**South China University of Technology**

# The Experiment Report of Machine Learning

**SCHOOL:** SCHOOL OF SOFTWARE ENGINEERING

**SUBJECT:** SOFTWARE ENGINEERING

Author:
Shoukai Xu and Yaofu Chen

Supervisor:
Qingyao Wu

Student ID：
201530611180 and 20153061380
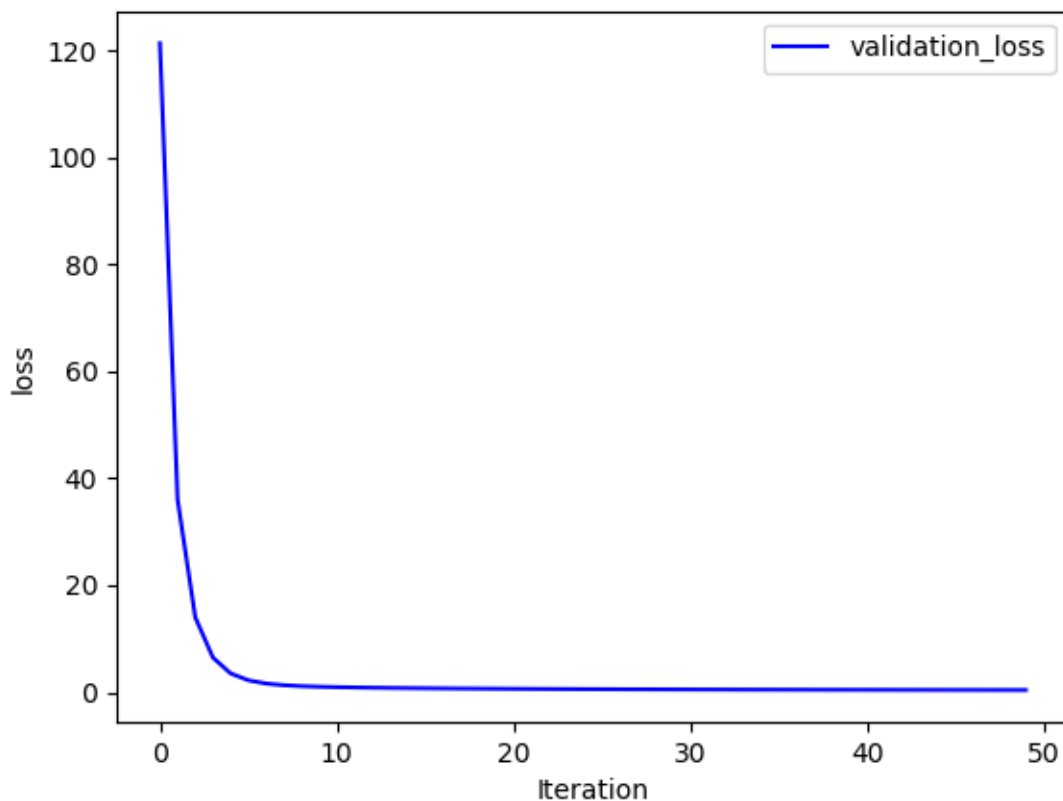and 201530613535

Grade:
Undergraduate

December 9, 2017

# Recommender System Based on Matrix Decomposition

**ABSTRACT**—

Our project is based on matrix decomposition and we take several steps to complete the task:
- Read Data
- Spilt Data
- Initialization
- Training
- Draw The Loss Figure

Here is the brief look at the result:

## I. INTRODUCTION

### i. Motivation

    a)   Explore the construction of recommended system.
    b)   Understand the principle of matrix decomposition.
    c)   Be familiar to the use of gradient descent.
    d)   Construct a recommendation system under small-scale dataset, cultivate engineering ability.

## ii. Dataset
   a) Utilizing MovieLens-100k dataset.
   b) u.data -- Consisting 10,000 comments from 943 users out of 1682 movies. At least, each user comment 20 videos. Users and movies are numbered consecutively from number 1 respectively. The data is sorted randomly
   c) u1.base / u1.test are train set and validation set respectively, seperated from dataset u.data with proportion of 80% and 20%. It also make sense to train set and validation set from u1.base / u1.test to u5.base / u5.test.
   d) You can also construct train set and validation set according to your own evaluation method.

## iii. Environment
Python3，at least containing following python package: sklearn，numpy，jupyter，matplotlib. An advice is installing anaconda3 directly, which already contains the python package mentioned above.


# II. METHOD AND THEORY
## i. ALS
   a) Read the data set and divide it (or use u1.base / u1.test to u5.base / u5.test directly). Populate the original scoring matrix $R_{nusers,nitems}$ against the raw data, and fill 0 for null values.
   b) Initialize the user factor matrix $P_{nusers,K}$ and the item (movie) factor matrix $Q_{nitem,K}$, where $K$ is the number of potential features.
   c) Determine the loss function and the hyperparameter learning rate ita and the penalty factor lamda.
   d) Use alternate least squares optimization method to decompose the sparse user score matrix, get the user factor matrix and item (movie) factor matrix:
      1) With fixd item factor matrix, find the loss partial derivative of each row (column) of the user factor matrices, ask the partial derivative to be zero and update the user factor matrices.
      2) With fixd user factor matrix, find the loss partial derivative of each row (column) of the item factor matrices, ask the partial derivative to be zero and update the item.
      3) Calculate the $L_{validation}$ on the validation set, comparing with the $L_{validation}$ of the previous iteration to determine if it has converged.

## ii. SGD
   a) Read the data set and divide it (or use u1.base / u1.test to u5.base / u5.test directly). Populate the original scoring matrix $R_{nusers,nitems}$ against the raw data, and fill 0 for null values.
   b) Initialize the user factor matrix $P_{nusers,K}$ and the item (movie) factor matrix $Q_{nitem,K}$, where $K$ is the number of potential features.
   c) Determine the loss function and the hyperparameter learning rate ita and the penalty factor lamda.
   d) Use alternate least squares optimization method to decompose the sparse user score matrix, get the user factor matrix and item (movie) factor matrix:
      1) S elect a sample from scoring matrix randomly;
      2) C alculate this sample's loss gradient of specific row(column) of user factor matrix and item factor matrix;
      3) Use SGD to update the specific row(column) of $P_{n\_users,K}$ and $Q_{n-items,K}$;
      4) Calculate the $L_{validation}$ on the validation set, comparing with the $L_{validation}$ of the previous iteration to determine if it has converged.

e) Repeat step 4. several times, get a satisfactory user factor matrix Q and an item factor matrix ,
   Draw a Lvalidation curve with varying iterations.

f) T he final score prediction matrix Rn_users,n_items is obtained by multiplying the user factor
   matrix Pn_users,K and the transpose of the item factor matrix Qn_items,K.

# III. EXPERIMENT

## i. Code

```
1  {
2    "cells": [
3      {
4        "cell_type": "code",
5        "execution_count": 2,
6        "metadata": {
7          "collapsed": true
8        },
9        "outputs": [
10         {
11           "data": {
12             "image/png": "iVBORw0KGgoAAAANSUhEUgAAAYgAAAEKCAYAAAAIO8L1AAAABHNCSVQICAgIfAhkiAAAAAlwSFlz\nAAALEgAACxIB0t1+/AAAAD10RVh0U29mdHdhcmUAbWF0cG
13             "text/plain": [
14               "<matplotlib.figure.Figure at 0x1a115ac9d30>"
15             ]
16           },
17           "metadata": {},
18           "output_type": "display_data"
19         }
20       ],
21       "source": [
22         "import numpy as np\n",
23         "import pandas as pd\n",
24         "import matplotlib.pyplot as plt\n",
25         "from sklearn.model_selection import train_test_split\n",
26         "\n",
27         "#读取数据\n",
28         "def get_data():\n",
29         "    header = ['user_id', 'item_id', 'rating', 'timestamp']\n",
30         "    df = pd.read_csv('ml-100k/u.data',sep='\\t',names=header,encoding='latin-1')\n",
31         "    return df\n",
32         "#分割数据\n",
33         "data = get_data()\n",
34         "n_users = data.user_id.unique().shape[0]\n",
35         "n_items = data.item_id.unique().shape[0]\n",
36         "#print(n_users)\n",
37         "#print(n_items)\n",
38         "train_data,test_data=train_test_split(data,test_size=0.25)\n",
39         "train_data= pd.DataFrame(train_data)\n",
40         "test_data= pd.DataFrame(test_data)\n",
41         "#全零初始化\n",
42         "train_matrix = np.zeros((n_users, n_items))\n",
43         "for line in train_data.itertuples():\n",
44         "    train_matrix[line[1]-1, line[2]-1] = line[3]\n",
45         "\n",
46         "test_matrix = np.zeros((n_users,n_items))\n",
47         "for line in test_data.itertuples():\n",
48         "    test_matrix[line[1]-1, line[2]-1] = line[3]\n",
49         "#print(train_matrix)\n",
50         "def train_and_test():\n",
51         "    k = 150\n",
52         "    penalty_factor = 0.01\n",
53         "    learning_rate = 0.0001\n",
54         "    epochs = 50\n",
55         "\n",
56         "    users = np.shape(train_matrix)[0]\n",
57         "    items = np.shape(train_matrix)[1]\n",
58         "    P = np.random.rand(users,k)\n",
59         "    Q = np.random.rand(k,items)\n",
60         "    validation_loss = []\n",
61         "\n",
```

```
62    "    for step in range(epochs):\n",
63    "        # 训练P矩阵和Q矩阵\n",
64    "        for u in range(users):\n",
65    "            for i in range(items):\n",
66    "                if train_matrix[u][i] > 0:\n",
67    "                    error = train_matrix[u][i] - np.dot(P[u], Q[:, i])\n",
68    "                    P[u] += 2 * learning_rate * error * Q[:, i] - 2 * penalty_factor * P[u]\n",
69    "                    Q[:, i] += 2 * learning_rate * error * P[u] - 2 * penalty_factor * Q[:, i]\n",
70    "        loss = np.sum(np.power(test_matrix - np.dot(P, Q), 2)) / (users * items) + penalty_factor * (\n",
71    "        np.sum(np.power(P, 2)) + np.sum(np.power(Q, 2))) / k\n",
72    "        validation_loss.append(loss)\n",
73    "\n",
74    "    return validation_loss\n",
75    "\n",
76    "loss = train_and_test()\n",
77    "#画出loss曲线图\n",
78    "plt.xlabel('Iteration')\n",
79    "plt.ylabel('loss')\n",
80    "plt.plot(loss, label=\"validation_loss\",color=\"blue\",linewidth=1.5)\n",
81    "plt.legend(loc='upper right')\n",
82    "plt.show()"
83    ]
84    }
85  ],
86  "metadata": {
87   "kernelspec": {
88    "display_name": "Python 2",
89    "language": "python",
90    "name": "python2"
91   },
92   "language_info": {
93    "codemirror_mode": {
94     "name": "ipython",
95     "version": 2
96    },
97    "file_extension": ".py",
98    "mimetype": "text/x-python",
99    "name": "python",
100   "nbconvert_exporter": "python",
101   "pygments_lexer": "ipython2",
102   "version": "2.7.6"
103  }
104  },
105  "nbformat": 4,
106  "nbformat_minor": 0
107 }
```

## ii. Train of Thought

a) Content Based

For each item, you need to create a profile for the item. For each user, create a user profile. For example, if you search for programming and search for food, then your user profile will show that you are right Programming, food is interested. So if we treat the item profile simply as the question mark for each question then you are probably interested in programming or food labeling issues.

b) Collaborative Filtering

The essence of cf is that he speculates what x wants to buy using a user's buying habits similar to user x. The breakdown of the cf algorithm is user-based and item-based.

c) Matrix Factorization

We imagine a matrix M, where each row represents a user, and each column represents a movie. Then M [i] [j] is the rating of user i for movie j. It may or may not exist (if i has not seen j). Assuming there are m users, n movies, then this matrix is the m * n matrix. Now, let's break it down into the product of two matrices of (m * f) * (f * n). (F is the number of movie features) The first is called user matrix and the second is item matrix. So what is the representation of each row in the user matrix: the degree of user's preference for each feature. For example, if we define features as "comedy, horror, martial arts," then the user u may be <1.0, 0.2, -0.2 in the user matrix line, which illustrates u's favorite comedy, hate terror and hate martial arts. Then in the item matrix, each movie corresponds to the extent of their "contained" features. For example, a pancake man might be <1.0, 0, 0.2>, and a horror film might be <0, 1, 0.5>. Then we count the dot product of these two vectors and get M [u] 1-0.04 = 0.96 M [u] [Horror] = 0.2-0.1 = 0.1.

# IV.  CONCLUTION