# The Experiment Report of Machine Learning

**SCHOOL:** SCHOOL OF SOFTWARE ENGINEERING

**SUBJECT:** SOFTWARE ENGINEERING

Author:
Siyuan Xiao, Li Zhang, Shengyan Wen

Supervisor:
Qingyao Wu

Student ID：
201721045886, 201721045909, 201721045893

Grade:
Postgraduate

December 24, 2017

# Recommender System Based on Matrix Decomposition

**Abstract—Nowadays, recommender system is a popular application that applying data mining techniques. In this field, Matrix Decomposition is used frequently.**

## I. INTRODUCTION

This report will talk about the whole experiment I have made for a recommender system based on Matrix Factorization. Its content is organized as follow:

1)   Section II contains the experiment steps.
2)   Section III contains the code for the experiment.
3)   Section IV makes conclusion for the experiment result.

## II. METHODS AND THEORY

Here we use the MovieLens-100k dataset which is placed in the requirement. The u1 database is used and it is already splitted into training set and validation set by the official.

Then the experiment will be performed by the following steps:

1)   Load the data text into memory and do the transformation to get the specific column, and the cells that lack of data is set to 0 (or the average value).
2)   For ALS method, we randomly initiate the Q matrix.
3)   After that, we calculate P and Q alternatively until the iteration number is arrived or the RMSE has little change.
4)   For SGD method, we randomly initiate the two matrices.
5)   In each iteration, we randomly pick one row in P and Q, respectively.
6)   We use these two selected rows to generate a prediction, then calculate the error and gradient for P and Q, respectively.
7)   Similarly, we iterate until it converges or reaches the given number of iterations

8)   Regardless of which method we use, do the whole prediction and calculate RMSE for the valid validation data.

## III. EXPERIMENT

Here I placed the critical code for the experiment:

1)   Preprocessing:

```
1.   # created by Swain, 2017-12-19, 11:29

2.   import numpy as np
3.   from numpy.linalg import inv
4.   import matplotlib.pyplot as plot
5.   from numpy import random

6.   #load u1 from MovieLens-100k dataset
7.   n_user = 943
8.   n_item = 1682
9.   trainpath = './u1.base'
10. testpath = './u1.test'

11. data_train = np.loadtxt(trainpath)
12. data_test = np.loadtxt(testpath)

13. def make_R(data):
14.     count = 0
15.     sum = 0
16.     R = np.zeros([n_user, n_item])
17.     for i in range(0, data.shape[0]):
18.         R[int(data[i, 0]) - 1, int(data[i, 1]) - 1] =
    int(data[i, 2])
19.         count = count + 1
20.         sum = sum + int(data[i, 2])
21.
22.     avg = sum / count
23.     '''for i in range(0, R_train.shape[0]):
24.         for j in range(0, R_train.shape[1]):
25.             if R_train[i, j] == 0:
26.                 R_train[i, j] = avg'''
27.     return R
28.
29. R_train = make_R(data_train)
30. R_vali = make_R(data_train)

31. Np = np.sum(R_train, axis=1)
32. Nq = np.sum(R_train, axis=0)

33. random.seed(24)
```

2) ALS:

```
1.   #ALS
2.   k = 10
3.   lambda_ = 0.03
4.   iteration = 1000

5.   P = random.random([n_user, k]) * 5
6.   Q = random.random([n_item, k]) * 5

7.   Lambda = lambda_ * np.eye(k)

8.   R = R_train

9.   for n in range(0, iteration):
10.      print ("start the " + str(n + 1) + "-th
     iteration")
11.
12.      #calculate P matrix
13.      for i in range(0, n_user):
14.         P[i] = np.dot(inv(np.dot(Q.T, Q) +
     Lambda), np.dot(Q.T, R[i]))
15.
16.      #calculate Q matrix
17.      for i in range(0, n_item):
18.         Q[i] = np.dot(inv(np.dot(P.T, P) +
     Lambda), np.dot(P.T, R[:, i]))
19.
20.      #do prediction
21.      R_predict = np.dot(P, Q.T)
22.
23.      #calculate error
24.      RMSE = 0
25.      count = 0
26.      for i in range(0, n_user):
27.         for j in range(0, n_item):
28.            if R_vali[i, j] != 0:
29.               RMSE = RMSE + np.power(R_vali[i,
     j] - R_predict[i, j], 2)
30.               count = count + 1
31.      RMSE = np.sqrt(RMSE / count)
32.      print ("  RMSE = " + str(RMSE))
33.

34.      if n > 0 and np.abs(RMSE - last_RMSE) <
     0.0001:
35.         break;
36.      last_RMSE = RMSE
```

3) SGD:

```
1.   #SGD
2.   k = 10
3.   lr = 0.05
4.   lambda_p = 0.03
5.   lambda_q = 0.03
6.   iteration = 1000000
7.
8.   #P = np.ones([n_user, k])
```

```
9.   #Q = np.ones([n_item, k])
10.  P = random.random([n_user, k])
11.  Q = random.random([n_item, k])
12.
13.  R = R_train
14.
15.  for n in range(0, iteration):
16.      #randomly pick a row from P and Q
17.      u = random.randint(0, n_user - 1)
18.      i = random.randint(0, n_item - 1)
19.
20.      #calculate error and gradient
21.      error = R[u, i] - np.dot(P[u], Q.T[:, i])
22.      deriv_p = lambda_p * P[u] - error * Q[i]
23.      deriv_q = lambda_q * Q[i] - error * P[u]
24.
25.      #do updation
26.      P[u] = P[u] - lr * deriv_p
27.      Q[i] = Q[i] - lr * deriv_q
28.
29.      print (str(n + 1) + "-th iteration: error = " +
     str(error))
30.
31.  #do prediction and calculate error
32.  R_predict = np.dot(P, Q.T)
33.  RMSE = 0
34.  count = 0
35.  for i in range(0, n_user):
36.      for j in range(0, n_item):
37.         if R_vali[i, j] != 0:
38.            RMSE = RMSE + np.power(R_vali[i, j] -
     R_predict[i, j], 2)
39.            count = count + 1
40.  RMSE = np.sqrt(RMSE / count)
41.  print ("  RMSE = " + str(RMSE))
```

## IV. CONCLUSION

For ALS method, it ends with the 52-nd iteration. For SGD method, we do 100,000 iterations. Here is the experiment result gained:

1) For ALS:

```
start the 51-th iteration
  RMSE = 2.45483996434
start the 52-th iteration
  RMSE = 2.45474236985
```

2) For SGD:

```
9999-th iteration: error = -0.318996351454
10000-th iteration: error = -0.0622771554502

  RMSE = 3.19995670944
```

Then we can draw a conclusion according to the experiment:

1)   ALS converges with far less iterations but may be likely to have overfitting problem.
2)   SGD should run much more iterations since the data used in each iteration just generate one prediction.
3)   If we set the undetermined cells of the training dataset to the average value (while the data shown is in case that setting to 0) of all training data, the RMSE of prediction can be reduced significantly.