# SGD Algorithm to predict movie ratings

**There will be some functions that start with the word "grader" ex: grader_matrix(), grader_mean(), grader_dim() etc, you should not change those function definition.**

**Every Grader function has to return True.**

1. Download the data from _here_

2. The data will be of this format, each data point is represented as a triplet of user_id, movie_id and rating

| user_id | movie_id | rating |
|---------|----------|--------|
| 77 | 236 | 3 |
| 471 | 208 | 5 |
| 641 | 401 | 4 |
| 31 | 298 | 4 |
| 58 | 504 | 5 |
| 235 | 727 | 5 |

# Task 1

**Predict the rating for a given (user_id, movie_id) pair**

**Predicted rating $\hat{y}_{ij}$ for user i, movied j pair is calcuated as** $\hat{y}_{ij} = \mu + b_i + c_j + u_i^T v_j$ **, here we will be finding the best values of** $b_i$ **and** $c_j$ **using SGD algorithm with the optimization problem for N users and M movies is defined as**
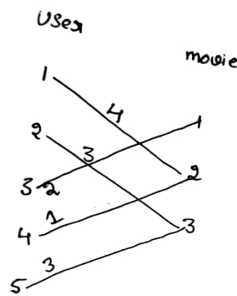
$$L = \min_{b,c,\{u_i\}_{i=1}^N,\{v_j\}_{j=1}^M} \alpha \qquad ($$

- $\mu$ : scalar mean rating
- $b_i$ : scalar bias term for user $i$
- $c_j$ : scalar bias term for movie $j$
- $u_i$ : K-dimensional vector for user $i$
- $v_j$ : K-dimensional vector for movie $j$

**\*. We will be giving you some functions, please write code in that functions only.**

**\*. After every function, we will be giving you expected output, please make sure that you get that output.**

1. Construct adjacency matrix with the given data, assuming its graph and the weight of each edge is the rating given by user to the movie



you can construct this matrix like $A[i][j] = r_{ij}$ here $i$ is user_id, $j$ is movie *id and $r_{ij}$ is rating given by user $i$ to the movie $j$*

Hint : you can create adjacency matrix using  csr_matrix

1. We will Apply SVD decomposition on the Adjaceny matrix  link1,  link2 and get three matrices $U, \sum, V$ such that $U \times \sum \times V^T$ ,
   $$= A$$
   if $A$ is of dimensions $N \times M$ then
   **U is of** $N \times k$,
   $\sum$ is of $k \times k$ and
   $V$ is $M \times k$ dimensions.

   \*. So the matrix $U$ can be represented as matrix representation of users, where each row $u_i$ represents a k-dimensional vector for a user

   \*. So the matrix $V$ can be represented as matrix representation of movies, where each row $v_j$ represents a k-dimensional vector for a movie.

2. Compute $\mu$ , $\mu$ represents the mean of all the rating given in the dataset.(write your code in **def m_u()**)
3. For each unique user initilize a bias value $B_i$ to zero, so if we have $N$ users $B$ will be a $N$ dimensional vector, the $i^{th}$ value of the $B$ will corresponds to the bias term for $i^{th}$ user (write your code in **def initialize()**)
4. For each unique movie initilize a bias value $C_j$ zero, so if we have $M$ movies $C$ will be a $M$ dimensional vector, the $j^{th}$ value of the $C$ will corresponds to the bias term for $j^{th}$ movie (write your code in **def initialize()**)
5. Compute dL/db_i (Write you code in **def derivative_db()**)
6. Compute dL/dc_j(write your code in **def derivative_dc()**
7. Print the mean squared error with predicted ratings.

```
for each epoch:
    for each pair of (user, movie):
        b_i =  b_i - learning_rate * dL/db_i
        c_j =  c_j - learning_rate * dL/dc_j
predict the ratings with formula
```

$$\hat{y}_{ij} = \mu + b_i + c_j$$
$$+ \text{dot\_product}$$
$$(u_i, v_j)$$

1. you can choose any learning rate and regularization term in the range $10^{-3}$ to $10^2$
2. **bonus:** instead of using SVD decomposition you can learn the vectors $u_i$, $v_j$ with the help of SGD algo similar to $b_i$ and $c_j$

In [278]:

```
import warnings
warnings.filterwarnings("ignore")
```

In [ ]:

# Task 2

As we know U is the learned matrix of user vectors, with its i-th row as the vector ui for user i. Each row of U can be seen as a "feature vector" for a particular user.

The question we'd like to investigate is this: do our computed per-user features that are optimized for predicting movie ratings contain anything to do with gender?

The provided data file **user_info.csv** contains an is_male column indicating which users in the dataset are male. Can you predict this signal given the features U?

> **Note 1** : there is no train test split in the data, the goal of this assignment is to give an intution about how to do matrix factorization with the help of SGD and application of truncated SVD. for better understanding of the collabarative fillerting please check netflix case study.
>
> **Note 2** : Check if scaling of $U, V$ matrices improve the metric

**Reading the csv file**

In [279]:

```
from sklearn.metrics import confusion_matrix
from sklearn.metrics import plot_confusion_matrix
from sklearn import tree
from google.colab import files
import io
import pandas as pd
```

In [280]:

```
uploaded = files.upload ()
```

Choose File   No file selected

Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

Saving ratings_train.csv to ratings_train (1).csv

In [281]:

```
data=pd.read_csv(io.BytesIO(uploaded['ratings_train.csv']))
data.head()
```

Out[281]:

|   | user_id | item_id | rating |
|---|---------|---------|--------|
| 0 | 772 | 36 | 3 |
| 1 | 471 | 228 | 5 |
| 2 | 641 | 401 | 4 |
| 3 | 312 | 98 | 4 |
| 4 | 58 | 504 | 5 |

In [282]:

```
data.shape
```

Out[282]:

```
(89992, 3)
```

### Create your adjacency matrix

In [283]:

```
users =data['user_id'].unique()
len(users)
items = data['item_id'].unique()
len(items)
ratings = data['rating']
ratings.shape
```

Out[283]:

```
(89992,)
```

In [284]:

```
#csr_matrix((data, (row_ind, col_ind)), [shape=(M, N)])
from scipy.sparse import csr_matrix
adjacency_matrix = csr_matrix((data['rating'],(data['user_id'],data['item_id'])))# write
your code of adjacency matrix here
```

In [285]:

```
adjacency_matrix.shape
```

Out[285]:

```
(943, 1681)
```

In [286]:

```
adj = adjacency_matrix.toarray()
```

In [287]:

```
adj[0][0]
```

Out[287]:

```
5
```

### Grader function - 1

```
def grader_matrix(matrix):
  assert(matrix.shape==(943,1681))
    return True
grader_matrix(adjacency_matrix)
```

```
True
```

**The unique items in the given csv file are 1662 only . But the id's vary from 0-1681 but they are not continuous and hence you'll get matrix of size 943x1681.**

## SVD decompostion

**Sample code for SVD decompostion**

```
from sklearn.utils.extmath import randomized_svd
import numpy as np
matrix = np.random.random((20, 10))
U, Sigma, VT = randomized_svd(matrix, n_components=5,n_iter=5, random_state=None)
print(U.shape)
print(Sigma.shape)
print(VT.T.shape)
```

```
(20, 5)
(5,)
(10, 5)
```

**Write your code for SVD decompostion**

```
# Please use adjacency_matrix as matrix for SVD decompostion
# You can choose n_components as your choice

from sklearn.utils.extmath import randomized_svd
import numpy as np
U, Sigma, VT = randomized_svd(adjacency_matrix, n_components=50,n_iter=5, random_state=N
one)
print(U.shape)
print(Sigma.shape)
print(VT.T.shape)
```

```
(943, 50)
(50,)
(1681, 50)
```

## Compute mean of ratings

```
def m_u(ratings):
    '''In this function, we will compute mean for all the ratings'''
    # you can use mean() function to do this
    # check this (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataF
rame.mean.html) link for more details.
    Mean = np.round(np.mean(ratings),3)


    return Mean
```

```
mu=m_u(data['rating'])
```

```
print(mu)
```

```
3.529
```

In [293]:

```
def grader_mean(mu):
    assert(np.round(mu,3)==3.529)
    return True
mu=m_u(data['rating'])
grader_mean(mu)
```

Out[293]:

```
True
```

## Initialize
## $B_i$ and
## $C_j$

**Hint : Number of rows of adjacent matrix corresponds to user dimensions( $B_i$), number of columns of adjacent matrix corresponds to movie dimensions ($C_j$)**

In [294]:

```
def initialize(dim):
    '''In this function, we will initialize bias value 'B' and 'C'.'''
    # initalize the value to zeros
    # return output as a list of zeros
    Bi = np.zeros(dim)
    return Bi
```

In [295]:

```
dim=adjacency_matrix.shape[0]# give the number of dimensions for b_i (Here b_i correspond
s to users)
b_i=initialize(dim)
b_i.sum()
```

Out[295]:

```
0.0
```

In [296]:

```
dim= adjacency_matrix.shape[1]# give the number of dimensions for c_j (Here c_j correspon
ds to movies)
c_j=initialize(dim)
c_j.sum()
```

Out[296]:

```
0.0
```

In [297]:

```
def grader_dim(b_i,c_j):
    assert(len(b_i)==943 and np.sum(b_i)==0)
    assert(len(c_j)==1681 and np.sum(c_j)==0)
    return True
grader_dim(b_i,c_j)
```

Out[297]:

```
True
```

## Compute dL/db_i

In [298]:

```python
data['rating'].shape
```

Out[298]:

```
(89992,)
```

In [299]:

```python
def derivative_db(user_id,item_id,rating,U,V,mu,alpha,b,c):
    '''In this function, we will compute dL/db_i'''
    alpha = 0.01
    p1 = 2*b-2*(rating-mu-b-c-np.dot(U[user_id],V[:,item_id]))
    #print(mu+np.sum(np.dot(U,V)))

    return p1
```

## Grader function -4

In [300]:

```python
def grader_db(value):
    assert(np.round(value,3)==-0.931)
    return True
```

In [301]:

```python
U1, Sigma, V1 = randomized_svd(adjacency_matrix, n_components=2,n_iter=5, random_state=2
4)
# Please don't change random state
# Here we are considering n_componets = 2 for our convinence
alpha=0.01
b = 0
c = 0
value=derivative_db(312,98,4,U1,V1,mu,alpha,b,c)
print(np.round(value,3))
grader_db(value)
```

```
-0.932

---------------------------------------------------------------------------
AssertionError                            Traceback (most recent call last)
<ipython-input-301-227f898d2330> in <module>()
      7 value=derivative_db(312,98,4,U1,V1,mu,alpha,b,c)
      8 print(np.round(value,3))
----> 9 grader_db(value)

<ipython-input-300-4216b7779205> in grader_db(value)
      1 def grader_db(value):
----> 2     assert(np.round(value,3)==-0.931)
      3     return True

AssertionError:
```

## Compute dL/dc_j

In [302]:

```python
def derivative_dc(user_id,item_id,rating,U,V,mu,alpha,b,c):
    '''In this function, we will compute dL/dc_j'''
    alpha = 0.01
    p1 = 2*c-2*(rating-mu-b-c-np.dot(U[user_id],V[:,item_id]))

    return p1
```

In [303]:

```
def grader_dc(value):
    assert(np.round(value,3)==-2.929)
    return True
U1, Sigma, V1 = randomized_svd(adjacency_matrix, n_components=2,n_iter=5, random_state=2
4)
# Please don't change random state
# Here we are considering n_componets = 2 for our convinence
r=0.01
value=derivative_dc(58,504,5,U1,V1,mu,alpha,b,c)
print(round(value,3))
grader_dc(value)
#print(value)

U1.shape
```

-2.93

```
---------------------------------------------------------------------------
AssertionError                            Traceback (most recent call last)
<ipython-input-303-3e58a37f03fc> in <module>()
      8 value=derivative_dc(58,504,5,U1,V1,mu,alpha,b,c)
      9 print(round(value,3))
---> 10 grader_dc(value)
     11 #print(value)
     12

<ipython-input-303-3e58a37f03fc> in grader_dc(value)
      1 def grader_dc(value):
----> 2     assert(np.round(value,3)==-2.929)
      3     return True
      4 U1, Sigma, V1 = randomized_svd(adjacency_matrix, n_components=2,n_iter=5, random_
state=24)
      5 # Please don't change random state

AssertionError:
```

## Compute MSE (mean squared error) for predicted ratings

**for each epoch, print the MSE value**

```
    for each epoch:

        for each pair of (user, movie):

            b_i =  b_i - learning_rate * dL/db_i

            c_j =  c_j - learning_rate * dL/dc_j

    predict the ratings with formula
```

$$\hat{y}_{ij} = \mu$$
$$+ b_i$$
$$+ c_j$$
$$+ \text{dot\_product}(u_i, v_j)$$

In [304]:

```
adj = adjacency_matrix.toarray()
```

In [305]:

```
from tqdm import tqdm
```

```
Y = []
alpha = 0.01
i = 50
MSE = []
learning_rate = 0.01
mu=m_u(data['rating'])
y_ = np.zeros((943,1681))
b_i_new = b_i
c_j_new = c_j
mse = 0
count  = 0
while(i>0):
  mse = 0
  for j in range(len(users)):
    for k in range(len(items)):
      u = users[j]
      m = items[k]
      if adj[j][k] != 0:
        db = derivative_db(u,m,adj[j][k],U,VT,mu,alpha,b_i[j],c_j[k]) # finding db
        dc = derivative_dc(u,m,adj[j][k],U,VT,mu,alpha,b_i[j],c_j[k]) #finding dc
        b_i_new[j] = b_i[j] - learning_rate*db
        c_j_new[k] = c_j[k] - learning_rate*dc
        y_[i][j] = b_i[j]+c_j[k]+mu+np.dot(U[users[j]],VT[:,items[k]])
        b_i[j] = b_i_new[j]
        c_j[k] = c_j_new[k]
        mse+=(adj[j][k]-(b_i_new[j]+c_j_new[k]+mu+np.dot(U[users[j]],VT[:,items[k]])))**
2 #finding mean squared error
  MSE.append(mse/89992)
  i-=1
```

In [306]:

```
MSE
```

Out[306]:

```
[0.9498344969653528,
 0.9008469577201043,
 0.8921194628385596,
 0.8884975370456183,
 0.8865746098491997,
 0.8854042095024072,
 0.8846229699960352,
 0.884065906667918,
 0.8836488090083722,
 0.8833247279512504,
 0.8830655941570061,
 0.882853644222159,
 0.8826770959848079,
 0.8825278219682089,
 0.8824000277469866,
 0.8822894652363465,
 0.8821929450024265,
 0.8821080230228298,
 0.8820327931261032,
 0.881965745620444,
 0.881905668628241,
 0.8818515777161359,
 0.8818026647252092,
 0.8817582599131273,
 0.8817178035058312,
 0.881680824019563,
 0.8816469215320519,
 0.8816157546248791,
 0.8815870300856735,
 0.8815604947094978,
 0.881535928715931,
 0.8815131404210277,
 0.8814919618940983,
 0.8814722453932639,
 0.8814538604221718,
 0.881436691285356,
```

```
0.8814055998175827,
0.8813915033094053,
0.8813782716124058,
0.881365838153029,
0.8813541428046836,
0.8813431311249714,
0.8813327536990094,
0.8813229655719547,
0.8813137257573905,
0.8813049968094545,
0.8812967444495341,
0.8812889372396809,
0.8812815462956763]
```

**Plot epoch number vs MSE**

- **epoch number on X-axis**
- **MSE on Y-axis**

In [307]:

```python
import matplotlib.pyplot as plt
```

In [308]:

```python
plt.plot(np.arange(1,51),MSE)
plt.xlabel('epoch')
plt.ylabel('MSE')
plt.title('MSE v/s epoches')
```

Out[308]:

```
Text(0.5, 1.0, 'MSE v/s epoches')
```



# Task 2

- **For this task you have to consider the user_matrix U and the user_info.csv file.**
- **You have to consider is_male columns as output features and rest as input features. Now you have to fit a model by posing this problem as binary classification task.**
- **You can apply any model like Logistic regression or Decision tree and check the performance of the model.**
- **Do plot confusion matrix after fitting your model and write your observations how your model is performing in this task.**
- **Optional work- You can try scaling your U matrix.Scaling means changing the values of n_componenets while performing svd and then check your results.**

```
U
```

```
array([[ 0.0662257 ,  0.00788853, -0.01253125, ...,  0.01367393,
        -0.01599038,  0.07419343],
       [ 0.01364432, -0.04889502,  0.05655371, ..., -0.01525794,
         0.00837367, -0.01568815],
       [ 0.00543826, -0.0251278 ,  0.02002774, ..., -0.02052443,
         0.02072003, -0.02445033],
       ...,
       [ 0.00738924, -0.02597375,  0.0063433 , ...,  0.02178487,
        -0.01543472,  0.00302407],
       [ 0.02499924,  0.00447791,  0.02605644, ...,  0.03279804,
        -0.02790097, -0.04015734],
       [ 0.04337341, -0.00281487, -0.0607779 , ..., -0.02570051,
        -0.0559467 ,  0.07182758]])
```

```
data1 = pd.read_csv('user_info.csv.txt')
data1.shape
```

```
(943, 4)
```

```
target = data1['is_male']
target
```

```
0      1
1      0
2      1
3      1
4      0
      ..
938    0
939    1
940    1
941    0
942    1
Name: is_male, Length: 943, dtype: int64
```

```
data1 = data1.drop(['is_male'],axis = 1)
```

```
data2 = pd.DataFrame(U)
```

```
data1.shape
```

```
(943, 3)
```

```
data2.shape
```

```
(943, 50)
```

```
In [316]:
```

```python
data = pd.concat([data1,data2],axis = 1)
```

```
In [317]:
```

```python
data
```

```
Out[317]:
```

| | user_id | age | orig_user_id | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 24 | 1 | 0.066226 | 0.007889 | -0.012531 | -0.086164 | 0.024869 | 0.006658 | 0.080034 | -0.027573 | 0.067700 | 0.0 |
| 1 | 1 | 53 | 2 | 0.013644 | -0.048895 | 0.056554 | 0.015810 | -0.012037 | 0.017731 | 0.010700 | -0.010228 | 0.028445 | 0.0 |
| 2 | 2 | 23 | 3 | 0.005438 | -0.025128 | 0.020028 | 0.032832 | 0.035080 | 0.001921 | 0.007691 | -0.000993 | -0.021173 | 0.0 |
| 3 | 3 | 24 | 4 | 0.005704 | -0.018211 | 0.010898 | 0.021867 | 0.013920 | -0.014181 | 0.012242 | -0.009123 | -0.012769 | 0.0 |
| 4 | 4 | 33 | 5 | 0.034122 | 0.009005 | -0.044054 | -0.016049 | 0.004326 | -0.021503 | 0.095574 | 0.079511 | -0.017195 | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 938 | 938 | 26 | 939 | 0.010350 | -0.038006 | 0.006501 | -0.013989 | 0.051223 | 0.001718 | 0.037136 | 0.010857 | 0.010762 | 0.0 |
| 939 | 939 | 32 | 940 | 0.031624 | -0.007730 | 0.032983 | 0.013862 | 0.023619 | -0.008443 | 0.054688 | -0.031091 | -0.015142 | 0.0 |
| 940 | 940 | 20 | 941 | 0.007389 | -0.025974 | 0.006343 | -0.017067 | -0.007397 | -0.020780 | 0.015469 | 0.015052 | 0.000977 | 0.0 |
| 941 | 941 | 48 | 942 | 0.024999 | 0.004478 | 0.026056 | 0.077343 | -0.000767 | -0.038300 | -0.010409 | -0.016338 | -0.011159 | 0.0 |
| 942 | 942 | 22 | 943 | 0.043373 | -0.002815 | -0.060778 | -0.031584 | 0.039834 | 0.006366 | -0.040937 | 0.069160 | 0.005817 | 0.0 |

**943 rows × 53 columns**

```
In [318]:
```

```python
data.isna().sum()
```

```
Out[318]:
```

```
user_id         0
age             0
orig_user_id    0
0               0
1               0
2               0
3               0
4               0
5               0
6               0
7               0
8               0
9               0
10              0
11              0
12              0
13              0
14              0
15              0
16              0
17              0
18              0
19              0
20              0
```

```
21              0
22              0
23              0
24              0
25              0
26              0
27              0
28              0
29              0
30              0
31              0
32              0
33              0
34              0
35              0
36              0
37              0
38              0
39              0
40              0
41              0
42              0
43              0
44              0
45              0
46              0
47              0
48              0
49              0
dtype: int64
```

In [319]:

```python
data.duplicated().sum()
```

Out[319]:

```
0
```

In [320]:

```python
target.value_counts()
```

Out[320]:

```
1    670
0    273
Name: is_male, dtype: int64
```

In [323]:

```python
from sklearn.model_selection import train_test_split

X_train , X_test ,y_train, y_test = train_test_split(data,target,test_size = 0.33,random
_state = 10)
```

In [324]:

```python
print(X_train.shape,y_train.shape)
print(X_test.shape,y_test.shape)
```

```
(631, 53) (631,)
(312, 53) (312,)
```

In [325]:

```python
from sklearn.model_selection import RandomizedSearchCV
from sklearn.linear_model import LogisticRegression
```

In [326]:

```python
Logistic = LogisticRegression(random_state = 10)
```

```
param = [{'C': [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
model = RandomizedSearchCV(Logistic,param,random_state = 10)
```

In [327]:

```
model.fit(X_train,y_train)
```

Out[327]:

```
RandomizedSearchCV(estimator=LogisticRegression(random_state=10),
                   param_distributions=[{'C': [0.0001, 0.01, 1, 100, 10000]}],
                   random_state=10)
```

In [328]:

```
model.best_estimator_
```

Out[328]:

```
LogisticRegression(C=10000, random_state=10)
```

In [329]:

```
model1 = LogisticRegression(C = 10000,random_state=10)
```

In [330]:
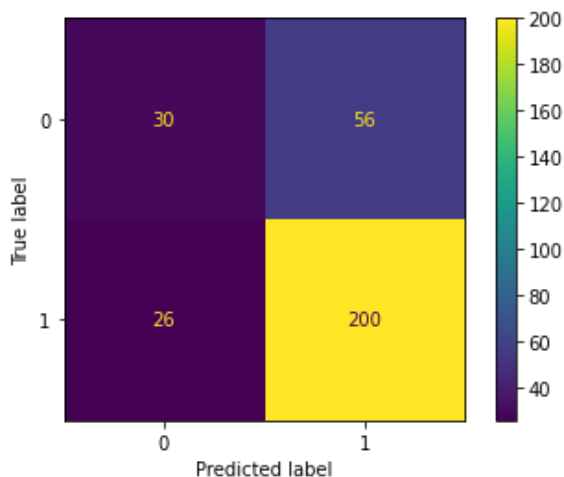
```
model1.fit(X_train,y_train)
```

Out[330]:

```
LogisticRegression(C=10000, random_state=10)
```

In [331]:

```
y_pred = model1.predict(X_test)
```

In [333]:

```
plot_confusion_matrix(model1,X_test,y_test)
plt.show()

roc_auc_score(y_test, model1.predict_proba(X_test)[:, 1])
```



Out[333]:

```
0.7430541263634493
```