random_state=None, reg_alpha=None,
                                     reg_lambda=None, scale_pos_weight=None,
                                     subsample=None, tree_method=None,
                                     validate_parameters=False, verbosity=None),
             iid='deprecated', n_jobs=-1,
             param_grid={'max_depth': [1, 2, 3, 4, 5, 6, 7],
                          'n_estimators': [10, 50, 100, 200, 300, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring=make_scorer(rmse), verbose=0)
```

In [22]:

```
results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['mean_test_score'])
results.head(10)
```

Out[22]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max
10	41.490293	0.266281	0.196482	0.000007	
11	61.353465	0.035904	0.297704	0.066322	
9	29.683720	0.074301	0.225896	0.009475	
15	37.251222	0.177525	0.176029	0.001496	
14	20.549191	0.163064	0.190491	0.017952	
19	13.498538	0.030419	0.157579	0.004987	
20	24.495521	0.069062	0.166308	0.003744	
8	16.334844	0.018450	0.192485	0.007979	
13	10.681448	0.079787	0.162067	0.001496	
16	51.755690	0.256325	0.217418	0.010971	

In [23]:

```
# initialize Our first XGBoost model...
xgb_bsl = xgb.XGBRegressor(silent=False, n_jobs=18, random_state=15, max_dept
train_results, test_results = run_xgboost(xgb_bsl, x_train, y_train, x_test,

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

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

Training the model..

Done. Time taken : 0:00:18.915509

Done

Evaluating the model with TRAIN data...

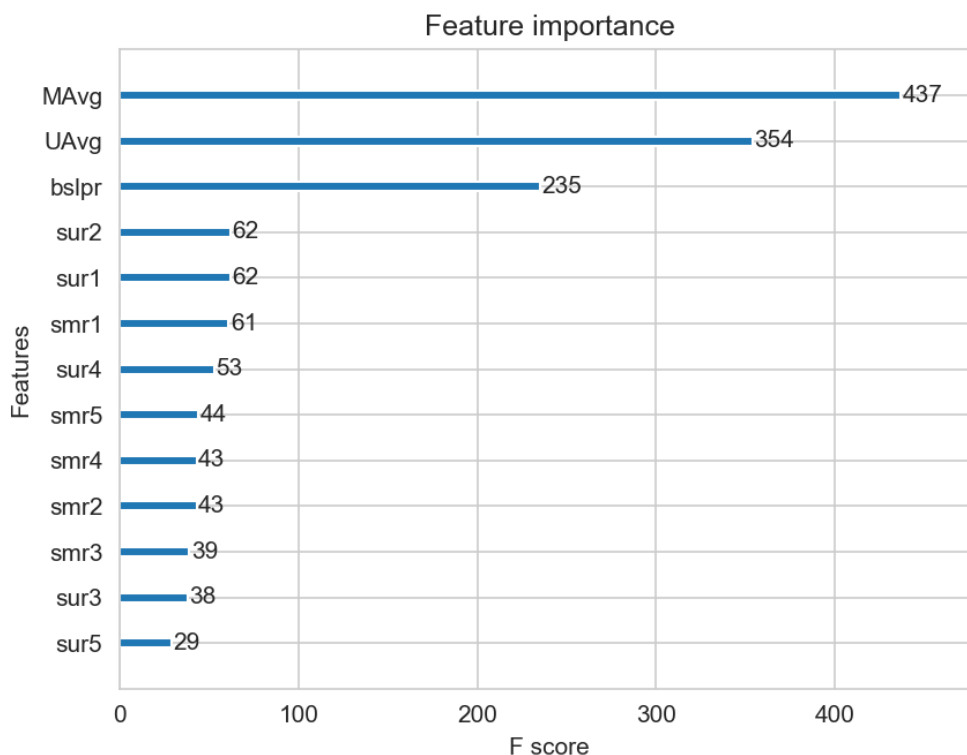
Evaluating Test data

TEST DATA

RMSE : 1.080313427655194

MAPE : 34.798780412271

<IPython.core.display.Javascript object>



4.4.4 Surprise KNNBaseline predictor

In [24]:

```
from surprise import KNNBaseline
```

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

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

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

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

- _Notations follows same as above (user user based predicted rating) _

4.4.4.1 Surprise KNNBaseline with user user similarities

In [25]:

```
# we specify , how to compute similarities and what to consider with sim_opti
sim_options = {'user_based' : True,
               'name': 'pearson_baseline',
               'shrinkage': 100,
               'min_support': 2
              }
# we keep other parameters like regularization parameter and learning_rate as
bsl_options = {'method': 'sgd'}

knn_bsl_u = KNNBaseline(k=40, sim_options = sim_options, bsl_options = bsl_op
knn_bsl_u_train_results, knn_bsl_u_test_results = run_surprise(knn_bsl_u, tra

# 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 : 0:01:26.285205

Evaluating the model with train data..

time taken : 0:03:24.803363

Train Data

RMSE : 0.381999862577986

MAPE : 10.70788648951525

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.236369

Test Data

RMSE : 1.0777732968071752

MAPE : 35.067375899427915

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:04:51.325934

4.4.4.2 Surprise KNNBaseline with movie movie similarities

In [26]:

```
# we specify , how to compute similarities and what to consider with sim_opti
# 'user_based' : Fals => this considers the similarities of movies instead of

sim_options = {'user_based' : False,
               'name': 'pearson_baseline',
               'shrinkage': 100,
               'min_support': 2
              }
# we keep other parameters like regularization parameter and Learning_rate as
bsl_options = {'method': 'sgd'}

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

knn_bsl_m_train_results, knn_bsl_m_test_results = run_surprise(knn_bsl_m, tra

# 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 : 0:00:02.737941

Evaluating the model with train data..

time taken : 0:00:17.510540

Train Data

RMSE : 0.40392992817004597

MAPE : 10.992348015616841

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.101757

Test Data

RMSE : 1.077805003654094

MAPE : 35.06784453290757

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:00:20.351202

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

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

__Preparing Train data __

In [27]:

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

reg_train.head(2)
```

Out[27]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3
0	174683	10	3.540551	4.0	5.0	4.0	3.0	4.0	3.0	3.0	3.0
1	233949	10	3.540551	4.0	4.0	5.0	5.0	5.0	2.0	3.0	4.0

__Preparing Test data __

In [28]:

```
reg_test_df['knn_bs1_u'] = models_evaluation_test['knn_bs1_u']['predictions']
reg_test_df['knn_bs1_m'] = models_evaluation_test['knn_bs1_m']['predictions']

reg_test_df.head(2)
```

Out[28]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur!
0	808635	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551
1	898730	71	3.540551	3.540551	3.540551	3.540551	3.540551	3.540551

In [29]:

```
# 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']
```

In [30]:

```
parameters = {
    'max_depth' : [1, 2, 3, 4, 5, 6, 7],
    'n_estimators' : [10, 50, 100, 200, 300, 500]
}
xgbdt = xgb.XGBRegressor(n_jobs=18)
clf = GridSearchCV(xgbdt, parameters, cv=2, scoring = scorer, return_train_score=True)
clf.fit(x_train, y_train)
```

Out[30]:

```
GridSearchCV(cv=2, error_score=nan,
             estimator=XGBRegressor(base_score=None, booster=None,
                                     colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None, gamma=None,
                                     gpu_id=None, importance_type=None,
                                     interaction_constraints=None,
                                     learning_rate=None, max_depth=None,
                                     min_child_weight=None, missing=nan,
                                     monotone_constraints=None, n_estimators=300,
                                     objective='reg:squarederror',
                                     random_state=None, reg_alpha=None,
                                     reg_lambda=None, scale_pos_weight=None,
                                     subsample=None, tree_method=None,
                                     validate_parameters=False,
                                     verbosity=None),
             iid='deprecated', n_jobs=-1,
             param_grid={'max_depth': [1, 2, 3, 4, 5, 6, 7],
                          'n_estimators': [10, 50, 100, 200, 300, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring=make_scorer(rmse), verbose=0)
```

In [31]:

```
results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['mean_test_score'])
results.head(10)
```

Out[31]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max
10	46.928867	0.032426	0.203456	0.011968	
9	33.774400	0.271287	0.201461	0.002992	
11	70.032739	0.342584	0.253838	0.008494	
8	17.990468	0.059839	0.189505	0.011980	
14	23.482383	0.105718	0.204952	0.002493	
13	12.862469	0.032666	0.180517	0.002992	
5	50.378150	0.069814	0.233875	0.025432	
19	15.399467	0.032912	0.175032	0.008477	
4	32.702264	0.056848	0.244346	0.005984	
15	42.361581	0.215424	0.251827	0.000499	

In [32]:

```
# declare the model
xgb_knn_bsl = xgb.XGBRegressor(silent=False, n_jobs=18, random_state=15, max_
train_results, test_results = run_xgboost(xgb_knn_bsl, x_train, y_train, x_te

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

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

Training the model..

Done. Time taken : 0:00:12.891729

Done

Evaluating the model with TRAIN data...

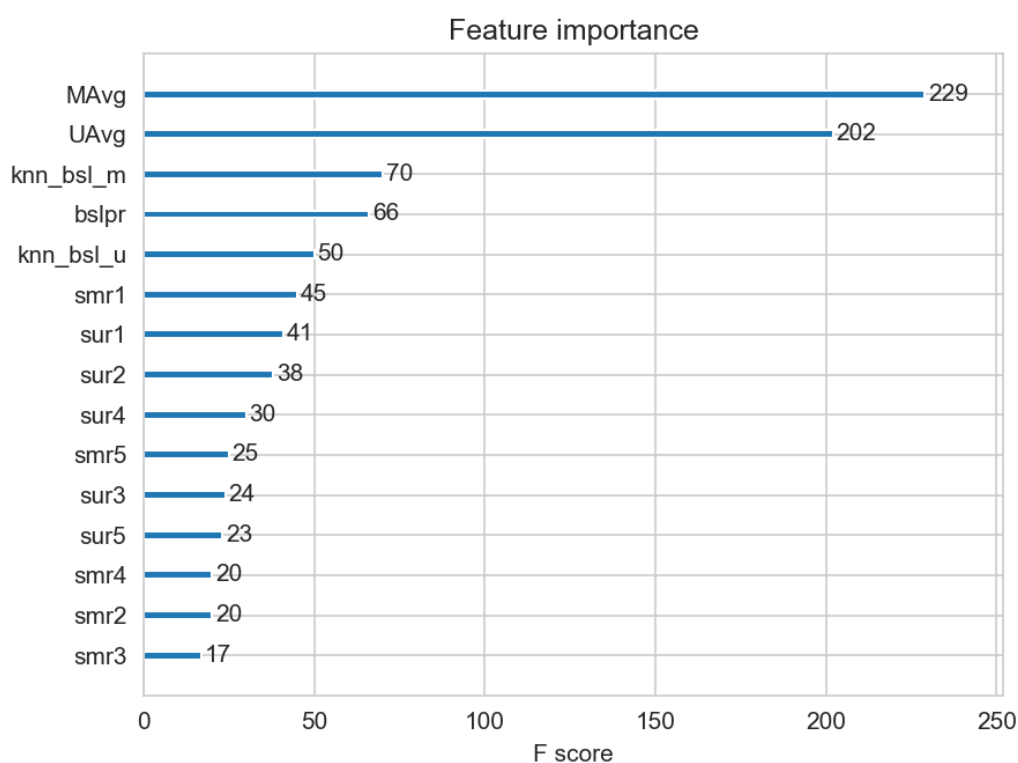
Evaluating Test data

TEST DATA

RMSE : 1.079729034819609

MAPE : 34.870587042998615

<IPython.core.display.Javascript object>



4.4.6 Matrix Factorization Techniques

4.4.6.1 SVD Matrix Factorization User Movie interactions

In [33]:

```
from surprise import SVD
```

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

- __ Predicted Rating : __
 - - $\hat{r}_{ui} = \mu + b_u + b_i + q_i^T p_u$
 - q_i - Representation of item(movie) in latent factor space
 - p_u - Representation of user in new latent factor space
- A BASIC MATRIX FACTORIZATION MODEL in [https://datajobs.com/data-science-repo/Recommender-Systems-\[Netflix\].pdf](https://datajobs.com/data-science-repo/Recommender-Systems-[Netflix].pdf) (<https://datajobs.com/data-science-repo/Recommender-Systems-%5BNetflix%5D.pdf>)
- Optimization problem with user item interactions and regularization (to avoid overfitting)
 - - $$\sum_{r_{ui} \in R_{train}} (r_{ui} - \hat{r}_{ui})^2 + \lambda (b_i^2 + b_u^2 + ||q_i||^2 + ||p_u||^2)$$

In [34]:

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

# 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 : 0:00:11.015325

Evaluating the model with train data..

time taken : 0:00:01.942734

Train Data

RMSE : 0.6604275824313268

MAPE : 19.91824112446639

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.112426

Test Data

RMSE : 1.0777301884881167

MAPE : 34.99321365749742

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:00:13.070485

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

In [35]:

```
from surprise import SVDpp
```

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

- __ Predicted Rating : __

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

- 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)**

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

In [36]:

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

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

Training the model...

```
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 : 0:04:12.630981

Evaluating the model with train data..

time taken : 0:00:11.559573

Train Data

RMSE : 0.6242943697858871

MAPE : 18.218713271480112

adding train results in the dictionary..

Evaluating for test data...

time taken : 0:00:00.209414

Test Data

RMSE : 1.0781584680497336

MAPE : 34.90014147266236

storing the test results in test dictionary...

Total time taken to run this algorithm : 0:04:24.400965

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

Preparing Train data

In [37]:

```
# 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[37]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2
0	174683	10	3.540551	4.0	5.0	4.0	3.0	4.0	3.0	3.0
1	233949	10	3.540551	4.0	4.0	5.0	5.0	5.0	2.0	3.0

2 rows × 21 columns

__Preparing Test data __

In [38]:

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

reg_test_df.head(2)
```

Out[38]:

	user	movie	GAvg	sur1	sur2	sur3	sur4
0	808635	71	3.540551	3.540551	3.540551	3.540551	3.540551
1	898730	71	3.540551	3.540551	3.540551	3.540551	3.540551

2 rows × 8 columns

In [39]:

```
# 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']
```

In [41]:

```
from sklearn.model_selection import GridSearchCV

parameters = {
    'max_depth' : [1, 2, 3, 4, 5, 6, 7],
    'n_estimators' : [10, 50, 100, 200, 300, 500]
}
xgbdt = xgb.XGBRegressor(n_jobs=18)
clf = GridSearchCV(xgbdt, parameters, cv=2, scoring = scorer, return_train_score=True)
clf.fit(x_train, y_train)
```

Out[41]:

```
GridSearchCV(cv=2, error_score=nan,
             estimator=XGBRegressor(base_score=None, booster=None,
                                     colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None, gamma=None,
                                     gpu_id=None, importance_type=None,
                                     interaction_constraints=None,
                                     learning_rate=None, max_depth=None,
                                     max_delta_step=None, max_leaf_child_weight=None,
                                     min_child_weight=None, missing=nan, monotone_constraints=None,
                                     n_estimators=None, n_jobs=None, num_parallel_tree=None,
                                     objective='reg:squarederror', random_state=None, reg_alpha=None,
                                     reg_lambda=None, scale_pos_weight=None, subsample=None, tree_method=None,
                                     validate_parameters=False, verbosity=None),
             iid='deprecated', n_jobs=-1,
             param_grid={'max_depth': [1, 2, 3, 4, 5, 6, 7],
                          'n_estimators': [10, 50, 100, 200, 300, 500]}),
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring=make_scorer(rmse), verbose=0)
```


In [42]:

```
results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['mean_test_score'])
results.head(10)
```

Out[42]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max
9	36.209414	0.047872	0.216422	0.000998	
10	51.207968	0.038896	0.281248	0.018950	
8	19.215764	0.000498	0.231880	0.017453	
11	77.429720	0.134152	0.270277	0.006981	
14	25.781187	0.107213	0.191488	0.017952	
5	54.589355	0.059341	0.273768	0.026429	
19	17.534304	0.071808	0.189992	0.007480	
13	13.865048	0.146608	0.179520	0.004987	
4	34.913402	0.066323	0.224400	0.017952	
15	47.232098	0.289241	0.226908	0.032399	

In [44]:

```
xgb_final = xgb.XGBRegressor(silent=False, n_jobs=18, random_state=15, max_de  
train_results, test_results = run_xgboost(xgb_final, x_train, y_train, x_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 : 0:00:09.164001

Done

Evaluating the model with TRAIN data...

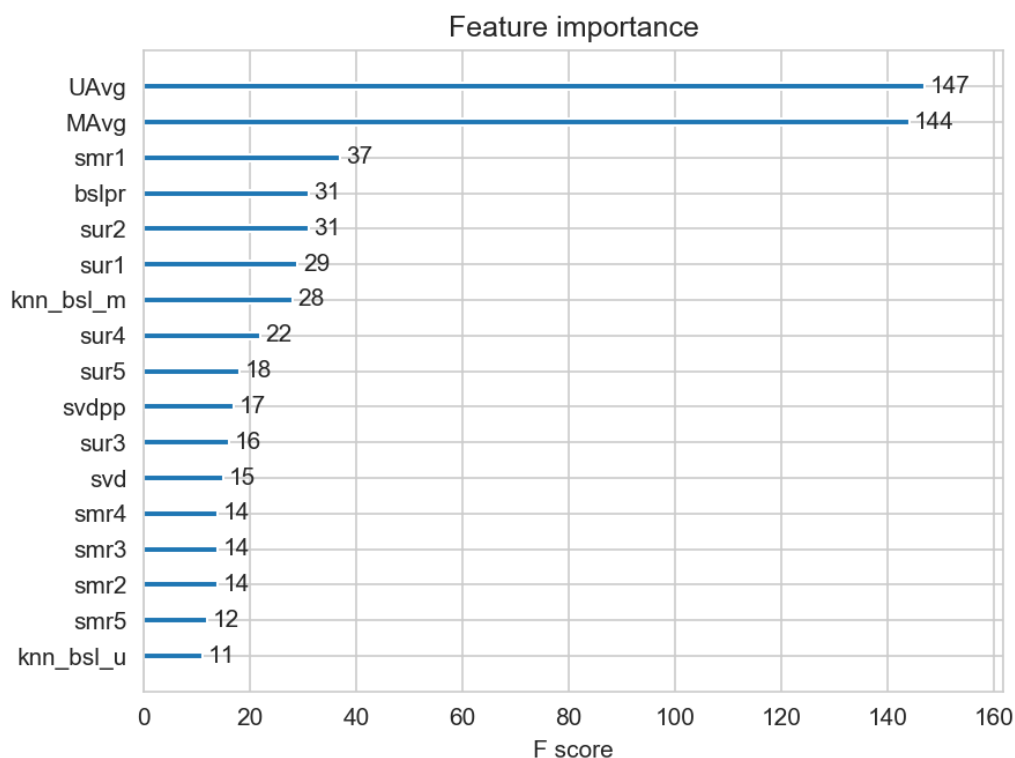
Evaluating Test data

TEST DATA

RMSE : 1.079329619935183

MAPE : 34.91638845389938

<IPython.core.display.Javascript object>



4.4.8 XgBoost with Surprise Baseline + Surprise KNNbaseline + MF Techniques

In [45]:

```
# 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']
```

In [46]:

```
parameters = {
    'max_depth' : [1, 2, 3, 4, 5, 6, 7],
    'n_estimators' : [10, 50, 100, 200, 300, 500]
}
xgbdt = xgb.XGBRegressor(n_jobs=18)
clf = GridSearchCV(xgbdt, parameters, cv=2, scoring = scorer, return_train_score=True)
clf.fit(x_train, y_train)
```

Out[46]:

```
GridSearchCV(cv=2, error_score=nan,
             estimator=XGBRegressor(base_score=None, booster=None,
                                     colsample_bylevel=None,
                                     colsample_bynode=None,
                                     colsample_bytree=None, gamma=None,
                                     gpu_id=None, importance_type=None,
                                     interaction_constraints=None,
                                     learning_rate=None, max_depth=None,
                                     min_child_weight=None, missing=nan,
                                     monotone_constraints=None, n_estimators=300,
                                     objective='reg:squarederror',
                                     random_state=None, reg_alpha=None,
                                     reg_lambda=None, scale_pos_weight=None,
                                     subsample=None, tree_method=None,
                                     validate_parameters=False,
                                     verbosity=None),
             iid='deprecated', n_jobs=-1,
             param_grid={'max_depth': [1, 2, 3, 4, 5, 6, 7],
                         'n_estimators': [10, 50, 100, 200, 300, 500]},
             pre_dispatch='2*n_jobs', refit=True, return_train_score=True,
             scoring=make_scorer(rmse), verbose=0)
```

In [49]:

```
results = pd.DataFrame.from_dict(clf.cv_results_)
results = results.sort_values(['mean_test_score'])
results.head(10)
```

Out[49]:

	mean_fit_time	std_fit_time	mean_score_time	std_score_time	param_max
1	4.844545	0.027427	0.085771	0.002992	
2	8.586039	0.046875	0.094248	0.003491	
3	15.501574	0.036901	0.108709	0.001995	
4	22.153798	0.025930	0.115192	0.006483	
5	34.369704	0.487702	0.132898	0.007234	
7	6.553974	0.028424	0.093749	0.009974	
8	11.935603	0.062341	0.114700	0.002003	
13	8.200598	0.128157	0.084773	0.004987	
19	10.101579	0.130651	0.081061	0.008256	
9	22.416111	0.080784	0.111953	0.005239	

In [50]:

```
xgb_all_models = xgb.XGBRegressor(silent=False, n_jobs=18, random_state=15, r
train_results, test_results = run_xgboost(xgb_all_models, x_train, y_train, >

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

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

Training the model..

Done. Time taken : 0:00:01.353204

Done

Evaluating the model with TRAIN data...

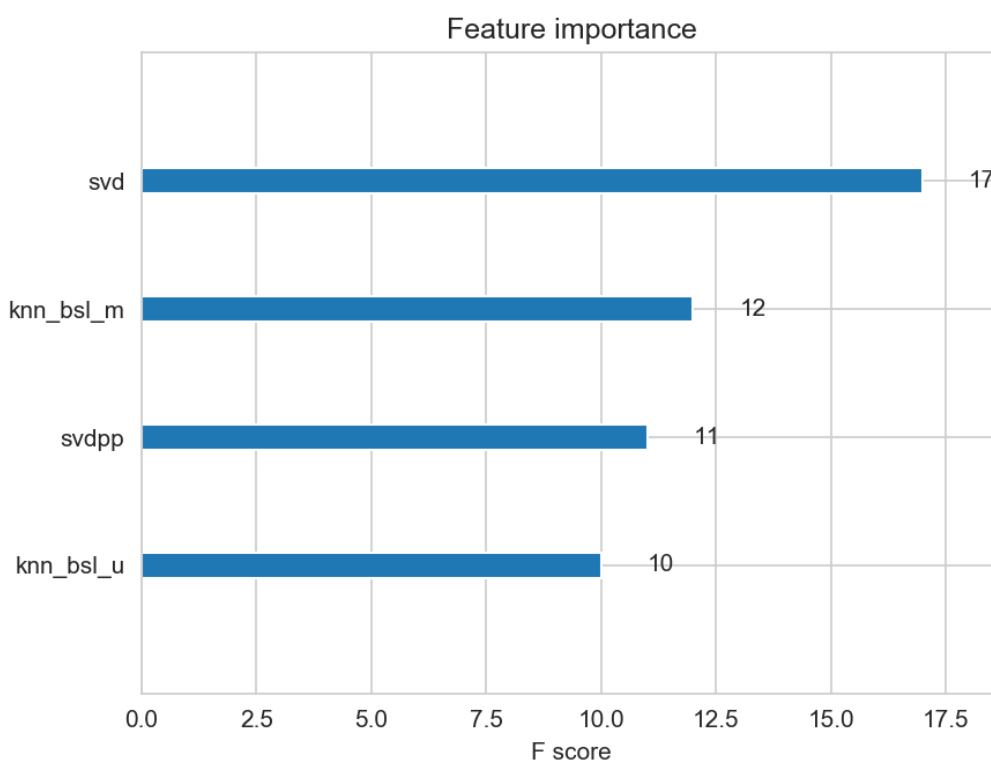
Evaluating Test data

TEST DATA

RMSE : 1.084507707899312

MAPE : 35.46699593611435

<IPython.core.display.Javascript object>



4.5 Comparision between all models

In [51]:

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

Out[51]:

```
svd                1.0777301884881167
knn_bsl_u          1.0777732968071752
knn_bsl_m          1.077805003654094
svdpp              1.0781584680497336
xgb_final          1.079329619935183
first_algo         1.0795164032626583
xgb_knn_bsl        1.079729034819609
xgb_bsl            1.080313427655194
bsl_algo           1.0818422475962859
xgb_all_models     1.084507707899312
Name: rmse, dtype: object
```

Summary

- Used 15K Users and 1.5K Movies to train the models
- Performed Hyperparameter tuning using GridSearchCV
- Obtained best values for the hyperparameters max_depth, n_estimators
- As a result of hyperparameter tuning, we obtained better 'rmse' values