Objectives:

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

Constraints:

1. Some form of interpretability.

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

Double-click (or enter) to edit

# Preprocessing

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

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

    row = list()
```

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

    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
```

Double-click (or enter) to edit

```
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..
```

df.head()

|  | 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 |

df.describe()['rating']

```
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
```

## Checking for NaN values

```
# 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
```

Double-click (or enter) to edit

## Removing Duplicates

```
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..
```

Double-click (or enter) to edit

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

```python
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
```

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

```python
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")
```

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

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

```
    Training data
    --------------------------------------------------

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

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

```
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
```

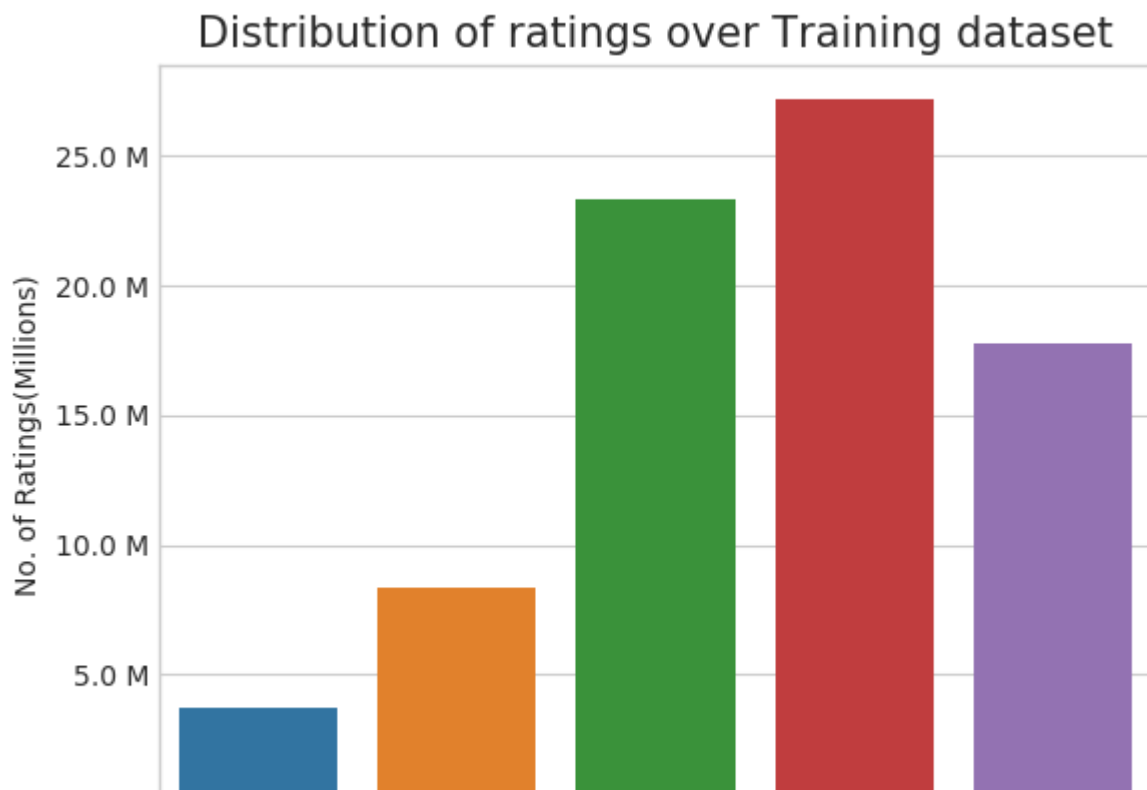# Exploratory Data Analysis on Train data

Double-click (or enter) to edit

```
# 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"
```

## Distribution of ratings

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

plt.show()
```

# Distribution of ratings over Training dataset



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

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

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

train_df.tail()
```

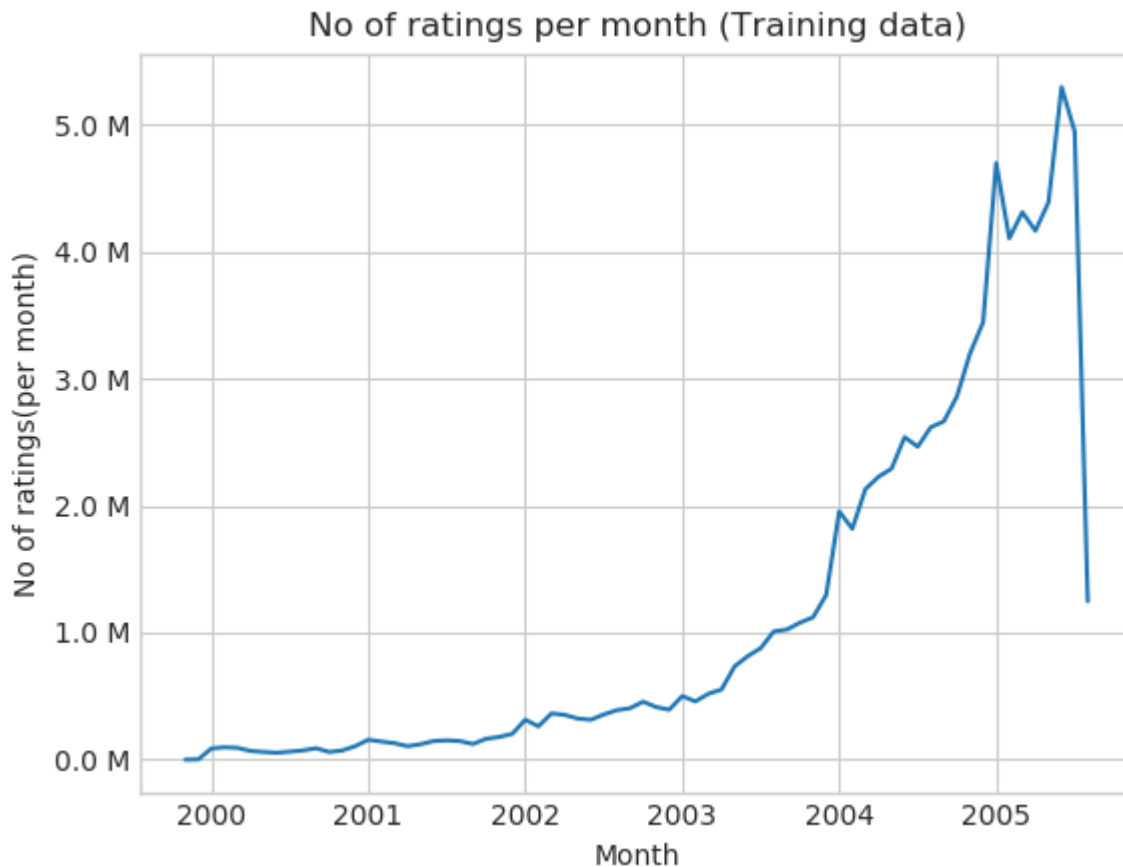|  | movie | user | rating | date | day_of_week |
|---|---|---|---|---|---|
| **80384400** | 12074 | 2033618 | 4 | 2005-08-08 | Monday |
| **80384401** | 862 | 1797061 | 3 | 2005-08-08 | Monday |
| **80384402** | 10986 | 1498715 | 5 | 2005-08-08 | Monday |
| **80384403** | 14861 | 500016 | 4 | 2005-08-08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

# Number of Ratings per a month

```
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)')
```

```
pit.yidvei( Nu or ratings(pei munti) )
ax.set_yticklabels([human(item, 'M') for item in ax.get_yticks()])
plt.show()
```



No of ratings per month (Training data)

Double-click (or enter) to edit

## Analysis on the Ratings given by user

```
no_of_rated_movies_per_user = train_df.groupby(by='user')['rating'].count().sort_values(ascen
```

```
no_of_rated_movies_per_user.head()
```

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

```
fig = plt.figure(figsize=plt.figaspect(.5))
```
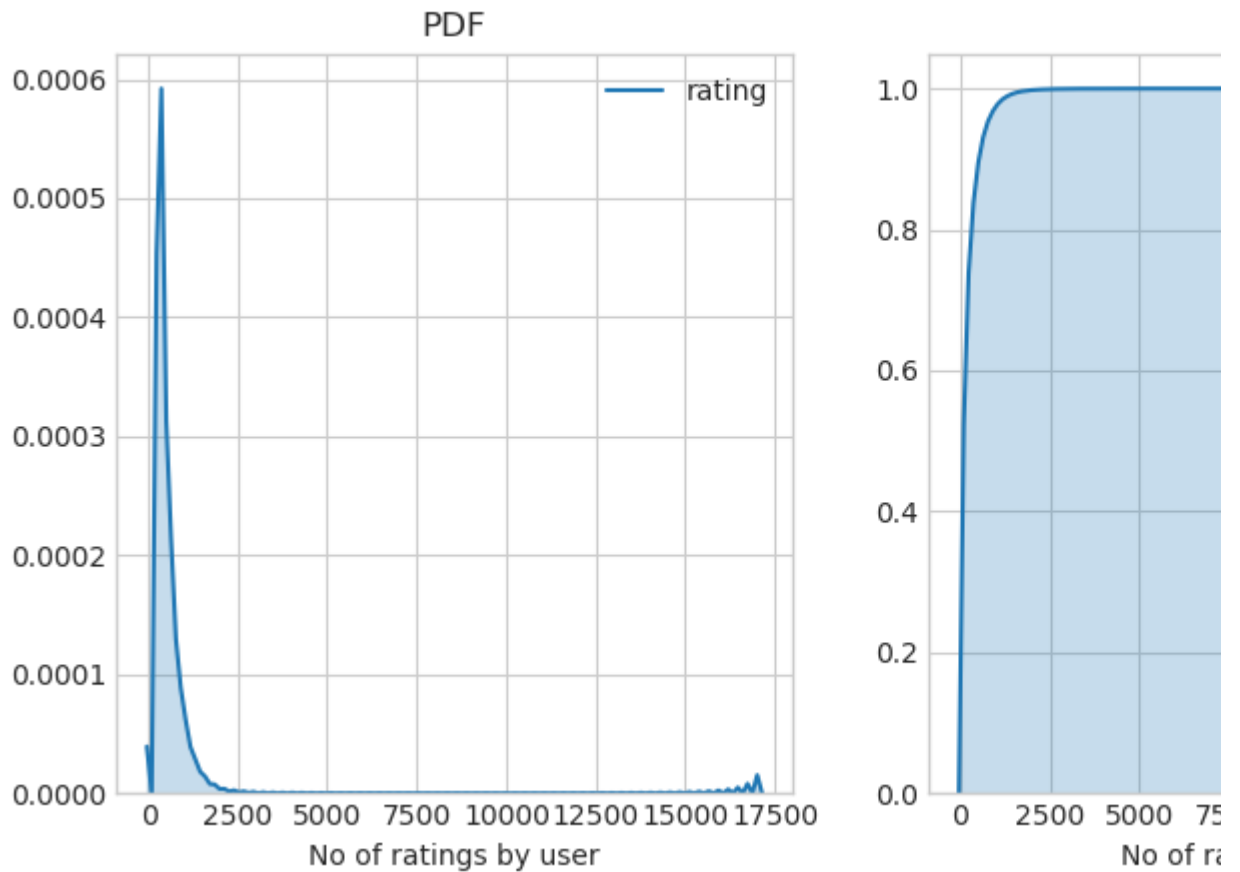
```
ax1 = plt.subplot(121)
```

```
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()
```



```
no_of_rated_movies_per_user.describe()
```

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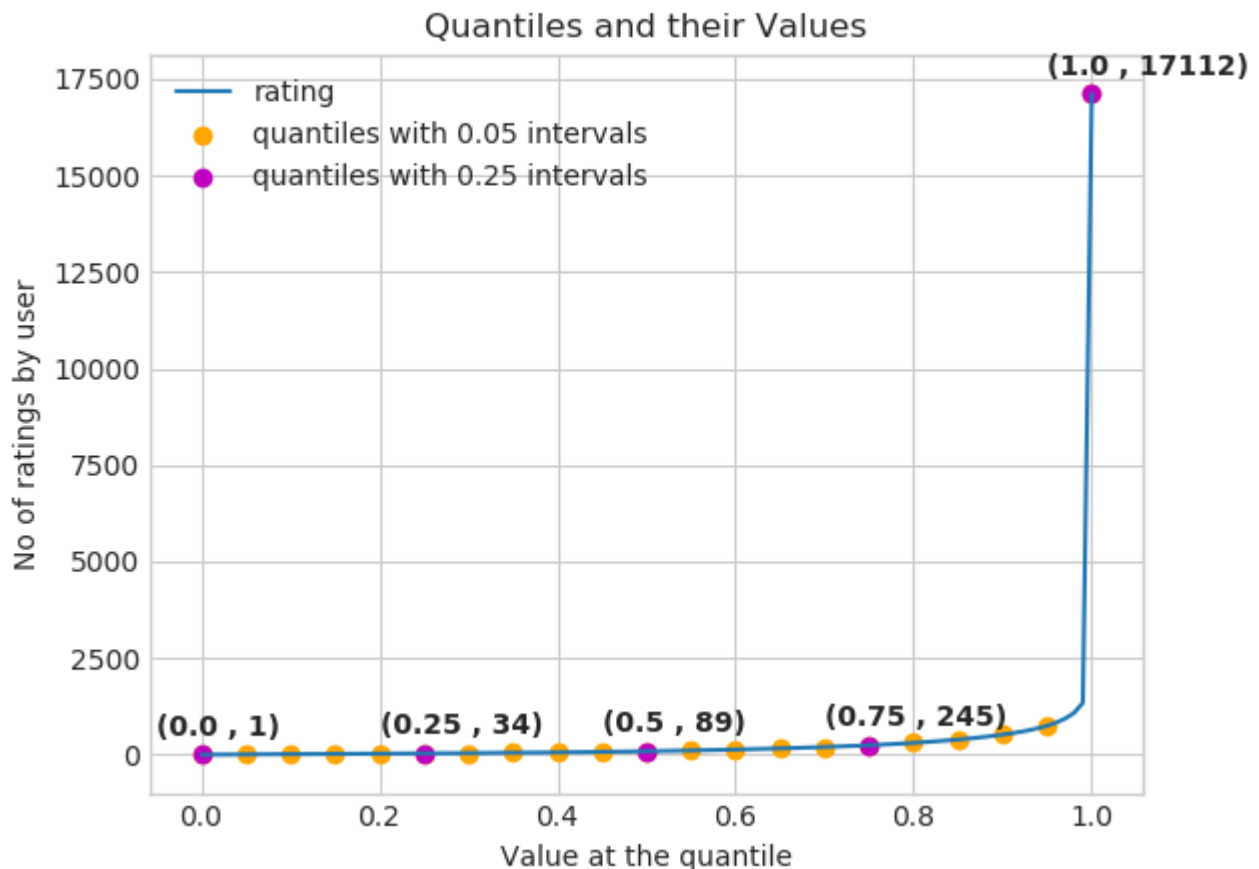
> *There, is something interesting going on with the quantiles..*

```
quantiles = no_of_rated_movies_per_user.quantile(np.arange(0,1.01,0.01), interpolation='highe
```

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

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


plt.show()
```



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

```
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**??

```
print('\n No of ratings at last 5 percentile : {}\n'.format(sum(no_of_rated_movies_per_user>=
```
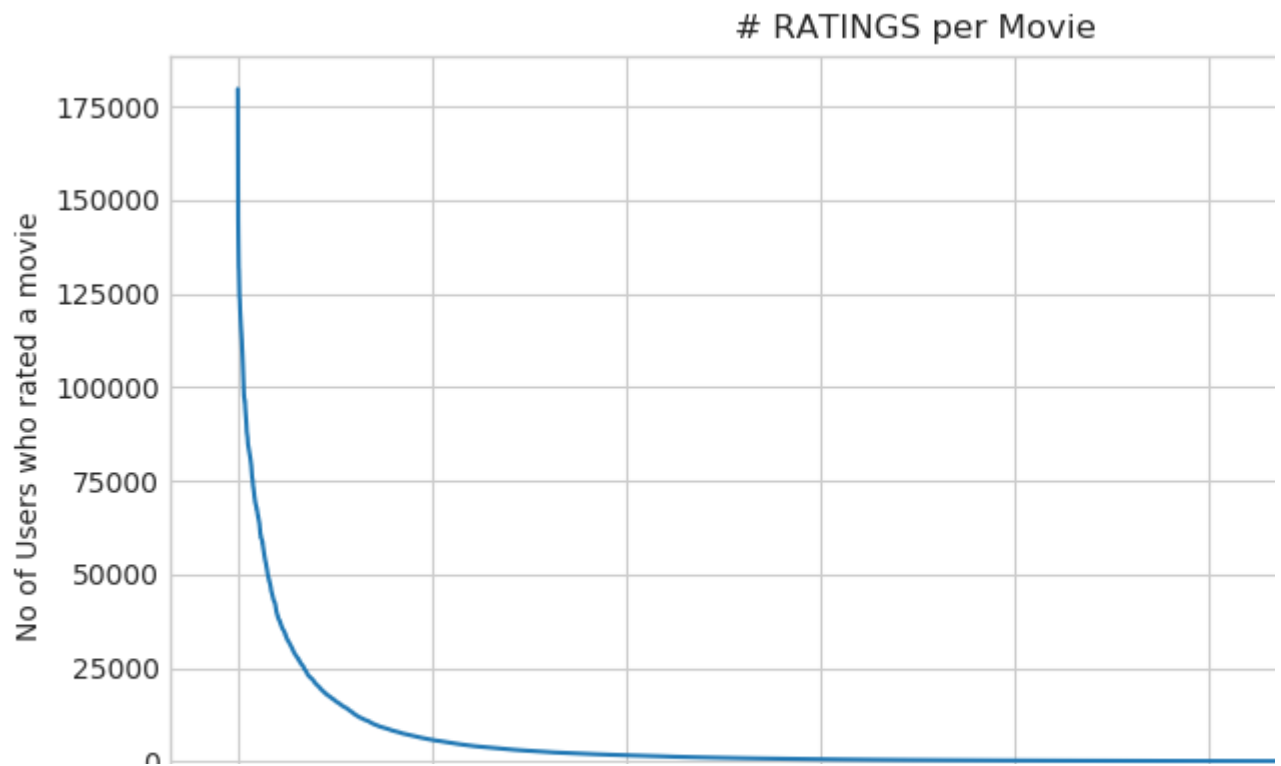
```
    No of ratings at last 5 percentile : 20305
```

# Analysis of ratings of a movie given by a user

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

fig = plt.figure(figsize=plt.figaspect(.5))
ax = plt.gca()
plt.plot(no_of_ratings_per_movie.values)
plt.title('# RATINGS per Movie')
plt.xlabel('Movie')
plt.ylabel('No of Users who rated a movie')
ax.set_xticklabels([])

plt.show()
```
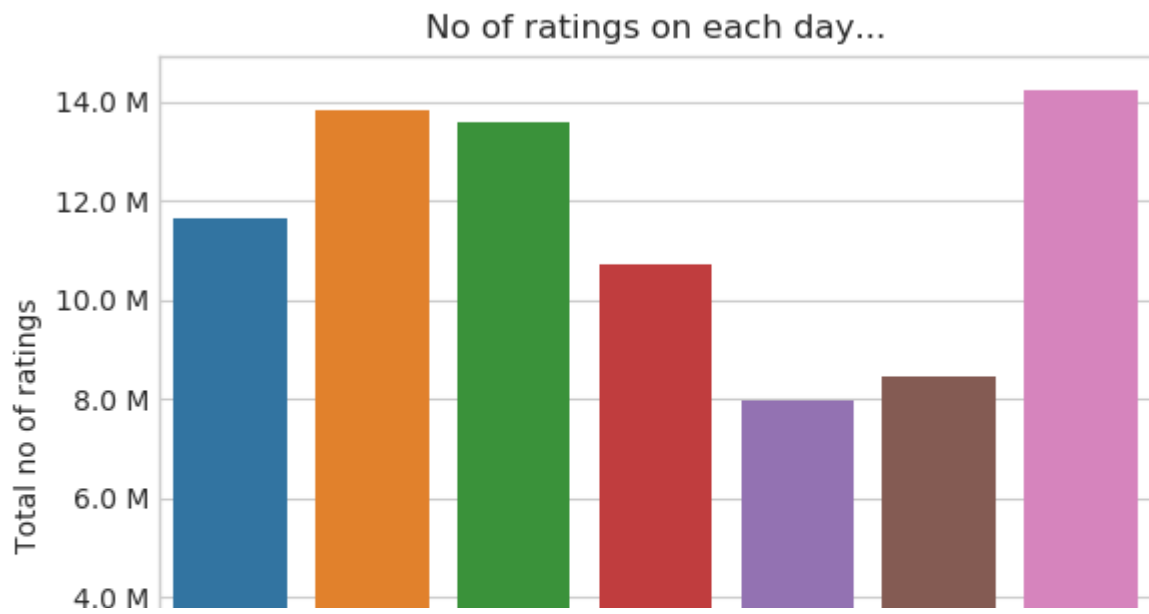
# RATINGS per Movie



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

    - There are some movies (which are very popular) which are rated by huge number of users.

    - But most of the movies(like 90%) got some hundereds of ratings.

# Number of ratings on each day of the week

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

No of ratings on each day...

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



0:01:10.003761

```
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
```

Double-click (or enter) to edit

# Creating sparse matrix from data frame

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

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

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

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

```
    Done..

    0:01:13.804969
```

## The Sparsity of Train Sparse Matrix

```
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 %
```

# Creating sparse matrix from test data frame

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

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

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

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

    0:00:18.566120
```

## The Sparsity of Test data Matrix

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

```
elem = test_sparse_matrix.count_nonzero()
```

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

> Sparsity Of Test matrix : 99.95731772988694 %

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

```python
# 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
    # creae a dictonary of users and their average ratigns..
    average_ratings = { i : sum_of_ratings[i]/no_of_ratings[i]
                                for i in range(u if of_users else m)
                                    if no_of_ratings[i] !=0}

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

## finding global average of all movie ratings

```python
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
```

> {'global': 3.582890686321557}

## finding average rating per user

```
train averages['user'] = get average ratings(train sparse matrix, of users=True)
```

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

       Average rating of user 10 : 3.3781094527363185

## finding average rating per movie

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

       AVerage rating of movie 15 : 3.3038461538461537
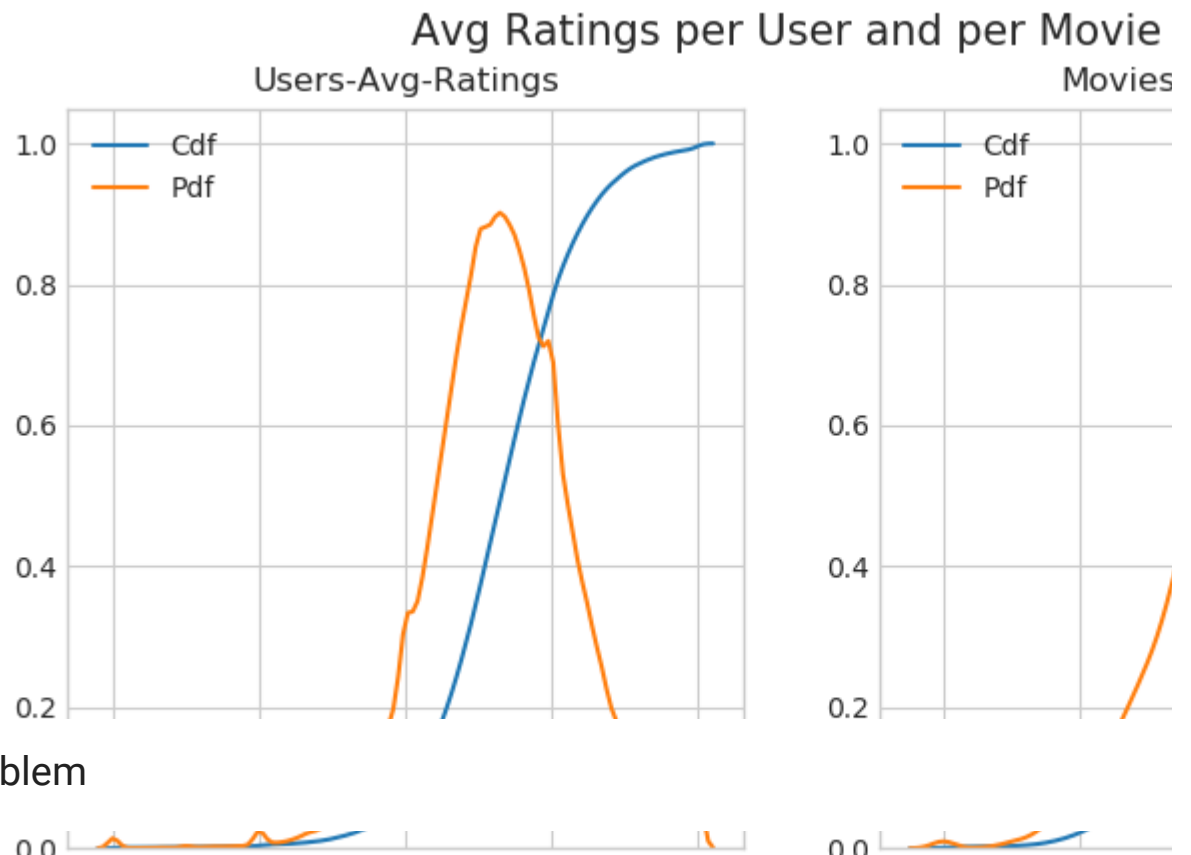
Double-click (or enter) to edit

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

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

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

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

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

## Avg Ratings per User and per Movie

### Users-Avg-Ratings / Movies



## Cold Start problem

## Cold Start problem with Users

```
0:00:35 003443
total_users = len(np.unique(df.user))
users_train = len(train_averages['user'])
new_users = total_users - users_train

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

```
    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.

## Cold Start problem with Movies

```
total_movies = len(np.unique(df.movie))
```

```
movies_train = len(train_averages['movie'])
new_movies = total_movies - movies_train

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

    Total number of Movies  : 17770

    Number of Users in Train data : 17424

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


Double-click (or enter) to edit

# Computing Similarity matrice

## Computing Movie-Movie Similarity matrix

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

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

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

    It seems you don't have that file. Computing movie_movie similarity...
    Done..
    Saving it to disk without the need of re-computing it again..
    Done..
    It's a  (17771, 17771)  dimensional matrix
    0:10:02.736054

```
m_m_sim_sparse.shape
```

```
(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.

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

```
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
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])
```

Double-click (or enter) to edit

# Finding most similar movies using similarity matrix

_ Does Similarity really works as the way we expected...? __
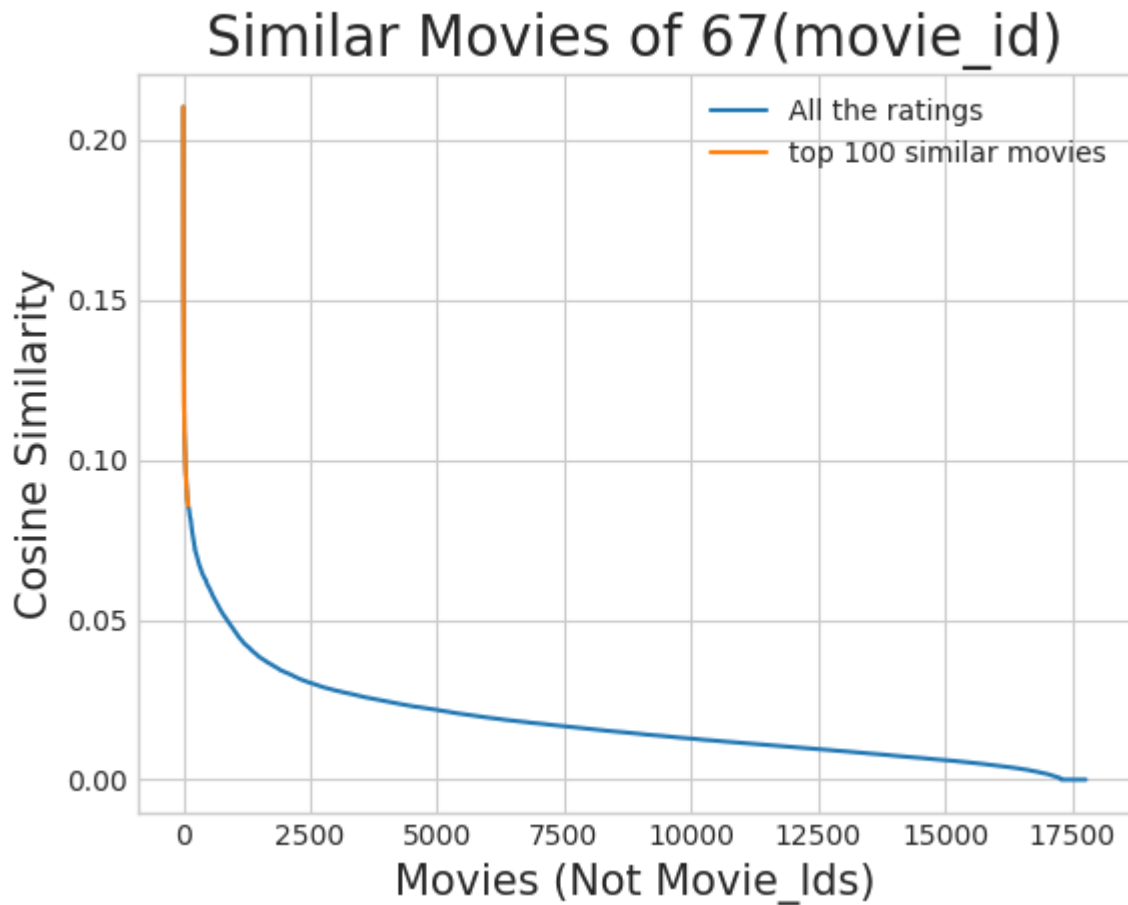_Let's pick some random movie and check for its similar movies....

```
# First Let's load the movie details into soe dataframe..
```

```python
# movie details are in 'netflix/movie_titles.csv'

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

movie_titles.head()
```

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

|  | 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'

```python
mv_id = 67

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

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

print("\nWe have {} movies which are similarto this  and we will get only top most..".format(
```

```
Movie -----> Vampire Journals

It has 270 Ratings from users.

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

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

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

similarities[similar_indices]

sim_indices = similarities.argsort()[::-1][1:] # It will sort and reverse the array and ignor
                                       # and return its indices(movie_ids)
```

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



**Top 10 similar movies**

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

| movie_id | year_of_release | title |
|---|---|---|
| 323 | 1999.0 | Modern Vampires |
| 4044 | 1998.0 | Subspecies 4: Bloodstorm |
| 1688 | 1993.0 | To Sleep With a Vampire |

Double-click (or enter) to edit

| | | |
|---|---|---|
| 12000 | 1998.0 | Dracula Rising |

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

| | | |
|---|---|---|
| 1900 | 1997.0 | Club Vampire |
| 13873 | 2001.0 | The Breed |
| 15867 | 2003.0 | Dracula II: Ascension |