



# 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](http://surprise.readthedocs.io/en/stable/getting_started.html) (we use many models from this library)

- installing surprise: <https://github.com/NicolasHug/Surprise#installation>
- Research paper: <http://courses.ischool.berkeley.edu/i290-dm/s11/SECURE/a1-koren.pdf> (most of our work was inspired by this paper)
- SVD Decomposition : <https://www.youtube.com/watch?v=P5mlg91as1c>

## 1.4 Real world/Business Objectives and constraints

Objectives:

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

Constraints:

1. Some form of interpretability.

## 2. Machine Learning Problem

### 2.1 Data

#### 2.1.1 Data Overview

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

Data files :

- combined\_data\_1.txt
- combined\_data\_2.txt
- combined\_data\_3.txt
- combined\_data\_4.txt
- movie\_titles.csv

The first line of each file [combined\_data\_1.txt, combined\_data\_2.txt, combined\_data\_3.txt, combined\_data\_4.txt] contains the movie id followed by a colon. Each subsequent line in the file corresponds to a rating from a customer and its date in the following format:

CustomerID,Rating,Date

MovieIDs range from 1 to 17770 sequentially.

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

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

Dates have the format YYYY-MM-DD.

## 2.1.2 Example Data point

1:

1488844,3,2005-09-06

822109,5,2005-05-13

885013,4,2005-10-19

30878,4,2005-12-26

823519,3,2004-05-03

893988,3,2005-11-17

124105,4,2004-08-05

1248029,3,2004-04-22

1842128,4,2004-05-09

2238063,3,2005-05-11

1503895,4,2005-05-19

2207774,5,2005-06-06

2590061,3,2004-08-12

2442,3,2004-04-14

543865,4,2004-05-28

1209119,4,2004-03-23

804919,4,2004-06-10

1086807,3,2004-12-28  
1711859,4,2005-05-08  
372233,5,2005-11-23  
1080361,3,2005-03-28  
1245640,3,2005-12-19  
558634,4,2004-12-14  
2165002,4,2004-04-06  
1181550,3,2004-02-01  
1227322,4,2004-02-06  
427928,4,2004-02-26  
814701,5,2005-09-29  
808731,4,2005-10-31  
662870,5,2005-08-24  
337541,5,2005-03-23  
786312,3,2004-11-16  
1133214,4,2004-03-07  
1537427,4,2004-03-29  
1209954,5,2005-05-09  
2381599,3,2005-09-12  
525356,2,2004-07-11  
1910569,4,2004-04-12  
2263586,4,2004-08-20  
2421815,2,2004-02-26  
1009622,1,2005-01-19  
1481961,2,2005-05-24  
401047,4,2005-06-03  
2179073,3,2004-08-29  
1434636,3,2004-05-01  
93986,5,2005-10-06  
1308744,5,2005-10-29  
2647871,4,2005-12-30  
1905581,5,2005-08-16  
2508819,3,2004-05-18  
1578279,1,2005-05-19  
1159695,4,2005-02-15

2588432,3,2005-03-31  
2423091,3,2005-09-12  
470232,4,2004-04-08  
2148699,2,2004-06-05  
1342007,3,2004-07-16  
466135,4,2004-07-13  
2472440,3,2005-08-13  
1283744,3,2004-04-17  
1927580,4,2004-11-08  
716874,5,2005-05-06  
4326,4,2005-10-29

## 2.2 Mapping the real world problem to a Machine Learning Problem

### 2.2.1 Type of Machine Learning Problem

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

The given problem is a Recommendation problem

It can also seen as a Regression problem

### 2.2.2 Performance metric

- Mean Absolute Percentage Error: [https://en.wikipedia.org/wiki/Mean\\_absolute\\_percentage\\_error](https://en.wikipedia.org/wiki/Mean_absolute_percentage_error)
- Root Mean Square Error: [https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](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]:

```
1 # this is just to know how much time will it take to run this entire ipython notebook
2 from datetime import datetime
3 # globalstart = datetime.now()
4 import pandas as pd
5 import numpy as np
6 import matplotlib
7 matplotlib.use('nbagg')
8
9 import matplotlib.pyplot as plt
10 plt.rcParams.update({'figure.max_open_warning': 0})
11
12 import seaborn as sns
13 sns.set_style('whitegrid')
14 import os
15 from scipy import sparse
16 from scipy.sparse import csr_matrix
17 import surprise
18 from sklearn.decomposition import TruncatedSVD
19 from sklearn.metrics.pairwise import cosine_similarity
20 import random
21 %matplotlib inline
```

### 3. Exploratory Data Analysis

#### 3.1 Preprocessing

##### 3.1.1 Converting / Merging whole data to required format: u\_i, m\_j, r\_ij

In [2]:

```
1 start = datetime.now()
2 if not os.path.isfile('data.csv'):
3     # Create a file 'data.csv' before reading it
4     # Read all the files in netflix and store them in one big file('data.csv')
5     # We're reading from each of the four files and appending each rating to a global file 'train.csv'
6     data = open('data.csv', mode='w')
7
8     row = list()
9     files=['data_folder/combined_data_1.txt','data_folder/combined_data_2.txt',
10           'data_folder/combined_data_3.txt', 'data_folder/combined_data_4.txt']
11    for file in files:
12        print("Reading ratings from {}...".format(file))
13        with open(file) as f:
14            for line in f:
15                del row[:] # you don't have to do this.
16                line = line.strip()
17                if line.endswith(':'):
18                    # All below are ratings for this movie, until another movie appears.
19                    movie_id = line.replace(':', '')
20                else:
21                    row = [x for x in line.split(',')]
22                    row.insert(0, movie_id)
23                    data.write(','.join(row))
24                    data.write('\n')
25            print("Done.\n")
26        data.close()
27    print('Time taken :', datetime.now() - start)
```

Time taken : 0:00:00.000809

In [3]:

```
1 %%time
2 print("creating the dataframe from data.csv file..")
3 df = pd.read_csv('data.csv', sep=',',
4                   names=['movie', 'user', 'rating', 'date'])
5 df.date = pd.to_datetime(df.date)
6 print('Done.\n')
7
8 # we are arranging the ratings according to time.
9 print('Sorting the dataframe by date..')
10 df.sort_values(by='date', inplace=True)
11 print('Done..')
```

creating the dataframe from data.csv file..

Done.

Sorting the dataframe by date..

Done..

CPU times: user 1min 33s, sys: 11.4 s, total: 1min 44s

Wall time: 1min 46s

In [4]:

```
1 df.head(5)
```

Out[4]:

	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 [5]: 1 df.describe()['rating']
```

```
Out[5]: count    1.004805e+08
mean     3.604290e+00
std      1.085219e+00
min     1.000000e+00
25%     3.000000e+00
50%     4.000000e+00
75%     4.000000e+00
max     5.000000e+00
Name: rating, dtype: float64
```

### 3.1.2 Checking for NaN values

```
In [6]: 1 %%time
2 # just to make sure that all Nan containing rows are deleted..
3 print("No of Nan values in our dataframe : ", sum(df.isnull().any()))
```

```
No of Nan values in our dataframe :  0
CPU times: user 7.26 s, sys: 200 ms, total: 7.46 s
Wall time: 7.46 s
```

### 3.1.3 Removing Duplicates

```
In [7]: 1 dup_bool = df.duplicated(['movie','user','rating'])
2 dups = sum(dup_bool) # by considering all columns..( including timestamp)
3 print("There are {} duplicate rating entries in the data..".format(dups))
```

```
There are 0 duplicate rating entries in the data..
```

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

In [8]:

```
1 print("Total data ")
2 print("-"*50)
3 print("\nTotal no of ratings :",df.shape[0])
4 print("Total No of Users   :", len(np.unique(df.user)))
5 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 [9]:

```
1 if not os.path.isfile('train.csv'):
2     # create the dataframe and store it in the disk for offline purposes..
3     df.iloc[:int(df.shape[0]*0.80)].to_csv("train.csv", index=False)
4
5 if not os.path.isfile('test.csv'):
6     # create the dataframe and store it in the disk for offline purposes..
7     df.iloc[int(df.shape[0]*0.80):].to_csv("test.csv", index=False)
8
9 train_df = pd.read_csv("train.csv", parse_dates=['date'])
10 test_df = pd.read_csv("test.csv")
```

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

In [10]:

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

Training data

---

```
Total no of ratings : 80384405
Total No of Users : 405041
Total No of movies : 17424
```

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

In [11]:

```
1 print("Test data ")
2 print("-"*50)
3 print("\nTotal no of ratings :",test_df.shape[0])
4 print("Total No of Users : ", len(np.unique(test_df.user)))
5 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 [12]: 1 # method to make y-axis more readable
2 def human(num, units = 'M'):
3     units = units.lower()
4     num = float(num)
5     if units == 'k':
6         return str(num/10**3) + " K"
7     elif units == 'm':
8         return str(num/10**6) + " M"
9     elif units == 'b':
10        return str(num/10**9) + " B"
```

### 3.3.1 Distribution of ratings

In [14]:

```
1 sns.set_style('darkgrid')
2 fig, ax=plt.subplots()
3 plt.rcParams["figure.figsize"] = (10,7)
4 plt.title('Distribution of ratings over Training dataset', fontsize=15)
5 sns.countplot(train_df.rating)
6 ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
7 ax.set_ylabel('No. of Ratings(Millions)')
8
9 plt.show()
```



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

In [15]:

```
1 # It is used to skip the warning ''SettingWithCopyWarning''..  
2 pd.options.mode.chained_assignment = None # default='warn'  
3  
4 train_df['day_of_week'] = train_df.date.dt.weekday_name  
5  
6 train_df.tail()
```

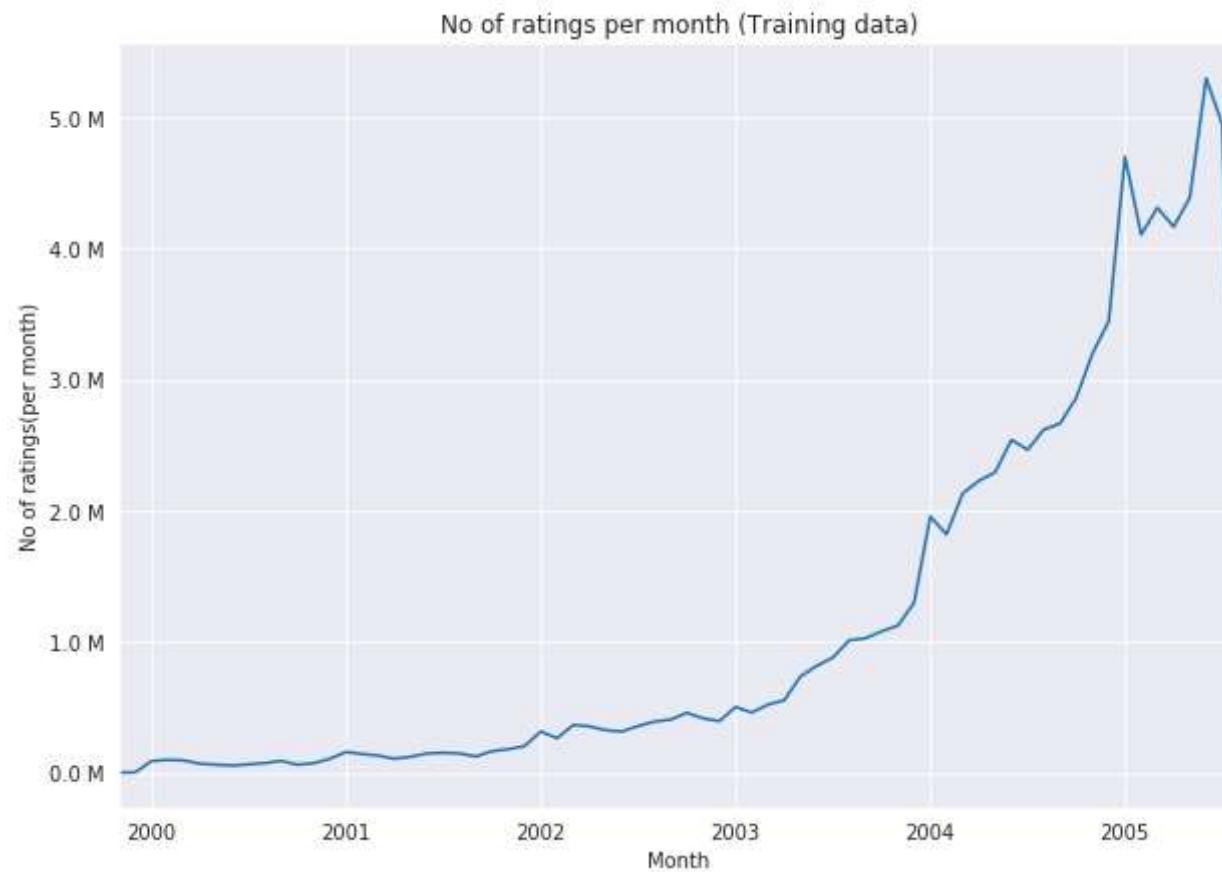
Out[15]:

	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 [16]:

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



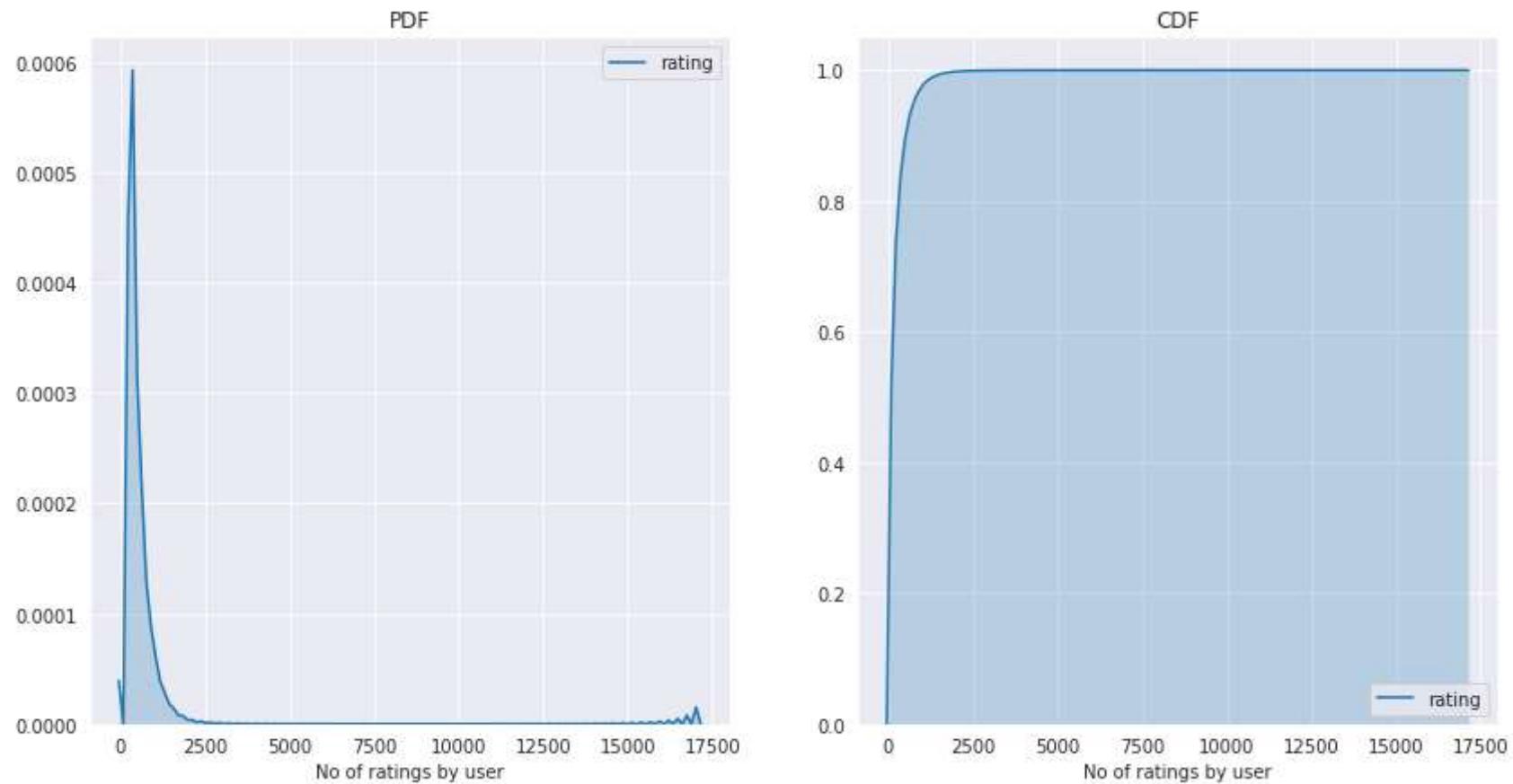
### 3.3.3 Analysis on the Ratings given by user

```
In [17]: 1 no_of_rated_movies_per_user = train_df.groupby(by='user')['rating'].count().sort_values(ascending=False)
          2
          3 no_of_rated_movies_per_user.head()
```

```
Out[17]: user
305344    17112
2439493   15896
387418    15402
1639792   9767
1461435   9447
Name: rating, dtype: int64
```

In [18]:

```
1 fig = plt.figure(figsize=plt.figaspect(.5))
2 sns.set_style('darkgrid')
3 ax1 = plt.subplot(121)
4 sns.kdeplot(no_of_rated_movies_per_user, shade=True, ax=ax1)
5 plt.xlabel('No of ratings by user')
6 plt.title("PDF")
7
8 ax2 = plt.subplot(122)
9 sns.kdeplot(no_of_rated_movies_per_user, shade=True, cumulative=True, ax=ax2)
10 plt.xlabel('No of ratings by user')
11 plt.title('CDF')
12
13 plt.show()
```



```
In [19]: 1 no_of_rated_movies_per_user.describe()
```

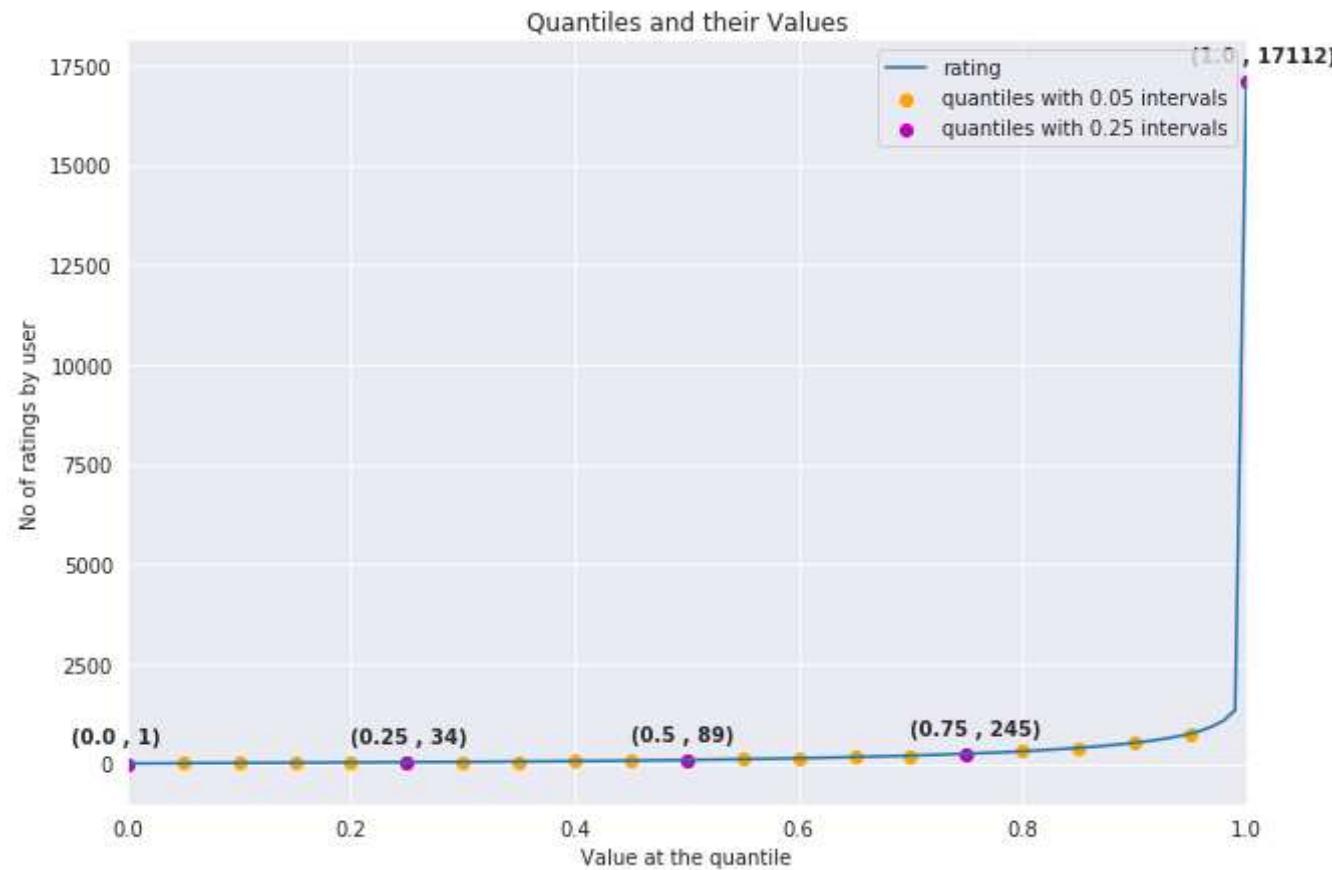
```
Out[19]: count    405041.000000
mean      198.459921
std       290.793238
min       1.000000
25%      34.000000
50%      89.000000
75%     245.000000
max     17112.000000
Name: rating, dtype: float64
```

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

```
In [20]: 1 quantiles = no_of_rated_movies_per_user.quantile(np.arange(0,1.01,0.01), interpolation='higher')
```

In [21]:

```
1 plt.title("Quantiles and their Values")
2 quantiles.plot()
3 # quantiles with 0.05 difference
4 plt.scatter(x=quantiles.index[::-5], y=quantiles.values[::-5], c='orange', label="quantiles with 0.05 intervals")
5 # quantiles with 0.25 difference
6 plt.scatter(x=quantiles.index[::-25], y=quantiles.values[::-25], c='m', label = "quantiles with 0.25 intervals")
7 plt.ylabel('No of ratings by user')
8 plt.xlabel('Value at the quantile')
9 plt.legend(loc='best')
10
11 # annotate the 25th, 50th, 75th and 100th percentile values....
12 for x,y in zip(quantiles.index[::-25], quantiles[::-25]):
13     plt.annotate(s="{} , {}".format(x,y), xy=(x,y), xytext=(x-0.05, y+500)
14             ,fontweight='bold')
15
16
17 plt.show()
```



```
In [22]: 1 quantiles[::5]
```

```
Out[22]: 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??**

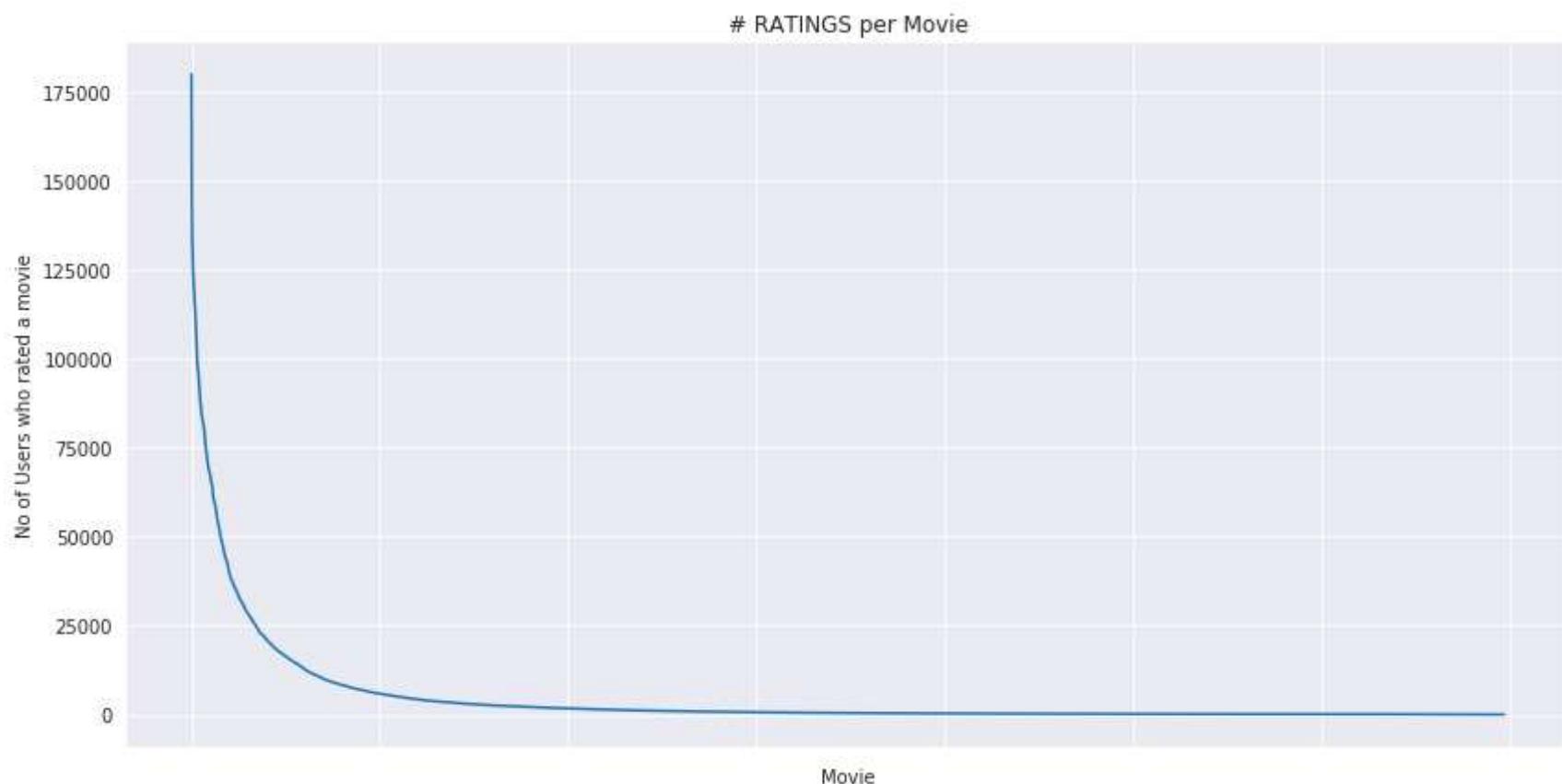
```
In [23]: 1 print('\n No of ratings at last 5 percentile : {}{}'.format(sum(no_of_rated_movies_per_user>= 749)))
```

No of ratings at last 5 percentile : 20305

### 3.3.4 Analysis of ratings of a movie given by a user

In [24]:

```
1 no_of_ratings_per_movie = train_df.groupby(by='movie')['rating'].count().sort_values(ascending=False)
2
3 fig = plt.figure(figsize=plt.figaspect(.5))
4 ax = plt.gca()
5 plt.plot(no_of_ratings_per_movie.values)
6 plt.title('# RATINGS per Movie')
7 plt.xlabel('Movie')
8 plt.ylabel('No of Users who rated a movie')
9 ax.set_xticklabels([])
10
11 plt.show()
```



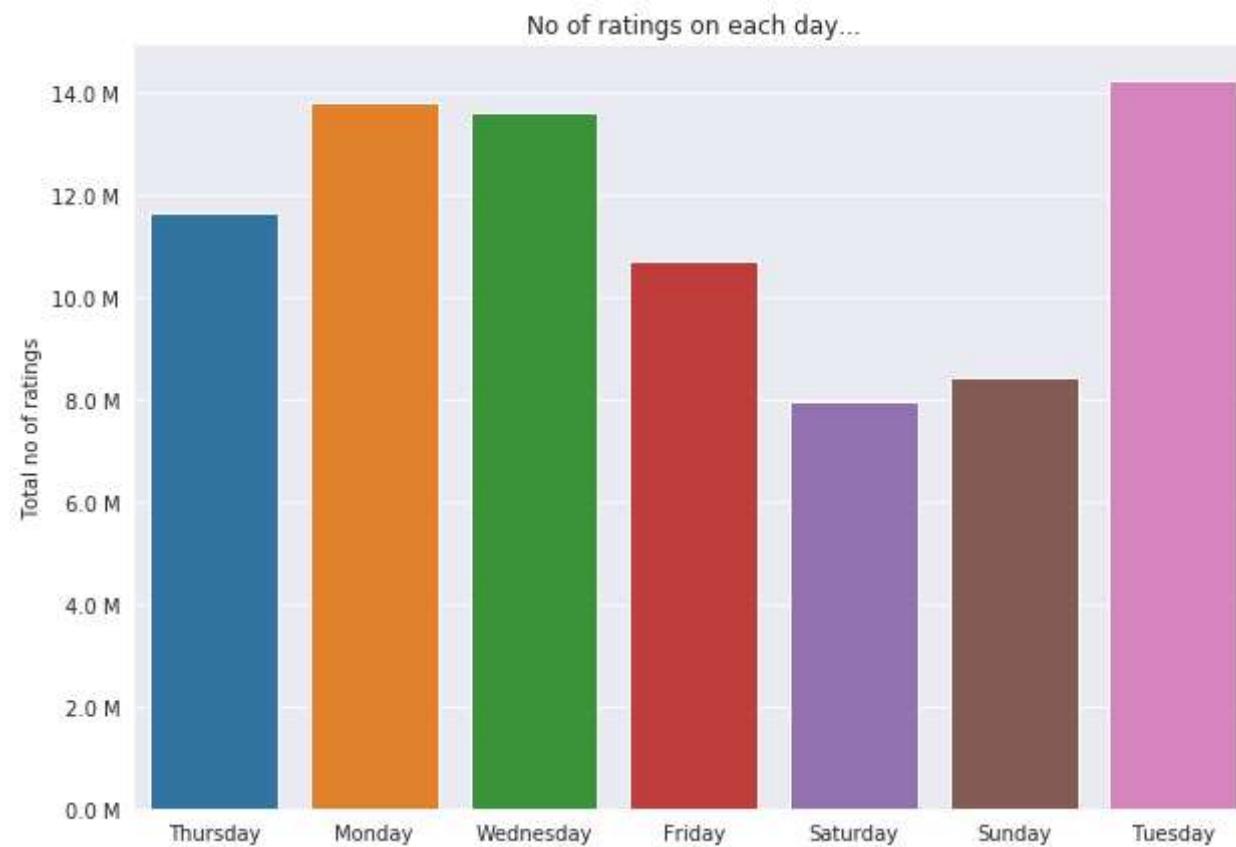
- It is very skewed.. just like number of ratings given per user.

- There are some movies (which are very popular) which are rated by huge number of users.
- But most of the movies (like 90%) got some hundreds of ratings.

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

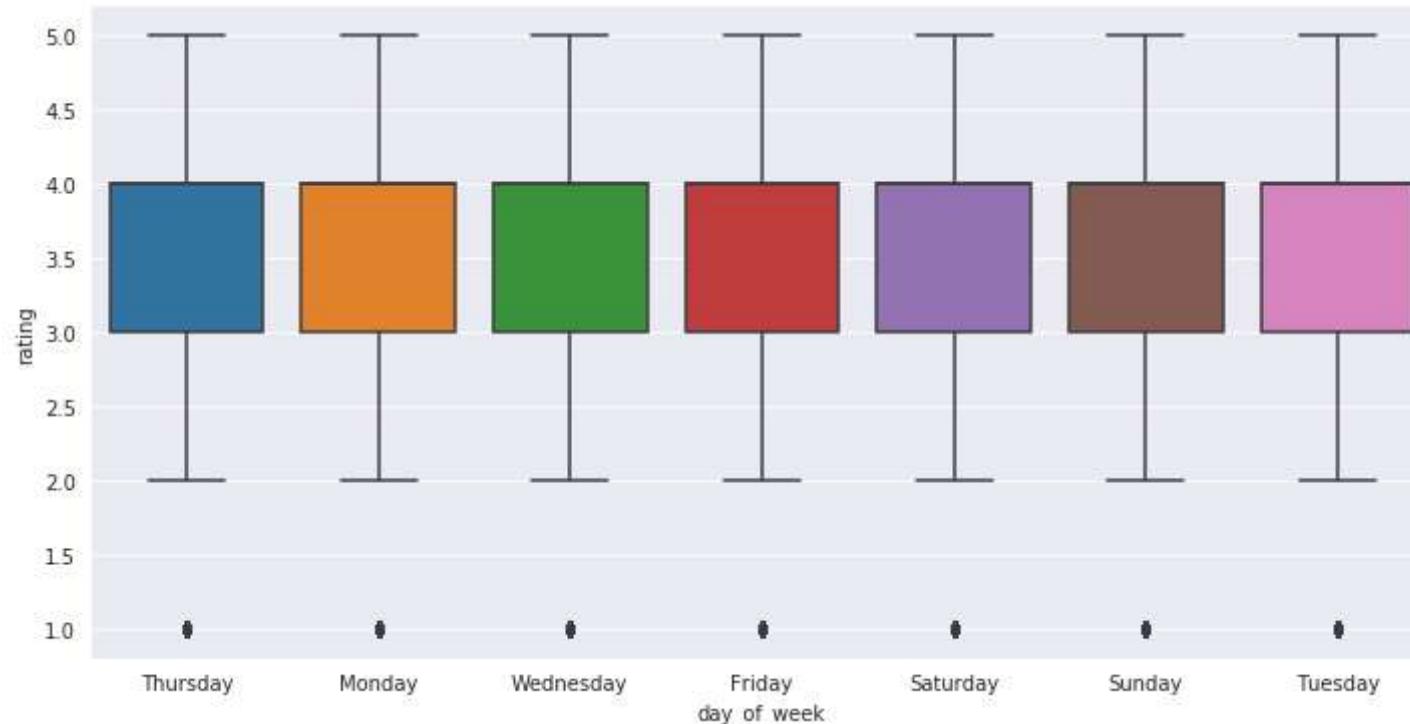
In [25]:

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



In [28]:

```
1 %%time
2 fig = plt.figure(figsize=(12,6))#plt.figaspect(.45))
3 sns.boxplot(y='rating', x='day_of_week', data=train_df)
4 plt.show()
```



CPU times: user 27.7 s, sys: 4.64 s, total: 32.4 s

Wall time: 31.9 s

```
In [29]:  
1 avg_week_df = train_df.groupby(by=['day_of_week'])['rating'].mean()  
2 print(" Average ratings")  
3 print("-"*30)  
4 print(avg_week_df)  
5 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]:

```

1 start = datetime.now()
2 if os.path.isfile('train_sparse_matrix.npz'):
3     print("It is present in your pwd, getting it from disk....")
4     # just get it from the disk instead of computing it
5     train_sparse_matrix = sparse.load_npz('train_sparse_matrix.npz')
6     print("DONE..")
7 else:
8     print("We are creating sparse_matrix from the dataframe..")
9     # create sparse_matrix and store it for after usage.
10    # csr_matrix(data_values, (row_index, col_index), shape_of_matrix)
11    # It should be in such a way that, MATRIX[row, col] = data
12    train_sparse_matrix = sparse.csr_matrix((train_df.rating.values, (train_df.user.values,
13                                              train_df.movie.values)),)
14
15    print('Done. It\'s shape is : (user, movie) : ',train_sparse_matrix.shape)
16    print('Saving it into disk for furthur usage..')
17    # save it into disk
18    sparse.save_npz("train_sparse_matrix.npz", train_sparse_matrix)
19    print('Done..\n')
20
21 print(datetime.now() - start)

```

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

DONE..

0:00:04.318880

### The Sparsity of Train Sparse Matrix

In [3]:

```

1 us,mv = train_sparse_matrix.shape
2 elem = train_sparse_matrix.count_nonzero()
3
4 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]:

```

1 start = datetime.now()
2 if os.path.isfile('test_sparse_matrix.npz'):
3     print("It is present in your pwd, getting it from disk....")
4     # just get it from the disk instead of computing it
5     test_sparse_matrix = sparse.load_npz('test_sparse_matrix.npz')
6     print("DONE..")
7 else:
8     print("We are creating sparse_matrix from the dataframe..")
9     # create sparse_matrix and store it for after usage.
10    # csr_matrix(data_values, (row_index, col_index), shape_of_matrix)
11    # It should be in such a way that, MATRIX[row, col] = data
12    test_sparse_matrix = sparse.csr_matrix((test_df.rating.values, (test_df.user.values,
13                                              test_df.movie.values)))
14
15    print('Done. It\'s shape is : (user, movie) : ',test_sparse_matrix.shape)
16    print('Saving it into disk for furthur usage..')
17    # save it into disk
18    sparse.save_npz("test_sparse_matrix.npz", test_sparse_matrix)
19    print('Done..\n')
20
21 print(datetime.now() - start)

```

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

DONE..

0:00:01.138138

### The Sparsity of Test data Matrix

In [5]:

```

1 us,mv = test_sparse_matrix.shape
2 elem = test_sparse_matrix.count_nonzero()
3
4 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 [6]: 1 # get the user averages in dictionary (key: user_id/movie_id, value: avg rating)
2
3 def get_average_ratings(sparse_matrix, of_users):
4
5     # average ratings of user/axes
6     ax = 1 if of_users else 0 # 1 - User axes, 0 - Movie axes
7
8     # ".A1" is for converting Column_Matrix to 1-D numpy array
9     sum_of_ratings = sparse_matrix.sum(axis=ax).A1
10    # Boolean matrix of ratings ( whether a user rated that movie or not)
11    is_rated = sparse_matrix!=0
12    # no of ratings that each user OR movie..
13    no_of_ratings = is_rated.sum(axis=ax).A1
14
15    # max_user and max_movie ids in sparse matrix
16    u,m = sparse_matrix.shape
17    # create a dictionary of users and their average ratings..
18    average_ratings = { i : sum_of_ratings[i]/no_of_ratings[i]
19                        for i in range(u if of_users else m)
20                        if no_of_ratings[i] !=0}
21
22    # return that dictionary of average ratings
23    return average_ratings
```

### 3.3.7.1 finding global average of all movie ratings

```
In [7]: 1 train_averages = dict()
2 # get the global average of ratings in our train set.
3 train_global_average = train_sparse_matrix.sum()/train_sparse_matrix.count_nonzero()
4 train_averages['global'] = train_global_average
5 train_averages
```

Out[7]: {'global': 3.582890686321557}

### 3.3.7.2 finding average rating per user

```
In [8]: 1 train_averages['user'] = get_average_ratings(train_sparse_matrix, of_users=True)
2 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 [9]: 1 train_averages['movie'] = get_average_ratings(train_sparse_matrix, of_users=False)
2 print('\n Average rating of movie 15 :',train_averages['movie'][15])
```

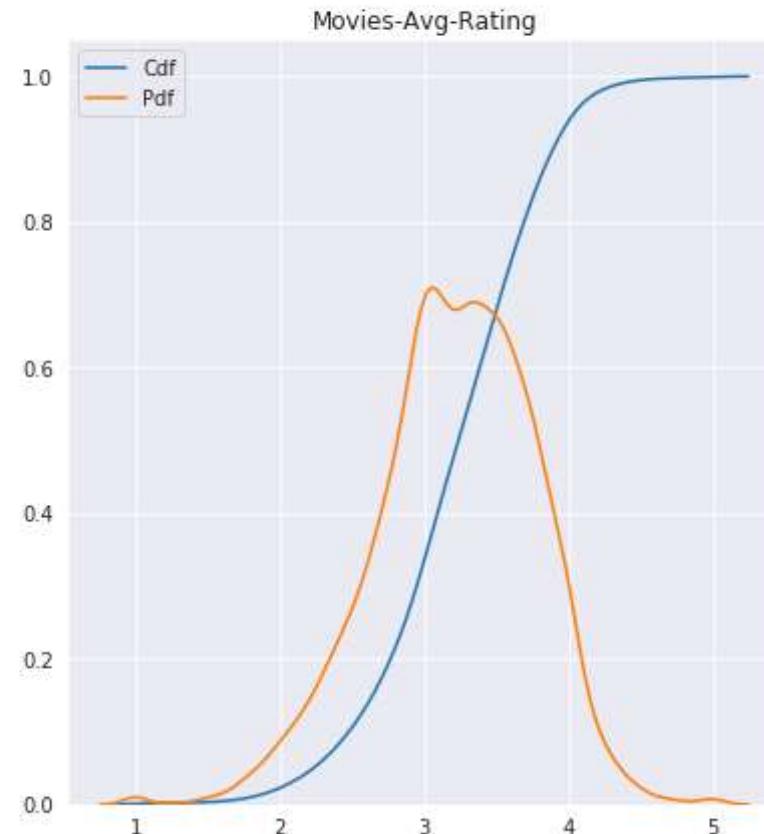
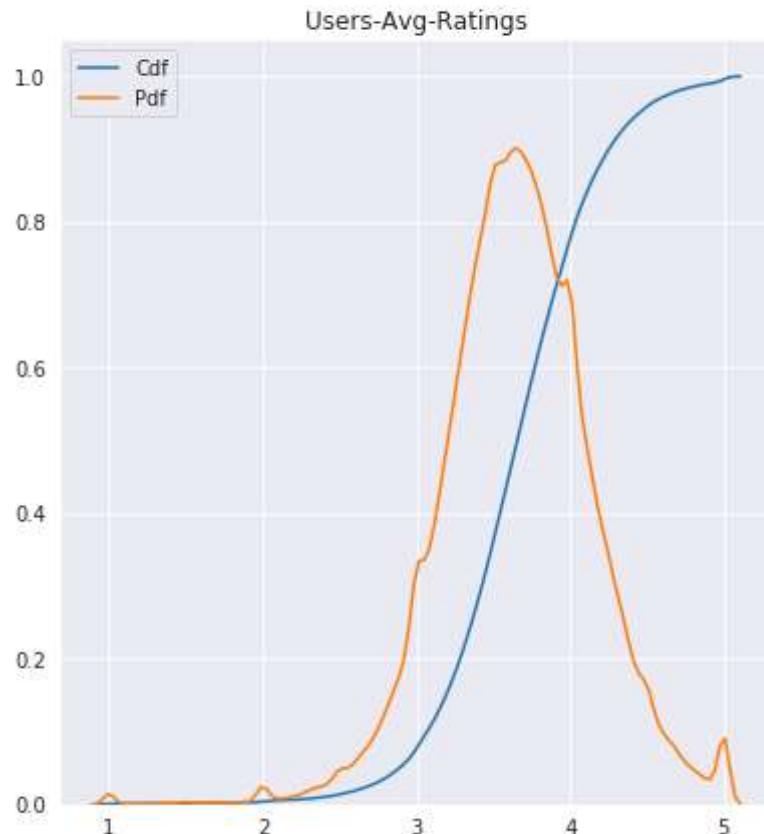
AVerage rating of movie 15 : 3.3038461538461537

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

In [40]:

```
1 %%time
2 # draw pdfs for average rating per user and average
3 fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))
4 fig.suptitle('Avg Ratings per User and per Movie', fontsize=15)
5
6 ax1.set_title('Users-Avg-Ratings')
7 # get the list of average user ratings from the averages dictionary..
8 user_averages = [rat for rat in train_averages['user'].values()]
9 sns.distplot(user_averages, ax=ax1, hist=False,
10             kde_kws=dict(cumulative=True), label='Cdf')
11 sns.distplot(user_averages, ax=ax1, hist=False, label='Pdf')
12
13 ax2.set_title('Movies-Avg-Rating')
14 # get the list of movie_average_ratings from the dictionary..
15 movie_averages = [rat for rat in train_averages['movie'].values()]
16 sns.distplot(movie_averages, ax=ax2, hist=False,
17             kde_kws=dict(cumulative=True), label='Cdf')
18 sns.distplot(movie_averages, ax=ax2, hist=False, label='Pdf')
19
20 plt.show()
```

### Avg Ratings per User and per Movie



CPU times: user 34.3 s, sys: 376 ms, total: 34.7 s  
Wall time: 34.2 s

#### 3.3.8 Cold Start problem

##### 3.3.8.1 Cold Start problem with Users

In [41]:

```
1 total_users = len(np.unique(df.user))
2 users_train = len(train_averages['user'])
3 new_users = total_users - users_train
4
5 print('\nTotal number of Users :', total_users)
6 print('\nNumber of Users in Train data :', users_train)
7 print("\nNo of Users that didn't appear in train data: {}({} %) \n ".format(new_users,
8                                         np.round((new_users/total_users)*100, 2)))
```

Total number of Users : 480189

Number of Users in Train data : 405041

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

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

### 3.3.8.2 Cold Start problem with Movies

In [42]:

```
1 total_movies = len(np.unique(df.movie))
2 movies_train = len(train_averages['movie'])
3 new_movies = total_movies - movies_train
4
5 print('\nTotal number of Movies :', total_movies)
6 print('\nNumber of Users in Train data :', movies_train)
7 print("\nNo of Movies that didn't appear in train data: {}({} %) \n ".format(new_movies,
8 np.round((new_movies/total_movies)*100, 2)))
```

Total number of Movies : 17770

Number of Users in Train data : 17424

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

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

## 3.4 Computing Similarity matrices

### 3.4.1 Computing User-User Similarity matrix

1. Calculating User User Similarity\_Matrix is **not very easy**(unless you have huge Computing Power and lots of time) because of number of users being large.
  - You can try if you want to. Your system could crash or the program stops with **Memory Error**

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

In [43]:

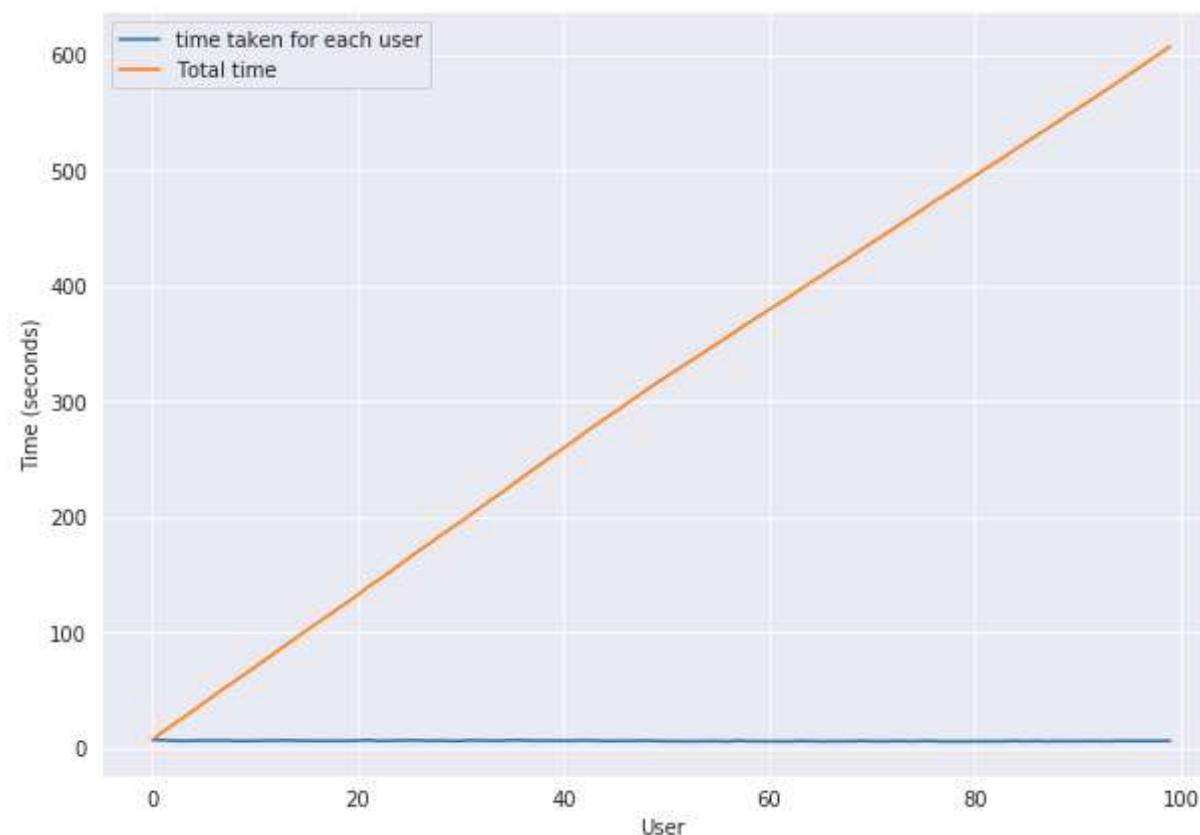
```
1 from sklearn.metrics.pairwise import cosine_similarity
2
3
4 def compute_user_similarity(sparse_matrix, compute_for_few=False, top = 100, verbose=False, verb_for_n_rows = 20,
5                             draw_time_taken=True):
6     no_of_users, _ = sparse_matrix.shape
7     # get the indices of non zero rows(users) from our sparse matrix
8     row_ind, col_ind = sparse_matrix.nonzero()
9     row_ind = sorted(set(row_ind)) # we don't have to
10    time_taken = list() # time taken for finding similar users for an user..
11
12    # we create rows, cols, and data lists.., which can be used to create sparse matrices
13    rows, cols, data = list(), list(), list()
14    if verbose: print("Computing top",top,"similarities for each user..")
15
16    start = datetime.now()
17    temp = 0
18
19    for row in row_ind[:top] if compute_for_few else row_ind:
20        temp = temp+1
21        prev = datetime.now()
22
23        # get the similarity row for this user with all other users
24        sim = cosine_similarity(sparse_matrix.getrow(row), sparse_matrix).ravel()
25        # We will get only the top ''top'' most similar users and ignore rest of them..
26        top_sim_ind = sim.argsort()[-top:]
27        top_sim_val = sim[top_sim_ind]
28
29        # add them to our rows, cols and data
30        rows.extend([row]*top)
31        cols.extend(top_sim_ind)
32        data.extend(top_sim_val)
33        time_taken.append(datetime.now().timestamp() - prev.timestamp())
34        if verbose:
35            if temp%verb_for_n_rows == 0:
36                print("computing done for {} users [ time elapsed : {} ]"
37                     .format(temp, datetime.now()-start))
38
39
40    # lets create sparse matrix out of these and return it
41    if verbose: print('Creating Sparse matrix from the computed similarities')
```

```
42 #return rows, cols, data
43
44 if draw_time_taken:
45     plt.plot(time_taken, label = 'time taken for each user')
46     plt.plot(np.cumsum(time_taken), label='Total time')
47     plt.legend(loc='best')
48     plt.xlabel('User')
49     plt.ylabel('Time (seconds)')
50     plt.show()
51
52 return sparse.csr_matrix((data, (rows, cols)), shape=(no_of_users, no_of_users)), time_taken
```

In [44]:

```
1 start = datetime.now()
2 u_u_sim_sparse, _ = compute_user_similarity(train_sparse_matrix, compute_for_few=True, top = 100,
3                                              verbose=True)
4 print("-"*100)
5 print("Time taken :",datetime.now()-start)
```

Computing top 100 similarities for each user..  
computing done for 20 users [ time elapsed : 0:02:06.209501 ]  
computing done for 40 users [ time elapsed : 0:04:13.294686 ]  
computing done for 60 users [ time elapsed : 0:06:13.808011 ]  
computing done for 80 users [ time elapsed : 0:08:09.528931 ]  
computing done for 100 users [ time elapsed : 0:10:07.287556 ]  
Creating Sparse matrix from the computed similarities



Time taken : 0:10:19.667435

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

- We have **405,041 users** in our training set and computing similarities between them..( **17K dimensional vector..**) is time consuming..
- From above plot, It took roughly **8.88 sec** for computing similar users for **one user**
- We have **405,041 users** with us in training set.
- $405041 \times 8.88 = 3596764.08 \text{ sec} = 59946.068 \text{ min} = 999.101133333 \text{ hours} = 41.629213889 \text{ days}\dots$ 
  - Even if we run on 4 cores parallelly (a typical system now a days), It will still take almost **10 and 1/2 days**.

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

```
In [0]: 1 from datetime import datetime
2 from sklearn.decomposition import TruncatedSVD
3
4 start = datetime.now()
5
6 # initialize the algorithm with some parameters..
7 # All of them are default except n_components. n_itr is for Randomized SVD solver.
8 netflix_svd = TruncatedSVD(n_components=500, algorithm='randomized', random_state=15)
9 trunc_svd = netflix_svd.fit_transform(train_sparse_matrix)
10
11 print(datetime.now()-start)
```

0:29:07.069783

Here,

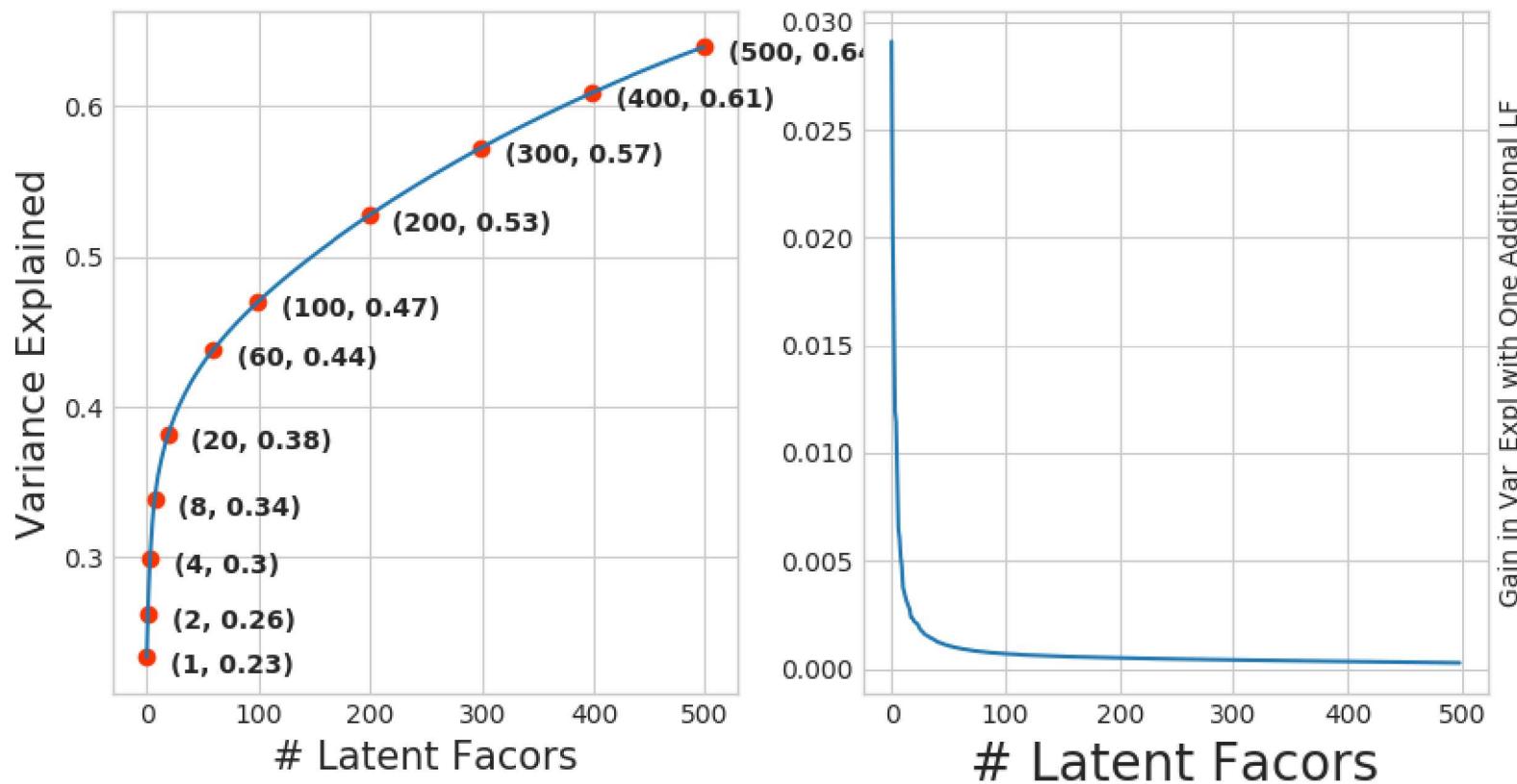
- $\Sigma \leftarrow (\text{netflix\_svd.singular\_values\_})$
- $V^T \leftarrow (\text{netflix\_svd.components\_})$
- $U$  is not returned. instead **Projection\_of\_X** onto the new vectorspace is returned.

- It uses **randomized svd** internally, which returns **All 3 of them saperately**. Use that instead..

```
In [0]: 1 expl_var = np.cumsum(netflix_svd.explained_variance_ratio_)
```

```
In [0]:  
1 fig, (ax1, ax2) = plt.subplots(nrows=1, ncols=2, figsize=plt.figaspect(.5))  
2  
3 ax1.set_ylabel("Variance Explained", fontsize=15)  
4 ax1.set_xlabel("# Latent Factors", fontsize=15)  
5 ax1.plot(expl_var)  
6 # annotate some (Latent factors, expl_var) to make it clear  
7 ind = [1, 2, 4, 8, 20, 60, 100, 200, 300, 400, 500]  
8 ax1.scatter(x = [i-1 for i in ind], y = expl_var[[i-1 for i in ind]], c='#ff3300')  
9 for i in ind:  
10     ax1.annotate(s = "({}, {})".format(i, np.round(expl_var[i-1], 2)), xy=(i-1, expl_var[i-1]),  
11                  xytext = (i+20, expl_var[i-1] - 0.01), fontweight='bold')  
12  
13 change_in_expl_var = [expl_var[i+1] - expl_var[i] for i in range(len(expl_var)-1)]  
14 ax2.plot(change_in_expl_var)  
15  
16  
17  
18 ax2.set_ylabel("Gain in Var_Expl with One Additional LF", fontsize=10)  
19 ax2.yaxis.set_label_position("right")  
20 ax2.set_xlabel("# Latent Factors", fontsize=10)  
21  
22 plt.show()
```

<IPython.core.display.Javascript object>



In [0]:

```
1 for i in ind:  
2     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 factor to it, the **\_gain in explained variance** with that addition is decreasing. (Obviously, because they are sorted that way).
- **LHS Graph:**
  - **x** --- ( No of latent factors ),
  - **y** --- ( The variance explained by taking x latent factors)
- **More decrease in the line (RHS graph) :**
  - We are getting more explained variance than before.
- **Less decrease in that line (RHS graph) :**
  - We are not getting benefitted from adding latent factor further. This is what is shown in the plots.
- **RHS Graph:**
  - **x** --- ( No of latent factors ),

- **y** --- ( Gain n Expl\_Var by taking one additional latent factor)

```
In [0]: 1 # Let's project our Original U_M matrix into into 500 Dimensional space...
2 start = datetime.now()
3 trunc_matrix = train_sparse_matrix.dot(netflix_svd.components_.T)
4 print(datetime.now() - start)
```

0:00:45.670265

```
In [0]: 1 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 [45]: 1 if not os.path.isfile('trunc_sparse_matrix.npz'):
2     # create that sparse sparse matrix
3     trunc_sparse_matrix = sparse.csr_matrix(trunc_matrix)
4     # Save this truncated sparse matrix for later usage..
5     sparse.save_npz('trunc_sparse_matrix', trunc_sparse_matrix)
6 else:
7     trunc_sparse_matrix = sparse.load_npz('trunc_sparse_matrix.npz')
```

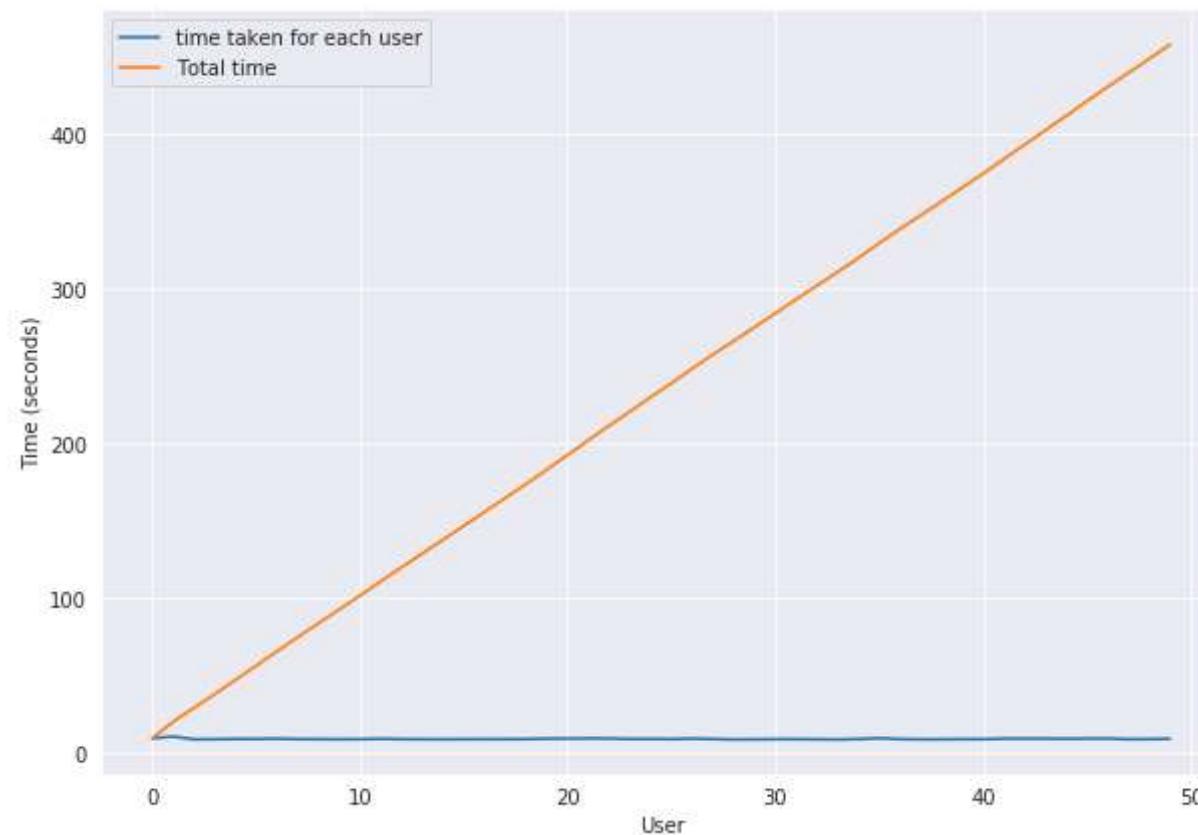
```
In [46]: 1 trunc_sparse_matrix.shape
```

Out[46]: (2649430, 500)

In [47]:

```
1 start = datetime.now()
2 trunc_u_u_sim_matrix, _ = compute_user_similarity(trunc_sparse_matrix, compute_for_few=True, top=50, verbose=True,
3                                                 verb_for_n_rows=10)
4 print("-"*50)
5 print("time:",datetime.now()-start)
```

Computing top 50 similarities for each user..  
computing done for 10 users [ time elapsed : 0:01:32.852714 ]  
computing done for 20 users [ time elapsed : 0:03:03.422646 ]  
computing done for 30 users [ time elapsed : 0:04:35.391511 ]  
computing done for 40 users [ time elapsed : 0:06:05.672636 ]  
computing done for 50 users [ time elapsed : 0:07:38.178973 ]  
Creating Sparse matrix from the computed similarities



time: 0:08:08.646406

: 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.
- $405041 \times 12.18 = 4933399.38$  sec == 82223.323 min == 1370.388716667 hours == 57.099529861 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 [13]:

```
1 %%time
2 if not os.path.isfile('m_m_sim_sparse.npz'):
3     print("It seems you don't have that file. Computing movie_movie similarity...")
4     start = datetime.now()
5     m_m_sim_sparse = cosine_similarity(X=train_sparse_matrix.T, dense_output=False)
6     print("Done..")
7     # store this sparse matrix in disk before using it. For future purposes.
8     print("Saving it to disk without the need of re-computing it again.. ")
9     sparse.save_npz("m_m_sim_sparse.npz", m_m_sim_sparse)
10    print("Done..")
11 else:
12     print("It is there, We will get it.")
13     m_m_sim_sparse = sparse.load_npz("m_m_sim_sparse.npz")
14     print("Done ...")
15
16 print("It's a ",m_m_sim_sparse.shape," dimensional matrix")
```

It is there, We will get it.

Done ...

It's a (17771, 17771) dimensional matrix

CPU times: user 24.1 s, sys: 3.68 s, total: 27.8 s

Wall time: 29.1 s

In [14]:

```
1 m_m_sim_sparse.shape
```

Out[14]: (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 [9]:

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

In [10]:

```
1 start = datetime.now()
2 similar_movies = dict()
3 for movie in movie_ids:
4     # get the top similar movies and store them in the dictionary
5     sim_movies = m_m_sim_sparse[movie].toarray().ravel().argsort()[:-1][1:]
6     similar_movies[movie] = sim_movies[:100]
7 print(datetime.now() - start)
8
9 # just testing similar movies for movie_15
10 similar_movies[15]
```

0:00:30.711678

Out[10]: 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 [61]: 1 # First Let's load the movie details into soe dataframe..
2 # movie details are in 'netflix/movie_titles.csv'
3
4 movie_titles = pd.read_csv("movie_titles.csv", sep=',', header = None,
5                             names=['movie_id', 'year_of_release', 'title'], verbose=True,
6                             index_col = 'movie_id', encoding = "ISO-8859-1")
7
8 movie_titles.head()
```

Tokenization took: 4.07 ms  
Type conversion took: 11.11 ms  
Parser memory cleanup took: 0.01 ms

Out[61]:

	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 [64]:

```
1 mv_id = 80
2
3 print("\nMovie ---->", movie_titles.loc[mv_id].values[1])
4
5 print("\nIt has {} Ratings from users.".format(train_sparse_matrix[:,mv_id].getnnz()))
6
7 print("\nWe have {} movies which are similarto this and we will get only top most..".format(m_m_sim_sparse[:,mv_id]
```

Movie ----> Winter Kills

It has 243 Ratings from users.

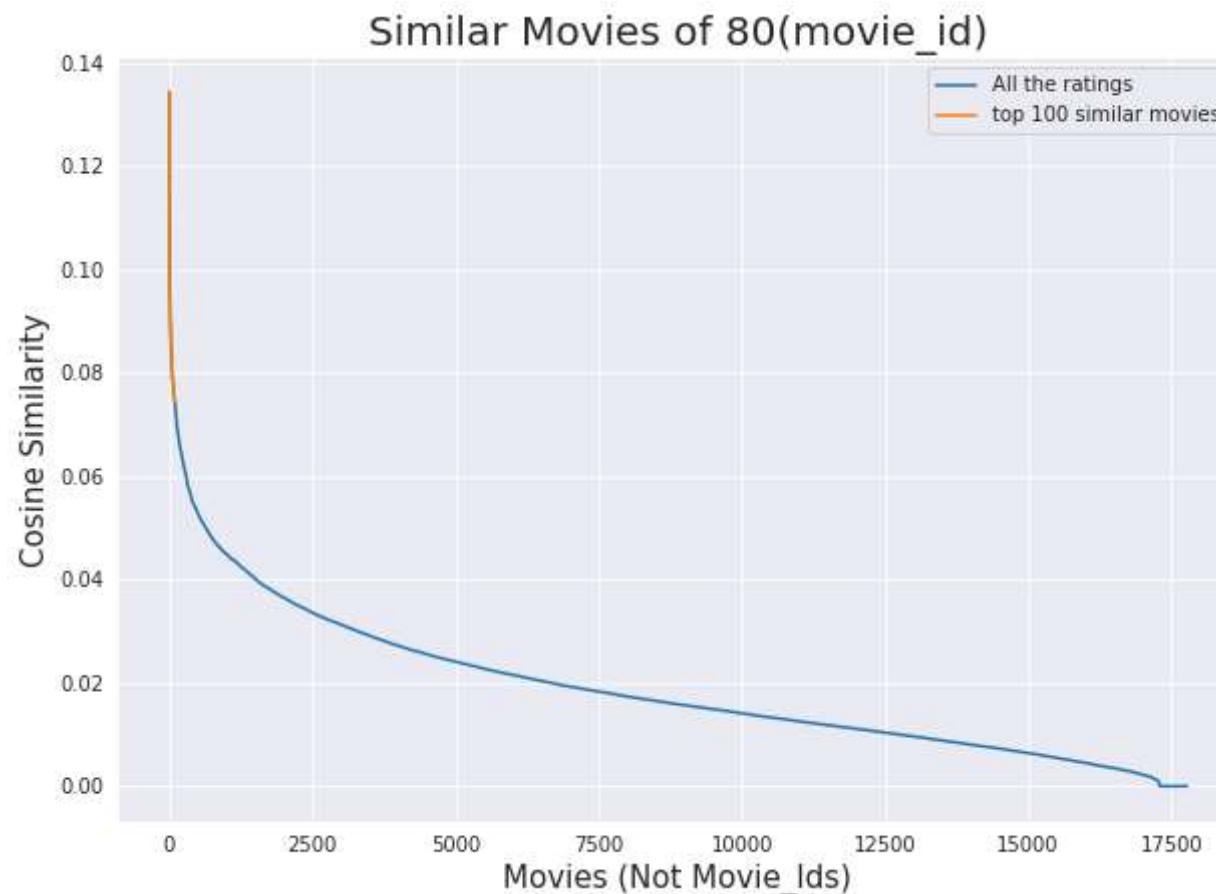
We have 17292 movies which are similarto this and we will get only top most..

In [65]:

```
1 similarities = m_m_sim_sparse[mv_id].toarray().ravel()
2
3 similar_indices = similarities.argsort()[:-1][1:]
4
5 similarities[similar_indices]
6
7 sim_indices = similarities.argsort()[:-1][1:] # It will sort and reverse the array and ignore its similarity (ie., 1
8 # and return its indices(movie_ids)
```

In [66]:

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



### Top 10 similar movies

```
In [67]: 1 movie_titles.loc[sim_indices[:10]]
```

Out[67]:

movie_id	year_of_release	title
13751	1981.0	Cutter's Way
6699	1972.0	Fat City
5263	1969.0	Medium Cool
15963	1966.0	Seconds
1354	1968.0	Targets
12304	1974.0	The Parallax View
11449	1955.0	The Big Knife
6141	1984.0	Flashpoint
7679	1968.0	Lady in Cement
2404	1991.0	Picture This

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

## 4. Machine Learning Models



In [15]:

```
1 def get_sample_sparse_matrix(sparse_matrix, no_users, no_movies, path, verbose = True):
2     """
3         It will get it from the ''path'' if it is present or It will create
4         and store the sampled sparse matrix in the path specified.
5     """
6
7     # get (row, col) and (rating) tuple from sparse_matrix...
8     row_ind, col_ind, ratings = sparse.find(sparse_matrix)
9     users = np.unique(row_ind)
10    movies = np.unique(col_ind)
11
12    print("Original Matrix : (users, movies) -- ({}, {})".format(len(users), len(movies)))
13    print("Original Matrix : Ratings -- {}".format(len(ratings)))
14
15    # It just to make sure to get same sample everytime we run this program..
16    # and pick without replacement....
17    np.random.seed(15)
18    sample_users = np.random.choice(users, no_users, replace=False)
19    sample_movies = np.random.choice(movies, no_movies, replace=False)
20    # get the boolean mask or these sampled_items in originl row/col_inds..
21    mask = np.logical_and( np.isin(row_ind, sample_users),
22                           np.isin(col_ind, sample_movies) )
23
24    sample_sparse_matrix = sparse.csr_matrix((ratings[mask], (row_ind[mask], col_ind[mask])),
25                                              shape=(max(sample_users)+1, max(sample_movies)+1))
26
27    if verbose:
28        print("Sampled Matrix : (users, movies) -- ({}, {})".format(len(sample_users), len(sample_movies)))
29        print("Sampled Matrix : Ratings -- ", format(ratings[mask].shape[0]))
30
31    print('Saving it into disk for furthur usage..')
32    # save it into disk
33    sparse.save_npz(path, sample_sparse_matrix)
34    if verbose:
35        print('Done..\n')
36
37    return sample_sparse_matrix
```

## 4.1 Sampling Data

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

In [16]:

```
1 %%time
2 path = "sample_train_sparse_matrix.npz"
3 if os.path.isfile(path):
4     print("It is present in your pwd, getting it from disk....")
5     # just get it from the disk instead of computing it
6     sample_train_sparse_matrix = sparse.load_npz(path)
7     print("DONE..")
8 else:
9     # get 10k users and 1k movies from available data
10    sample_train_sparse_matrix = get_sample_sparse_matrix(train_sparse_matrix, no_users=25000, no_movies=3000,
11                                path = path)
12
```

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

DONE..

CPU times: user 60 ms, sys: 0 ns, total: 60 ms

Wall time: 73.8 ms

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

In [17]:

```
1 start = datetime.now()
2
3 path = "sample_test_sparse_matrix.npz"
4 if os.path.isfile(path):
5     print("It is present in your pwd, getting it from disk....")
6     # just get it from the disk instead of computing it
7     sample_test_sparse_matrix = sparse.load_npz(path)
8     print("DONE..")
9 else:
10    # get 5k users and 500 movies from available data
11    sample_test_sparse_matrix = get_sample_sparse_matrix(test_sparse_matrix, no_users=5000, no_movies=500,
12                                                       path = "sample_test_sparse_matrix.npz")
13 print(datetime.now() - start)
```

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

Original Matrix : Ratings -- 20096102

Sampled Matrix : (users, movies) -- (5000 500)

Sampled Matrix : Ratings -- 7333

Saving it into disk for furthur usage..

Done..

0:00:18.504482

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

In [18]:

```
1 sample_train_averages = dict()
```

### 4.2.1 Finding Global Average of all movie ratings

```
In [19]: 1 # get the global average of ratings in our train set.  
2 global_average = sample_train_sparse_matrix.sum()/sample_train_sparse_matrix.count_nonzero()  
3 sample_train_averages['global'] = global_average  
4 sample_train_averages
```

```
Out[19]: {'global': 3.5875813607223455}
```

## 4.2.2 Finding Average rating per User

```
In [20]: 1 sample_train_averages['user'] = get_average_ratings(sample_train_sparse_matrix, of_users=True)  
2 print('\nAverage rating of user 1515220 :',sample_train_averages['user'][1515220])
```

```
Average rating of user 1515220 : 3.923076923076923
```

## 4.2.3 Finding Average rating per Movie

```
In [21]: 1 sample_train_averages['movie'] = get_average_ratings(sample_train_sparse_matrix, of_users=False)  
2 print('\n Average rating of movie 15153 :',sample_train_averages['movie'][15153])
```

```
AVerage rating of movie 15153 : 2.752
```

## 4.3 Featurizing data

```
In [22]: 1 print('\n No of ratings in Our Sampled train matrix is : {}'.format(sample_train_sparse_matrix.count_nonzero()))  
2 print('\n No of ratings in Our Sampled test matrix is : {}'.format(sample_test_sparse_matrix.count_nonzero()))
```

```
No of ratings in Our Sampled train matrix is : 856986
```

```
No of ratings in Our Sampled test matrix is : 7333
```

## 4.3.1 Featurizing data for regression problem

### 4.3.1.1 Featurizing train data

```
In [23]: 1 # get users, movies and ratings from our samples train sparse matrix  
2 sample_train_users, sample_train_movies, sample_train_ratings = sparse.find(sample_train_sparse_matrix)
```

In [6]:

```
1 #####  
2 # It took me almost 10 hours to prepare this train dataset.#  
3 #####  
4 start = datetime.now()  
5 if os.path.isfile('reg_train.csv'):  
6     print("File already exists you don't have to prepare again...")  
7 else:  
8     print('preparing {} tuples for the dataset..\n'.format(len(sample_train_ratings)))  
9     with open('reg_train.csv', mode='w') as reg_data_file:# if not exist bcz of mode w a empty .csv file created  
10        count = 0  
11        for (user, movie, rating) in zip(sample_train_users, sample_train_movies, sample_train_ratings):  
12            st = datetime.now()  
13            # print(user, movie)  
14            #----- Ratings of "movie" by similar users of "user" -----  
15            # compute the similar Users of the "user"  
16            user_sim = cosine_similarity(sample_train_sparse_matrix[user], sample_train_sparse_matrix).ravel()  
17            top_sim_users = user_sim.argsort()[:-1][1:] # we are ignoring 'The User' from its similar users.  
18            # get the ratings of most similar users for this movie  
19            top_ratings = sample_train_sparse_matrix[top_sim_users, movie].toarray().ravel()  
20            # we will make it's length "5" by adding movie averages to .  
21            top_sim_users_ratings = list(top_ratings[top_ratings != 0][:5])  
22            top_sim_users_ratings.extend([sample_train_averages['movie'][movie]]*(5 - len(top_sim_users_ratings)))  
23            # print(top_sim_users_ratings, end=" ")  
24  
25  
26            #----- Ratings by "user" to similar movies of "movie" -----  
27            # compute the similar movies of the "movie"  
28            movie_sim = cosine_similarity(sample_train_sparse_matrix[:, movie].T, sample_train_sparse_matrix.T).ravel()  
29            top_sim_movies = movie_sim.argsort()[:-1][1:] # we are ignoring 'The User' from its similar users.  
30            # get the ratings of most similar movie rated by this user..  
31            top_ratings = sample_train_sparse_matrix[user, top_sim_movies].toarray().ravel()  
32            # we will make it's length "5" by adding user averages to.  
33            top_sim_movies_ratings = list(top_ratings[top_ratings != 0][:5])  
34            top_sim_movies_ratings.extend([sample_train_averages['user'][user]]*(5 - len(top_sim_movies_ratings)))  
35            # print(top_sim_movies_ratings, end=" : -- ")  
36  
37            #-----prepare the row to be stores in a file-----#  
38            row = list()  
39            row.append(user)  
40            row.append(movie)  
41            # Now add the other features to this data...
```

```
42     row.append(sample_train_averages['global']) # first feature
43     # next 5 features are similar_users "movie" ratings
44     row.extend(top_sim_users_ratings)
45     # next 5 features are "user" ratings for similar_movies
46     row.extend(top_sim_movies_ratings)
47     # Avg_user rating
48     row.append(sample_train_averages['user'][user])
49     # Avg_movie rating
50     row.append(sample_train_averages['movie'][movie])
51
52     # finalley, The actual Rating of this user-movie pair...
53     row.append(rating)
54     count = count + 1
55
56     # add rows to the file opened..
57     reg_data_file.write(', '.join(map(str, row)))
58     reg_data_file.write('\n')
59     if (count)%25000 == 0:
60         # print(', '.join(map(str, row)))
61         print("Done for {} rows---- {} ".format(count, datetime.now() - start))
62
63
64 print(datetime.now() - start)
```

File already exists you don't have to prepare again...

0:00:00.000580

### Reading from the file to make a Train\_dataframe

In [7]:

```
1 reg_train = pd.read_csv('reg_train.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5', 'smr1', 'smr2', 'smr3', 'smr4', 'smr5', 'UAvg', 'MAvg', 'rating']
2 reg_train.head()
```

Out[7]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg	rating
0	174683	10	3.587581	5.0	5.0	3.0	4.0	4.0	3.0	5.0	4.0	3.0	2.0	3.882353	3.611111	5
1	233949	10	3.587581	4.0	4.0	5.0	1.0	3.0	2.0	3.0	2.0	3.0	3.0	2.692308	3.611111	3
2	555770	10	3.587581	4.0	5.0	4.0	4.0	5.0	4.0	2.0	5.0	4.0	4.0	3.795455	3.611111	4
3	767518	10	3.587581	2.0	5.0	4.0	4.0	3.0	5.0	5.0	4.0	4.0	3.0	3.884615	3.611111	5
4	894393	10	3.587581	3.0	5.0	4.0	4.0	3.0	4.0	4.0	4.0	4.0	4.0	4.000000	3.611111	4

- **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 [28]:

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

```
In [29]: 1 sample_train_averages['global']
```

```
Out[29]: 3.5875813607223455
```

In [2]:

```
1 start = datetime.now()
2
3 if os.path.isfile('reg_test.csv'):
4     print("It is already created...")
5 else:
6
7     print('preparing {} tuples for the dataset..\n'.format(len(sample_test_ratings)))
8     with open('reg_test.csv', mode='w') as reg_data_file:
9         count = 0
10        for (user, movie, rating) in zip(sample_test_users, sample_test_movies, sample_test_ratings):
11            st = datetime.now()
12
13            #----- Ratings of "movie" by similar users of "user" -----
14            #print(user, movie)
15            try:
16                # compute the similar Users of the "user"
17                user_sim = cosine_similarity(sample_train_sparse_matrix[user], sample_train_sparse_matrix).ravel()
18                top_sim_users = user_sim.argsort()[:-1][1:] # we are ignoring 'The User' from its similar users.
19                # get the ratings of most similar users for this movie
20                top_ratings = sample_train_sparse_matrix[top_sim_users, movie].toarray().ravel()
21                # we will make it's length "5" by adding movie averages to .
22                top_sim_users_ratings = list(top_ratings[top_ratings != 0][:5])
23                top_sim_users_ratings.extend([sample_train_averages['movie'][movie]]*(5 - len(top_sim_users_ratings)))
24                # print(top_sim_users_ratings, end="--")
25
26            except (IndexError, KeyError):
27                # It is a new User or new Movie or there are no ratings for given user for top similar movies...
28                ##### Cold Start Problem #####
29                top_sim_users_ratings.extend([sample_train_averages['global']]*(5 - len(top_sim_users_ratings)))
30                #print(top_sim_users_ratings)
31            except:
32                print(user, movie)
33                # we just want KeyErrors to be resolved. Not every Exception...
34                raise
35
36
37
38            #----- Ratings by "user" to similar movies of "movie" -----
39            try:
40                # compute the similar movies of the "movie"
41                movie_sim = cosine_similarity(sample_train_sparse_matrix[:, movie].T, sample_train_sparse_matrix.T).r
```

```
42     top_sim_movies = movie_sim.argsort()[:-1][1:] # we are ignoring 'The User' from its similar users.
43     # get the ratings of most similar movie rated by this user..
44     top_ratings = sample_train_sparse_matrix[user, top_sim_movies].toarray().ravel()
45     # we will make it's length "5" by adding user averages to.
46     top_sim_movies_ratings = list(top_ratings[top_ratings != 0][:5])
47     top_sim_movies_ratings.extend([sample_train_averages['user'][user]]*(5-len(top_sim_movies_ratings)))
48     #print(top_sim_movies_ratings)
49 except (IndexError, KeyError):
50     #print(top_sim_movies_ratings, end=" : -- ")
51     top_sim_movies_ratings.extend([sample_train_averages['global']]*(5-len(top_sim_movies_ratings)))
52     #print(top_sim_movies_ratings)
53 except :
54     raise
55
56     #-----prepare the row to be stores in a file-----
57 row = list()
58 # add usser and movie name first
59 row.append(user)
60 row.append(movie)
61 row.append(sample_train_averages['global']) # first feature
62 #print(row)
63 # next 5 features are similar_users "movie" ratings
64 row.extend(top_sim_users_ratings)
65 #print(row)
66 # next 5 features are "user" ratings for similar_movies
67 row.extend(top_sim_movies_ratings)
68 #print(row)
69 # Avg_user rating
70 try:
71     row.append(sample_train_averages['user'][user])
72 except KeyError:
73     row.append(sample_train_averages['global'])
74 except:
75     raise
76 #print(row)
77 # Avg_movie rating
78 try:
79     row.append(sample_train_averages['movie'][movie])
80 except KeyError:
81     row.append(sample_train_averages['global'])
82 except:
83     raise
```

```

84     #print(row)
85     # finalley, The actual Rating of this user-movie pair...
86     row.append(rating)
87     #print(row)
88     count = count + 1
89
90     # add rows to the file opened..
91     reg_data_file.write(', '.join(map(str, row)))
92     #print(', '.join(map(str, row)))
93     reg_data_file.write('\n')
94     if (count)%5000 == 0:
95         #print(', '.join(map(str, row)))
96         print("Done for {} rows---- {} ".format(count, datetime.now() - start))
97 print("",datetime.now() - start)

```

It is already created...

Reading from the file to make a test dataframe

In [3]:

```

1 reg_test_df = pd.read_csv('reg_test.csv', names = ['user', 'movie', 'GAvg', 'sur1', 'sur2', 'sur3', 'sur4', 'sur5',
2                                         'smr1', 'smr2', 'smr3', 'smr4', 'smr5',
3                                         'UAvg', 'MAvg', 'rating'], header=None)
4 reg_test_df.head(4)

```

Out[3]:

	user	movie	GAvg	sur1	sur2	sur3	sur4	sur5	smr1	smr2	smr3	smr4	smr5	UAvg	MAvg
0	808635	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
1	941866	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
2	1737912	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581
3	1849204	71	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581	3.587581

- **GAvg** : Average rating of all the ratings
- **Similar users rating of this movie:**
  - sur1, sur2, sur3, sur4, sur5 ( top 5 similiar 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.
-