

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

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

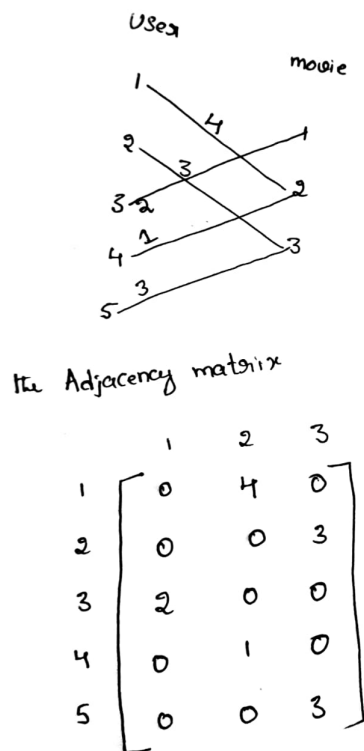
- (μ) : scalar mean rating
- (b_i) : scalar bias term for user (i)

- (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 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)

2. We will Apply SVD decomposition on the Adjacency matrix [link1](#)

(<https://stackoverflow.com/a/31528944/4084039>), [link2](#)

(<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 you 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 [1]: 1 import pandas as pd
        2 data=pd.read_csv('ratings_train.csv')
        3 data.head()
```

```
Out[1]:
```

	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 [2]: 1 min(data['user_id'])
```

```
Out[2]: 0
```

```
In [3]: 1 data.shape
```

```
Out[3]: (89992, 3)
```

Create your adjacency matrix

```
In [4]: 1 from scipy.sparse import csr_matrix
        2 adjacency_matrix = csr_matrix((data['rating'], (data['user_id'], data['item_i
```

```
In [5]: 1 adjacency_matrix.shape
```

```
Out[5]: (943, 1681)
```

Grader function - 1

```
In [6]: 1 def grader_matrix(matrix):
        2     assert(matrix.shape==(943,1681))
        3     return True
        4 grader_matrix(adjacency_matrix)
```

```
Out[6]: 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 decomposition

Sample code for SVD decomposition

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

(20, 5)

(5,)

(10, 5)

Write your code for SVD decomposition

```
In [8]: 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
5 U, Sigma, VT = randomized_svd(adjacency_matrix,
6                               n_components=150,
7                               n_iter=500,
8                               random_state=None)
```

Compute mean of ratings

```
In [9]: 1 def m_u(data):
2     return data['rating'].mean()
```

Grader function -2

```
In [10]: 1 def grader_mean(mu):
2     assert(np.round(mu, 3) == 3.529)
3     return True
4 mu = m_u(data)
5 grader_mean(mu)
```

Out[10]: 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 [32]: 1 b=[]
2 c=[]
3 def initialize(b,c):
4     b=np.zeros(adjacency_matrix.shape[0])
5     c=np.zeros(adjacency_matrix.shape[1])
6     return b,c
7 b,c=initialize(b,c)
8 print(b.shape,c.shape)
```

(943,) (1681,)

Grader function -3

```
In [12]: 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,c)
```

Out[12]: True

Compute dL/db_i

```
In [13]: 1 def derivative_db(user_id,item_id,rating,U,V,mu,alpha,b_i=np.zeros(adjacency_
2     '''In this function, we will compute dL/db_i'''
3     loss = (2*(alpha+1)*b_i[user_id]) - 2*(rating - mu - b_i[user_id] - c_j[
4     return loss
```

Grader function -4

```
In [14]: 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, ran
5 # Please don't change random state
6 # Here we are considering n_componets = 2 for our convinence
7 alpha=0.01
8 mu=m_u(data)
9 value=derivative_db(312,98,4,U1,V1,mu,alpha)
10 grader_db(value)
```

Out[14]: True

Compute dL/dc_j

```
In [15]: 1 def derivative_dc(user_id,item_id,rating,U,V,mu,alpha,b_i=np.zeros(adjacency_
2         '''In this function, we will compute dL/dc_j'''
3         loss = (2*(alpha+1)*c_j[item_id]) - 2*(rating - mu - b_i[user_id] - c_j[
4         return loss
```

Grader function - 5

```
In [16]: 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, ran
5 # Please don't change random state
6 # Here we are considering n_componets = 2 for our convinence
7 r=0.01
8 value=derivative_dc(58,504,5,U1,V1,mu,0.01)
9 grader_dc(value)
```

Out[16]: 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 [17]: 1 from sklearn.metrics import mean_squared_error as mse
2 def get_prediction(b_i,c_j,mu,df=data):
3     '''calculates net rmse'''
4     y_true = []
5     y_pred = []
6     for user,movie,rate in df[['user_id','item_id','rating']].values:
7         try:
8             y_hat = mu + b_i[user] + c_j[movie] + np.dot(U[user],VT[:,movie]).
9         except:
10            # handling cold start problem assigning global average for test u
11            y_hat = mu
12            y_true.append(rate)
13            y_pred.append(y_hat)
14    return mse(y_true,y_pred)

```

```

In [18]: 1 from tqdm import tqdm
2 def my_SGD(X, lr, alpha, u_mat, v_mat, epoch=30):
3     mu = m_u(data)
4     errors=[]
5     for i in tqdm(range(epoch)):
6         for user, movie, rating in zip(X.user_id.values, X.item_id.values, X.
7             b[user]=b[user] - lr * derivative_db(user,movie,rating,u_mat,v_mat,
8             c[movie]=c[movie] - lr * derivative_dc(user,movie,rating,u_mat,v_ma
9             error=get_prediction(b,c,mu)
10            errors.append(error)
11            # print('epoch: {0}, mse: {1} '.format(i+1, error))
12    return errors

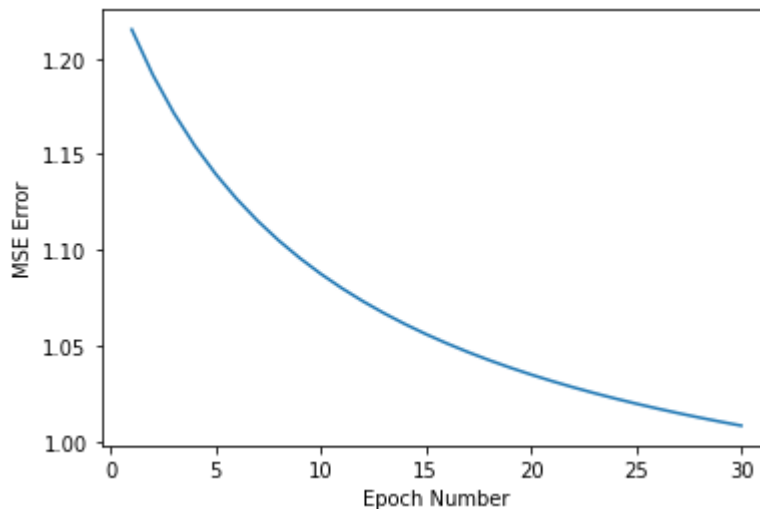
```

Plot epoch number vs MSE

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


```
In [19]: 1 epoch=30
2 my_errors = my_SGD(data, 0.0001,0.0001, U, VT,epoch)
3 import numpy as np
4 import matplotlib.pyplot as plt
5 x = np.arange(1,epoch+1)
6 y = my_errors
7 plt.plot(x, y)
8 plt.xlabel("Epoch Number")
9 plt.ylabel("MSE Error")
10 plt.show()
```

100%|██████████| 30/30 [01:06<00:00, 2.23s/it]



```
In [20]: 1 min(my_errors) #0.9318217355323781 1.008296586000108
```

Out[20]: 1.0082470622551105

```
In [21]: 1 max(my_errors)
```

Out[21]: 1.2150887920169118

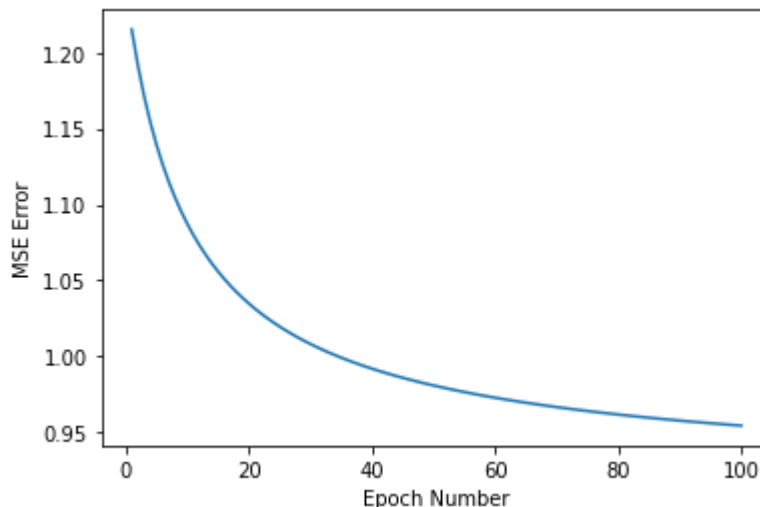
```
In [22]: 1 my_errors[29]
```

Out[22]: 1.0082470622551105

Before Running the code run the code block for initialisation of b,c.

```
In [33]: 1 epoch=100
2 my_errors = my_SGD(data, 0.0001,0.0001, U, VT,epoch)
3 import numpy as np
4 import matplotlib.pyplot as plt
5 x = np.arange(1,epoch+1)
6 y = my_errors
7 plt.plot(x, y)
8 plt.xlabel("Epoch Number")
9 plt.ylabel("MSE Error")
10 plt.show()
```

100%|██████████| 100/100 [03:56<00:00, 2.37s/it]



```
In [34]: 1 min(my_errors)
```

Out[34]: 0.954482225173638

```
In [35]: 1 max(my_errors)
```

Out[35]: 1.2150887920169118

```
In [36]: 1 my_errors[29]
```

Out[36]: 1.0082470622551105

Observation from the graph

First the b_i and c_j from random value reached to the value which provided minimum mse then due to constant changing of b_i and c_j the mse kept on increasing after reaching minima.

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_components while performing svd and then check your results.

```
In [23]: 1 data_male = pd.read_csv('user_info.csv')
          2 data_male.head()
```

```
Out[23]:
```

	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 [23]: 1
```

```
In [24]: 1 y_true = data_male.is_male.values
```

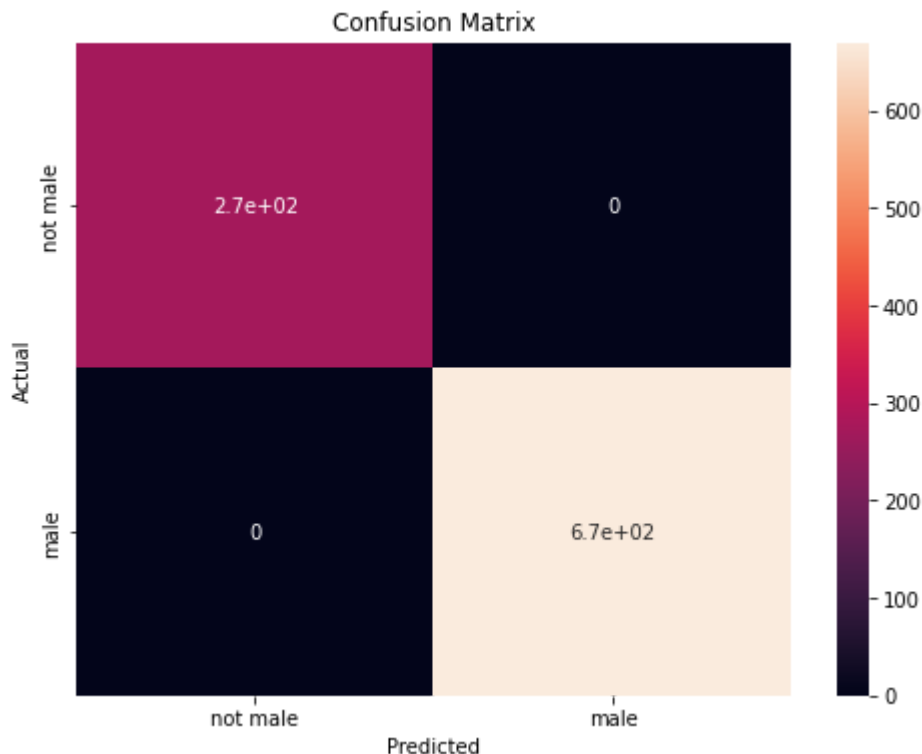
```
In [25]: 1 data_male['age']
          2 l=list(data_male['age'])
          3 l=np.array(l)
          4 l=l.reshape(-1,1)
```

```
In [26]: 1 from sklearn.tree import DecisionTreeClassifier
          2 clf = DecisionTreeClassifier()
          3 X=np.hstack((U,l))
          4 clf.fit(X, y_true)
          5 y_pred = clf.predict(X)
```

```
In [27]: 1 from sklearn.metrics import roc_auc_score
          2 acc = roc_auc_score(y_true, y_pred)
          3 print('accuracy: {0}'.format(acc))
```

accuracy: 1.0

```
In [28]: 1 from sklearn.metrics import confusion_matrix
2 from matplotlib import pyplot as plt
3 import seaborn as sns
4
5 c_matrix= confusion_matrix(y_true, y_pred)
6
7 df_cm = pd.DataFrame(c_matrix, index = [i for i in ['not male', 'male']],
8                        columns = [i for i in ['not male', 'male']])
9 plt.figure(figsize = (8,6))
10 sns.heatmap(df_cm, annot=True)
11 plt.title('Confusion Matrix')
12 plt.xlabel('Predicted')
13 plt.ylabel('Actual')
14 plt.show()
```



Observation Decision Tree Model

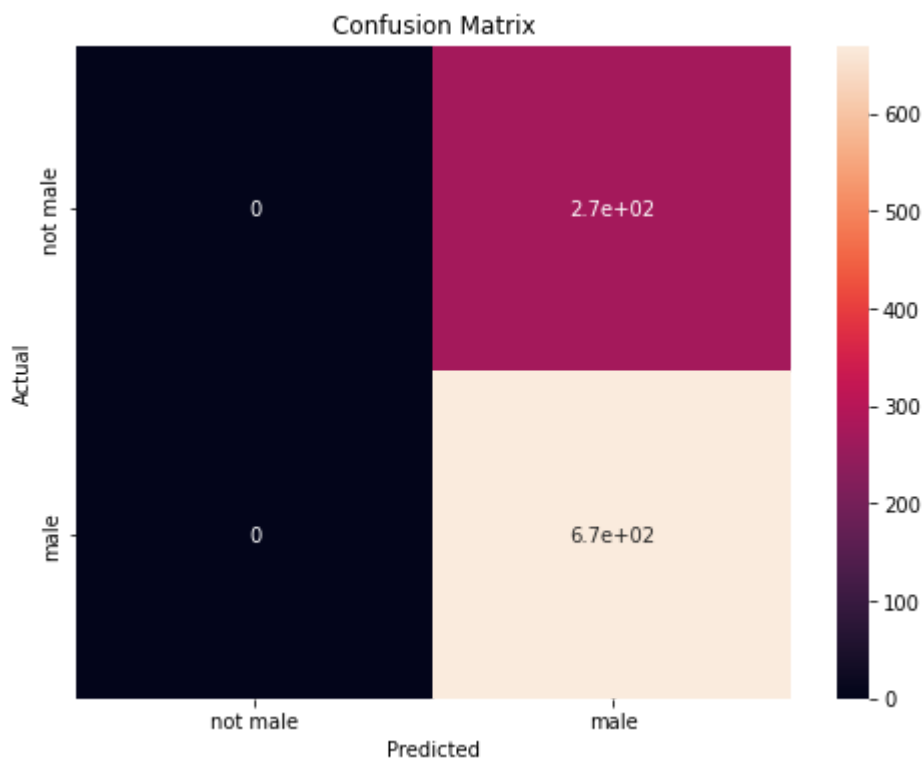
The model is highly accurate it is able to predict accurately male and non male members in the dataset in decision tree. I think this is because of default values of the hyperparameters i.e it can grow upto full depth and min_split is 2 only.

```
In [29]: 1 from sklearn.svm import SVC
2
3 clf = SVC(gamma='auto')
4 clf.fit(U, y_true)
5 y_pred = clf.predict(U)
```

```
In [30]: 1 from sklearn.metrics import roc_auc_score
2 acc = roc_auc_score(y_true, y_pred)
3 print('accuracy: {0}'.format(acc))
```

accuracy: 0.5

```
In [31]: 1 from sklearn.metrics import confusion_matrix
2 from matplotlib import pyplot as plt
3 import seaborn as sns
4
5 c_matrix= confusion_matrix(y_true, y_pred)
6
7 df_cm = pd.DataFrame(c_matrix, index = [i for i in ['not male', 'male']],
8                        columns = [i for i in ['not male', 'male']])
9 plt.figure(figsize = (8,6))
10 sns.heatmap(df_cm, annot=True)
11 plt.title('Confusion Matrix')
12 plt.xlabel('Predicted')
13 plt.ylabel('Actual')
14 plt.show()
```



Observation SVC Model

From the above SVC model we see that the model is highly inaccurate as it is not able to predict no male user.

There is no relation between U matrix and gender prediction.

