

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](https://drive.google.com/open?id=1-1z7iDB52cB6_Jp07Dqa-e0YSs-mivpq) (https://drive.google.com/open?id=1-1z7iDB52cB6_Jp07Dqa-e0YSs-mivpq).
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

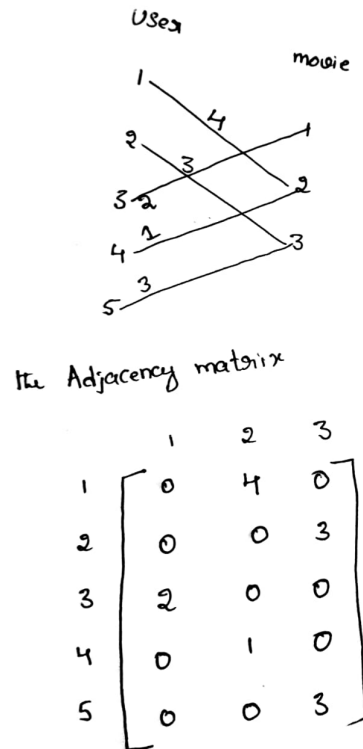
Predict the rating for a given (user_id, movie_id) pair

Predicted rating \hat{y}_{ij} for user i , movie j pair is calculated 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

$$L = \min_{b, c, \{u_i\}_{i=1}^N, \{v_j\}_{j=1}^M} \alpha \left(\sum_j \sum_k v_{jk}^2 + \sum_i \sum_k u_{ik}^2 + \sum_i b_i^2 + \sum_j c_j^2 \right) + \sum_{i,j \in \mathcal{I}^{\text{train}}} (y_{ij} - \mu - b_i - c_j - u_i^T v_j)^2$$

- (μ) : 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 [weighted un-directed bi-partited graph](https://en.wikipedia.org/wiki/Bipartite_graph) (https://en.wikipedia.org/wiki/Bipartite_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](https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html) (https://docs.scipy.org/doc/scipy/reference/generated/scipy.sparse.csr_matrix.html)

2. We will Apply SVD decomposition on the Adjacency matrix [link1](https://stackoverflow.com/a/31528944/4084039) (<https://stackoverflow.com/a/31528944/4084039>), [link2](https://machinelearningmastery.com/singular-value-decomposition-for-machine-learning/) (<https://machinelearningmastery.com/singular-value-decomposition-for-machine-learning/>) and get three matrices U , Σ , V such that $U \times \Sigma \times V^T = A$, if A is of dimensions $N \times M$ then
 - U is of $N \times k$,
 - Σ 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.

3. Compute μ , μ represents the mean of all the rating given in the dataset. (write your code in `def m_u()`)
4. For each unique user initialize 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()`)
5. For each unique movie initialize 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()`)

6. Compute dL/db_i (Write your code in `def derivative_db()`)
7. Compute dL/dc_j (write your code in `def derivative_dc()`)
8. 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

```

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

9. you can choose any learning rate and regularization term in the range 10^{-3} to 10^2
10. **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

Task 2

As we know U is the learned matrix of user vectors, with its i -th row as the vector u_i 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](https://drive.google.com/open?id=1PHFdJh_4gIPiLH5Q4UErH8GK71hTrzIY) (https://drive.google.com/open?id=1PHFdJh_4gIPiLH5Q4UErH8GK71hTrzIY) 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 intuition about how to do matrix factorization with the help of SGD and application of truncated SVD. for better understanding of the collaborative filtering please check netflix case study.

Note 2 : Check if scaling of U, V matrices improve the metric

Reading the csv file

In [3]:

```
1 from google.colab import files
2 uploaded = files.upload()
```

Choose Files No file chosen

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 [4]:

```
1 import pandas as pd
2 data=pd.read_csv('ratings_train.csv')
3 data.head()
```

Out[4]:

| | 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 [5]:

```
1 data.shape
```

Out[5]:

(89992, 3)

Create your adjacency matrix

In [0]:

```
1 from scipy.sparse import csr_matrix
2 adjacency_matrix =csr_matrix((data.rating.values, (data.user_id.values,
3 data.item_id.values)))
```

In [7]:

```
1 adjacency_matrix.shape
```

Out[7]:

(943, 1681)

Grader function - 1

In [8]:

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

Out[8]:

True

SVD decomposition

Sample code for SVD decomposition

In [9]:

```
1
2 from sklearn.utils.extmath import randomized_svd
3 import numpy as np
4 matrix = np.random.random((20, 10))
5 U, Sigma, VT = randomized_svd(matrix, n_components=5, n_iter=5, random_state=None)
6 print(U.shape)
7 print(Sigma.shape)
8 print(VT.T.shape)
```

(20, 5)

(5,)

(10, 5)

Write your code for SVD decomposition

In [10]:

```
1 # Please use adjacency_matrix as matrix for SVD decomposition
2 # You can choose n_components as your choice
3 from sklearn.utils.extmath import randomized_svd
4 import numpy as np
5 U, Sigma, VT = randomized_svd(adjacency_matrix, n_components=10, n_iter=5, random_state=None)
6 print(U.shape)
7 print(Sigma.shape)
8 print(VT.T.shape)
```

(943, 10)

(10,)

(1681, 10)

Compute mean of ratings

In [0]:

```

1 def m_u(ratings):
2     '''In this function, we will compute mean for all the ratings'''
3     # you can use mean() function to do this
4     # check this (https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.DataFrame.mean.html)
5     avg_rating=ratings.mean()
6
7     return avg_rating

```

Grader function -2

In [12]:

```

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

```

Out[12]:

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 [0]:

```

1 def initialize(dim):
2     '''In this function, we will initialize bias value 'B' and 'C'. '''
3     # initialize the value to zeros
4     # return output as a list of zeros
5     inta=np.zeros(dim)
6     return inta

```

In [0]:

```

1 dim=943 # give the number of dimensions for b_i (Here b_i corresponds to users)
2 b_i=initialize(dim)

```

In [0]:

```

1 dim=1681 # give the number of dimensions for c_j (Here c_j corresponds to movies)
2 c_j=initialize(dim)

```

Grader function -3

In [16]:

```

1 def grader_dim(b_i,c_j):
2     assert(len(b_i)==943 and np.sum(b_i)==0)
3     assert(len(c_j)==1681 and np.sum(c_j)==0)
4     return True
5 grader_dim(b_i,c_j)

```

Out[16]:

True

Compute dL/db_i

In [0]:

```

1 def derivative_db(user_id,item_id,rating,U,V,mu,alpha):
2     '''In this function, we will compute dL/db_i'''
3     reg=2*alpha*b_i[user_id]
4     loss=-2*(rating-mu-b_i[user_id]-c_j[item_id]-np.dot(U[user_id],V.T[item_id]))
5     der=reg+loss
6     return der
7

```

Grader function -4

In [18]:

```

1 def grader_db(value):
2     assert(np.round(value,3)==-0.931)
3     return True
4 U1, Sigma, V1 = randomized_svd(adjacency_matrix, n_components=2,n_iter=5, random_state=
5 # Please don't change random state
6 # Here we are considering n_componets = 2 for our convinence
7 alpha=0.01
8 value=derivative_db(312,98,4,U1,V1,mu,alpha)
9 grader_db(value)

```

Out[18]:

True

Compute dL/dc_j

In [0]:

```

1 def derivative_dc(user_id,item_id,rating,U,V,mu,alpha):
2     '''In this function, we will compute dL/dc_j'''
3     reg=2*alpha*c_j[item_id]
4     loss=-2*(rating-mu-b_i[user_id]-c_j[item_id]-np.dot(U[user_id],V.T[item_id]))
5     der=reg+loss
6     return der

```

Grader function - 5

In [20]:

```

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=
5 # Please don't change random state
6 # Here we are considering n_componets = 2 for our convinence
7 alpha=0.01
8 value=derivative_dc(58,504,5,U1,V1,mu,alpha)
9 grader_dc(value)

```

Out[20]:

True

Compute MSE (mean squared error) for predicted ratings

for each epoch, print the MSE value

for each epoch:

for each pair of (user, movie):

$b_i = b_i - \text{learning_rate} * dL/db_i$

$c_j = c_j - \text{learning_rate} * dL/dc_j$

predict the ratings with formula

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

In [21]:

```

1  from sklearn.metrics import mean_squared_error
2  rate=.01
3  y_act=data["rating"]
4  epochs=[]
5  mse=[]
6  for epoch in range(30):
7      epochs.append(epoch+1)
8      y_pred=[]
9      for user,item,rating in zip(data.iloc[:, 0], data.iloc[:, 1],data.iloc[:, 2]):
10         d_b=derivative_db(user,item,rating,U,VT,mu,alpha)
11         b_i[user]=b_i[user]-rate*d_b
12         d_c=derivative_dc(user,item,rating,U,VT,mu,alpha)
13         c_j[item]=c_j[item]-rate*d_c
14         for user,item,rating in zip(data.iloc[:, 0], data.iloc[:, 1],data.iloc[:, 2]):
15             pred=mu+b_i[user]+c_j[item]+np.dot(U[user],VT.T[item])
16             y_pred.append(pred)
17         m= mean_squared_error(y_act,y_pred)
18         mse.append(m)
19         print("--"+" " + "EPOCH"+" "+str(epoch+1))
20         print("MSE :",m)
21

```

```

-- EPOCH 1
MSE : 0.8884189183414518
-- EPOCH 2
MSE : 0.8618663902446467
-- EPOCH 3
MSE : 0.8522567714044784
-- EPOCH 4
MSE : 0.8476519295001488
-- EPOCH 5
MSE : 0.8450700701748326
-- EPOCH 6
MSE : 0.8434569326558222
-- EPOCH 7
MSE : 0.8423645653693673
-- EPOCH 8
MSE : 0.8415778603166821
-- EPOCH 9
MSE : 0.840983672873862
-- EPOCH 10
MSE : 0.8405179838861625
-- EPOCH 11
MSE : 0.8401422844595149
-- EPOCH 12
MSE : 0.8398321630607888
-- EPOCH 13
MSE : 0.8395714274218621
-- EPOCH 14
MSE : 0.8393489086179238
-- EPOCH 15
MSE : 0.8391566388613851
-- EPOCH 16
MSE : 0.8389887674421732
-- EPOCH 17
MSE : 0.8388408906000258
-- EPOCH 18
MSE : 0.8387096226557961
-- EPOCH 19

```

```
MSE : 0.8385923128016364
-- EPOCH 20
MSE : 0.8384868527152379
-- EPOCH 21
MSE : 0.8383915425055112
-- EPOCH 22
MSE : 0.8383049951488167
-- EPOCH 23
MSE : 0.838226066959076
-- EPOCH 24
MSE : 0.8381538060676555
-- EPOCH 25
MSE : 0.838087413620146
-- EPOCH 26
MSE : 0.8380262141215112
-- EPOCH 27
MSE : 0.8379696324746785
-- EPOCH 28
MSE : 0.8379171759921005
-- EPOCH 29
MSE : 0.8378684201538186
-- EPOCH 30
MSE : 0.8378229972238851
```

Plot epoch number vs MSE

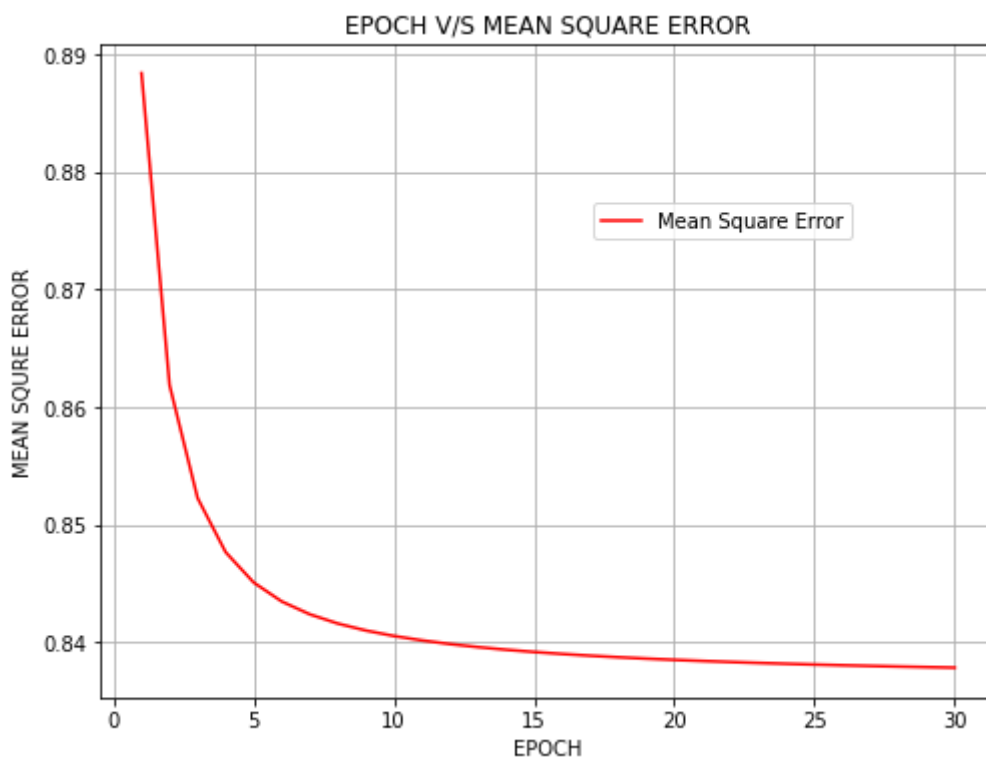
- epoch number on X-axis
- MSE on Y-axis

In [22]:

```
1 import matplotlib.pyplot as plt
2 x=epochs
3 y=mse
4 plt.figure(figsize=(8,6))
5 plt.plot(x,y,label='Mean Square Error',color="red")
6 plt.grid()
7 plt.xlabel("EPOCH")
8 plt.ylabel("MEAN SQUIRE ERROR")
9 plt.title("EPOCH V/S MEAN SQUARE ERROR")
10 plt.legend(loc=(.55,.7))
```

Out[22]:

<matplotlib.legend.Legend at 0x7ff456b11e80>



Task 2

In [24]:

```
1 from google.colab import files
2 uploaded = files.upload()
```

 No file chosen

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

Saving user_info.csv to user_info.csv

In [25]:

```
1 import pandas as pd
2 data1=pd.read_csv('user_info.csv')
3 data1.head()
```

Out[25]:

| | user_id | age | is_male | orig_user_id |
|---|---------|-----|---------|--------------|
| 0 | 0 | 24 | 1 | 1 |
| 1 | 1 | 53 | 0 | 2 |
| 2 | 2 | 23 | 1 | 3 |
| 3 | 3 | 24 | 1 | 4 |
| 4 | 4 | 33 | 0 | 5 |

In [0]:

```
1 X=U
2 Y=data1["is_male"]
```

In [0]:

```
1 from sklearn.linear_model import LogisticRegression
2 logreg = LogisticRegression(C=1e-4)
```

In [28]:

```
1 logreg.fit(X,Y)
```

Out[28]:

```
LogisticRegression(C=0.0001, class_weight=None, dual=False, fit_intercept=True,
intercept_scaling=1, l1_ratio=None, max_iter=100,
multi_class='auto', n_jobs=None, penalty='l2',
random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
warm_start=False)
```

In [29]:

```
1 from sklearn.metrics import accuracy_score
2 accuracy_score(Y,logreg.predict(X) )
```

Out[29]:

```
0.7104984093319194
```

In [30]:

```
1 from sklearn.metrics import confusion_matrix
2 confusion_matrix(Y,logreg.predict(X))
```

Out[30]:

```
array([[ 0, 273],
       [ 0, 670]])
```

- AFTER DOING SIMPLE LOGISTIC REGRESSION MODEL WE GET A ACCURACY OF 71% . SO WE CAN SAY THAT USER VECTOR (FEATURE VECTOR) OF RANDOMIZED SVD CONTAIN SOME AMOUNT OF GENDER INFORMATION.