

SGD Algorithm to predict movie ratings

1. Download the data from [here](#)  
2. the data will be of this formate, 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**

- $\mu$ : 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  $j$

then the 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

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

### TASK: 1 \_\_SGD Algorithm to minimize the loss \_\_ 1. for each unique user initialize a bias value  $b_i$  randomly, 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 2. for each unique movie initialize a bias value  $C_j$  randomly, 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 3. Construct adjacency matrix with the given data, assuming its [weighted undirected bi-partited 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  $j$  4. we will Apply SVD decomposition on the Adjacency matrix [link1](#), [link2](#) and get three matrices  $U, \Sigma, V^T$  such that  $U \times \Sigma \times V^T = A$ , if  $A$  is of dimensions  $N \times M$  then  $U$  is of  $N \times k$ ,  $\Sigma$  is of  $k \times k$  and  $V^T$  is  $M \times k$  dimensions.

5. 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  $i$ . 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  $j$ .  $\mu$  represents the mean of all the rating given in the dataset

```
8.
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

ŷij = μ + bi + cj + dot_product(ui, vj)

print the mean squared error with predicted ratings
```

- 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  $c_j$
- ### 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](#) 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

<https://aaic-forum.slack.com/archives/GFLSLF27P/p1579583159062500>

```
In [1]: import pandas as pd
import numpy as np
import random
import networkx as nx
from sklearn.metrics import mean_squared_error
from scipy.sparse import coo_matrix, csr_matrix
from sklearn.utils.extmath import randomized_svd
import matplotlib.pyplot as plt

data = pd.read_csv('ratings_train.csv')
data.shape

Out[1]: (89992, 3)

In [2]: sorted_unique_users = np.sort(data.user_id.unique())
sorted_unique_movies = np.sort(data.item_id.unique())

In [3]: #Create a sparse Matrix with user as a row, movies as a column and rating as a data and then convert
sparse Matrix to dense you will get adjacency matrix
z = coo_matrix((data['rating'], (data.user_id, data.item_id)), shape=(943, 1681))
matrix = z.todense()

In [4]: k = 5

# Initialized the size of matrix as no. of users and movies
matrix = np.random.random((sorted_unique_users.size, sorted_unique_movies.size))
U, Sigma, VT = randomized_svd(matrix, n_components=k, n_iter=5, random_state=None)
print(U.shape)
print(Sigma.shape)
print(VT.T.shape)

(943, 5)
(5,)
(1681, 5)

In [5]: sorted_unique_users = np.sort(data.user_id.unique())
sorted_unique_movies = np.sort(data.item_id.unique())
N = len(data)

epoch = 30
learning_rate = random.choice([0.001,0.01,0.1,1,10,100])
reg_term = random.choice([0.001,0.01,0.1,1,10,100])

B = np.random.rand(sorted_unique_users.size) #Bias value B for users
C = np.random.rand(max(sorted_unique_movies)+1) #Bias value B for users
mu = data.rating.mean() #mu represents the mean of all the rating given in the dataset
y = data['rating'].values
y_pred = np.zeros(data.shape[0])
mse_list = []

print("Epoch_no \tmse_score")
for e in range(epoch):
    for index, (user, movie, rating) in enumerate(data.values):
        b_i, c_j = B[user], C[movie]

        #Dbi is dL/dbi
        Dbi = 2 * reg_term * b_i - 2 * (rating- mu - b_i - c_j - np.dot(U[user], VT.T[movie]))
        b_i = b_i - learning_rate * Dbi

        #Dcj is dL/dc_j
        Dcj = 2 * reg_term * c_j - 2 * (rating - mu - b_i - c_j - np.dot(U[user], VT.T[movie]))
        c_j = c_j - learning_rate * Dcj

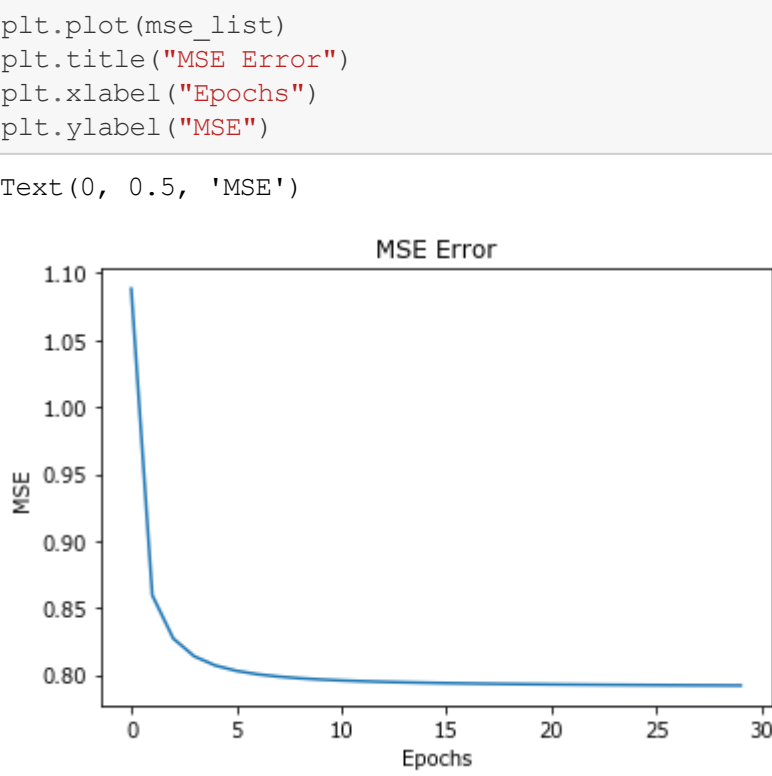
        #Updating b_i and c_j in B & C arrays
        B[user], C[movie] = b_i, c_j

        y_pred[index] = mu + b_i + c_j + np.dot(U[user], VT.T[movie])
        mse = mean_squared_error(y, y_pred)
        mse_list.append(mse)
        print(e, "\t",mse)

Epoch_no      mse_score
0              1.0883927951930348
1              0.8595241523017784
2              0.8271883058553882
3              0.8138901019480631
4              0.8070467228139505
5              0.8030457004935293
6              0.800492757789174
7              0.7987537257881843
8              0.7975068552218167
9              0.7965753426262686
10             0.7958557800022351
11             0.7952845081366725
12             0.794820596494635
13             0.79443671950269
14             0.7941140219507135
15             0.793839113771105
16             0.793602282904156
17             0.7933960981981261
18             0.7932151895986366
19             0.7930552215677726
20             0.792912828977716
21             0.7927853297872347
22             0.7926705599660774
23             0.7925667542004233
24             0.7924724582055307
25             0.7923864632634415
26             0.7923077566562033
27             0.7922354836478349
28             0.7921689179797982
29             0.792107438727257

In [6]: plt.plot(mse_list)
plt.title("MSE Error")
plt.xlabel("Epochs")
plt.ylabel("MSE")

Out[6]: Text(0, 0.5, 'MSE')
```



Updated ci to cj

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

```
In [7]: # Questions
# How to decide the right K
# To calculate v^2 jk over sum of all j and k, is it over all j or particular j
# How to calculate db_i and db_j

2 alpha bi - 2 (yij - mu - bi - ci - ui^Tvj)

2 alpha cj - 2 (yij - mu - bi - ci - ui^Tvj)
```

Task 2

```
In [8]: gender_data = pd.read_csv('user_info.csv')

In [9]: y = gender_data['is_male'].values
X = gender_data.drop('is_male', axis =1)

In [10]: new_X = pd.concat([X, pd.DataFrame(U, columns=list('ABCDE'))]), axis=1)

In [11]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(new_X, y, test_size = 0.25)

In [12]: from sklearn.tree import DecisionTreeClassifier
clf = DecisionTreeClassifier()
clf.fit(X_train, y_train)

Out[12]: DecisionTreeClassifier(class_weight=None, criterion='gini', max_depth=None,
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=1, min_samples_split=2,
min_weight_fraction_leaf=0.0, presort=False, random_state=None,
splitter='best')

In [13]: clf.score(X_test, y_test)

Out[13]: 0.597457627118644

In [14]: from sklearn.metrics import roc_auc_score
y_pred = clf.predict(np.array(X_test))
roc_auc_score(y_test, y_pred)

Out[14]: 0.5378151260504203
```

scaling U and VT

```
In [15]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
U_scaled = scaler.fit_transform(U)
VT_scaled = scaler.fit_transform(VT)

In [16]: sorted_unique_users = np.sort(data.user_id.unique())
sorted_unique_movies = np.sort(data.item_id.unique())
N = len(data)

epoch = 30
learning_rate = random.choice([0.001,0.01,0.1,1,10,100])
reg_term = random.choice([0.001,0.01,0.1,1,10,100])

B = np.random.rand(sorted_unique_users.size) #Bias value B for users
C = np.random.rand(max(sorted_unique_movies)+1) #Bias value B for users
mu = data.rating.mean() #mu represents the mean of all the rating given in the dataset
y = data['rating'].values
y_pred = np.zeros(data.shape[0])
mse_list = []

print("Epoch_no \tmse_score")
for e in range(epoch):
    for index, (user, movie, rating) in enumerate(data.values):
        b_i, c_j = B[user], C[movie]

        #Dbi is dL/dbi
        Dbi = 2 * reg_term * b_i - 2 * (rating- mu - b_i - c_j - np.dot(U_scaled[user], VT_scaled.T[movie]))
        b_i = b_i - learning_rate * Dbi

        #Dcj is dL/dc_j
        Dcj = 2 * reg_term * c_j - 2 * (rating - mu - b_i - c_j - np.dot(U_scaled[user], VT_scaled.T[movie]))
        c_j = c_j - learning_rate * Dcj

        #Updating b_i and c_j in B & C arrays
        B[user], C[movie] = b_i, c_j

        y_pred[index] = mu + b_i + c_j + np.dot(U_scaled[user], VT_scaled.T[movie])
        mse = mean_squared_error(y, y_pred)
        mse_list.append(mse)
        print(e, "\t",mse)

Epoch_no      mse_score
0              13.804552821506416
1              9.112838312454372
2              7.851924334330101
3              7.33083886135253
4              7.05213920659876
5              6.876639145285128
6              6.753900368315021
7              6.6619634616509701
8              6.589824974024179
9              6.531335075559372
10             6.48275367156266
11             6.441652347138368
12             6.406371433355669
13             6.3757294262230255
14             6.348565982964678
15             6.325096837228662
16             6.303940686561787
17             6.284987858384184
18             6.267917349101854
19             6.252468372875496
20             6.238426495170696
21             6.2256135321599935
22             6.21386005112602
23             6.203099706260709
24             6.193164894847303
25             6.183983800057747
26             6.175475632697321
27             6.167572715469684
28             6.160214582063829
29             6.15334869726591

In [17]: plt.plot(mse_list)
plt.title("MSE Error")
plt.xlabel("Epochs")
plt.ylabel("MSE")

Out[17]: Text(0, 0.5, 'MSE')
```

Scaling of U & V matrices does not improve the martric

