



南开大学  
Nankai University

南 开 大 学

计 网 学 院

大数据计算及应用实验报告

---

推荐系统编程作业

---

王隽毅、包烜灏、杨坤

年级：2019 级

专业：计算机科学与技术、信息安全

指导教师：杨征路

2022 年 6 月 10 日

# 目录

<b>一、 概述</b>	<b>1</b>
(一) 实验目的 . . . . .	1
(二) 选择语言 . . . . .	1
<b>二、 数据集说明</b>	<b>1</b>
(一) 数据集说明 . . . . .	1
(二) 数据集基本统计 . . . . .	1
(三) 数据集处理 . . . . .	2
<b>三、 算法细节</b>	<b>3</b>
(一) 奇异值分解 (SVD) . . . . .	3
1. 生成转置训练数据集文本 . . . . .	3
2. 导入数据集并计算矩阵 . . . . .	4
3. 计算 RMSE . . . . .	7
4. 计算 test.txt 中所需项的得分 . . . . .	8
5. 保存矩阵数据 . . . . .	8
(二) 基于物品属性的线性回归 . . . . .	9
1. 导入 itemAttribute.txt 数据集并处理物品属性值矩阵 . . . . .	9
2. 导入训练数据集并生成线性回归模型 . . . . .	10
3. 计算 test.txt 中所需项的得分 . . . . .	12
(三) UV 分解与线性回归的结果合并 . . . . .	13
(四) 基于用户的协同过滤 (UserCF) . . . . .	15
1. 计算每个用户评分均值 . . . . .	15
2. 预测与构建相似度矩阵 . . . . .	15
3. 计算 RMSE . . . . .	18
<b>四、 实验结果</b>	<b>19</b>
(一) RMSE . . . . .	19
(二) 训练时间 . . . . .	19
(三) 空间消耗 . . . . .	19
<b>五、 算法分析</b>	<b>20</b>
(一) 理论分析 . . . . .	20
(二) 实验结果分析 . . . . .	20
<b>六、 总结</b>	<b>20</b>

## 一、概述

### (一) 实验目的

本次实验目的在于采用大数据课程上学习到的三种不同算法——基于内容、奇异值分解和协同过滤。

数据集采用的是课程提供的 data-202205 数据集。

为了正确的实现推荐系统需要考虑 RMSE 值、训练时间、空间消耗的问题。

最终得到 test.txt 文件要求的对应用户对相应物品项的评分值。

### (二) 选择语言

为了方便进行矩阵相关的操作，经过小组成员讨论后，选择 Python 作为开发语言。

## 二、数据集说明

### (一) 数据集说明

数据集为课程提供的 data-202205 数据集。

其中包括 train.txt、itemAttribute.txt 两个数据集用于训练模型以及需要给出预测得分项的 test.txt 文件，以及两个格式说明文件。

### (二) 数据集基本统计

数据集的基本统计

用户总数：19835

train.txt 评分数：5001507

test.txt 评分数：119010

物品总数：624961

### (三) 数据集处理

使用 train\_split.py 脚本对 train.txt 进行处理。将 train.txt 中的数据按照 9:1 的比例分为训练集 \_train 和测试集 \_test。

train.txt 对于每个用户拥有的评分条目, 将其按照比例分配到训练集和测试集当中, 其中训练集的数据占原本 train.txt 的 9/10, 而测试集的数据占占原本 train.txt 的 1/10。

使用的处理数据集脚本实现如下:

划分训练集、测试集脚本

```
1 if __name__ == '__main__':
2     train_file = open(r".\train.txt", encoding='utf-8', mode='r')
3     for_train = open(r".\_train.txt", encoding='utf-8', mode='w')
4     for_test = open(r".\_test.txt", encoding='utf-8', mode='w')
5     while True:
6         line=train_file.readline()
7         if line=='':
8             break
9         item_num=int(line.split('|')[1])
10        user_num=int(line.split('|')[0])
11        test_num=item_num//10
12        for_test.write("%d|%d\n" %(user_num, test_num))
13        for _ in range(test_num):
14            line=train_file.readline()
15            for_test.write(line)
16        train_num = item_num-test_num
17        for_train.write("%d|%d\n" %(user_num, train_num))
18        sort_list=[]
19        for _ in range(train_num):
20            line = train_file.readline()
21            item=int(line.split('|')[1])
22            # score=int(line.split('|')[2])
23            sort_list.append((item, line))
24        sort_list=sorted(sort_list, key=lambda x:x[0])
25        for ele in sort_list:
26            line=ele[1]
27            for_train.write(line)
28        train_file.close()
29        for_train.close()
30        for_test.close()
```

使用脚本后, 生成训练集 \_train.txt 和测试集文件 \_test.txt, 使得训练集 \_train.txt 内包含评分总数达到 4509815 个, 测试集文件 \_test.txt 内评分总数到达 491692 个。

## 三、 算法细节

### (一) 奇异值分解 (SVD)

所谓 SVD 分解,即理论上,可以将任何一个矩阵分解为三个矩阵相乘,分别为左奇异矩阵、右奇异矩阵和奇异值矩阵,其中奇异值矩阵只有主对角线上的值非零,其值表示原矩阵的奇异值。

在本次实验中,根据训练数据集可以得到一个用户-物品的评分矩阵,成为效用矩阵  $M$ 。由于并不是每一个用户都对每一个物品进行了评分,故矩阵  $M$  为稀疏矩阵,含有值的矩阵元素只占较小部分。

本实验中将 SVD 方法应用于对矩阵  $M$  进行主成分分析。完整的 SVD 算法是对完整的稠密矩阵进行特征分解,而本次实验中得到的效用矩阵  $M$  为稀疏矩阵,大量矩阵元素的实际值未知,无法精准得到矩阵  $M$  的奇异值并分解,本实验并不需要用到效用矩阵的特征值,只需要对稀疏矩阵的未知值进行估计,故本次实验采用的方法为 SVD 分解的一个特例,UV 分解:只用到左右两个矩阵  $m$  行  $d$  列的  $U$  和  $d$  行  $m$  列的  $V$ ,使得两个矩阵的乘积  $UV$  和  $M$  在  $M$  的空白元素上较为接近,即使用  $d$  维特征近似地去刻画用户和物品之间的评分关系。

设定 epoch 数目,每次训练遍历一次训练集文件,直接将文件作为稀疏矩阵,不再整个读入内存中,根据读取的每个用户实际做出的评分,对逐个  $U$  和  $V$  矩阵的元素进行调整,通过计算 RMSE 值衡量 UV 分解对效用矩阵的拟合程度的好坏,使用简易的梯度下降来调整矩阵的数值,多次重复轮询循环,使得最终 RMSE 的值小于阈值或两次轮询之间的减小量小于阈值即矩阵分解收敛停止训练。

#### 1. 生成转置训练数据集文本

为了简化计算,需要生成以物品编号为行,以对该物品打出过评分的用户编号及相应评分作为列的转置训练集文本,实验使用的方法是先使用脚本将原始的文本数据集遍历,获取不同的物品及相应的评分用户记录在内存,再从内存中逐行写入转置训练集文本文件中。其中为防止内存中需保存的物品为索引的评分记录占用空间过大,在转置过程中应用了分块思想,使用 7 次循环,每次循环记录 100000 个物品的评分记录,再写入文本文件,再进入下一个循环继续记录。

使用的脚本实现如下:

划分训练集、测试集脚本

```
1 if __name__ == '__main__':
2     for_train = open(r".\train.txt", encoding='utf-8', mode='r')
3     transpose = open(r".\train_trans.txt", encoding='utf-8', mode='w')
4     for i in range(7):
5         l = [[] for _ in range(100000)]
6         for_train.seek(0,0)
7         while True:
8             line = for_train.readline()
9             if line == '':
10                 break
11             item_num = int(line.split('|')[1])
12             user = int(line.split('|')[0])
13             print("current user:{}".format(user))
14             for _ in range(item_num):
15                 line=for_train.readline()
16                 item=int(line.split(' ')[0])
```

```

17         score=int(line.split(' ')[1])
18         if item in range(i*100000,(i+1)*100000):
19             l[item%100000].append((user,score))
20     for index in range(len(l)):
21         if i==6 and index>24960:
22             break
23         if len(l[index])!=0:
24             transpose.write("{}|{}\n".format(index+i*100000,len(l[index]))
25             )
26             for ele in l[index]:
27                 transpose.write("{} {}{}\n".format(ele[0],ele[1]))
28         else:
29             transpose.write("{}|{}\n".format(index+i*100000, 0))
30     for_train.close()
31     transpose.close()

```

使用脚本后,生成转置训练集 train\_trans.txt。

## 2. 导入数据集并计算矩阵

这一部分在每个 epoch 当中重复进行,直至完成预先设定的 epoch 数。在对划分的训练集训练 100 次迭代进行数据分析,为生成最终的 test 中的评分,对整个 train.txt 进行 25 次迭代。

将本地的数据集和转置后的数据集文本读入,并初始化 UV 矩阵,令维度 d 为 1;实际在模型训练测试中尝试过 d 取值为 2 或 3 等更大值,最终发现当 d 为 1 时 RMSE 值最小,对其理论解释为由于稀疏矩阵中大量实际分数值未知,在试图用 d 作为各个物品的属性特征时试图将特征数量最小化即为将各个物品的共有属性数量最少化能增强 UV 矩阵拟合能力,降低拟合误差。

读取效用矩阵中关于特定用户的所有评分的数据集合,只有这些实际被用户赋予了评分的物品的分数项才参与 RMSE 误差值的计算,于是计算迭代过程为每次读取该行用户的一个评分值,分别在 U 和 V 矩阵中选取其对应的行和列的向量进行点乘得到平方误差,得到该行用户对应的所有评分过的物品在 V 中的对应列取出点乘后得到的误差平方和的导数,令其等于 0,计算得到此时 U 中对应行的各元素的调整后值存入 U 矩阵中,在读下一用户行继续更新 U。

实现如下:

### 读入用户数索引数据集并迭代 U 矩阵

```

1 train = open(uv.train_path, mode='r')
2 train.seek(0, 0)
3 while True:
4     line = train.readline()
5     if line == '':
6         break
7     try:
8         user = int(line.split('|')[0])
9     except ValueError:
10        print(line)
11        print(user)
12        exit(0)
13    print("current u_row:{}".format(user))
14    rate_num = int(line.split('|')[1])

```

```
15     u_line = user
16     cur_pos = train.tell()
17     for u_col in range(uv.d):
18         up = 0
19         down = 0
20         for _ in range(rate_num):
21             line = train.readline()
22             item = int(line.split(' ')[0])
23             score = int(line.split(' ')[1])
24             up += uv.v[u_col, item] * (score - (
25                 np.dot(uv.u[u_line, :], uv.v[:, item]) - uv.u[u_line,
26                     u_col] * uv.v[u_col, item]))
27             down += uv.v[u_col, item] ** 2
28         if down != 0:
29             z_point = up / down
30             if z_point > uv.u[u_line, u_col]:
31                 uv.u[u_line, u_col] += abs(z_point - uv.u[u_line, u_col]) * 1
32             else:
33                 uv.u[u_line, u_col] -= abs(z_point - uv.u[u_line, u_col]) * 1
34         if u_col != uv.d - 1:
35             train.seek(cur_pos, 0)
36     train.close()
```

对 V 进行类似上所述 U 的导数为 0 更新迭代计算，实现如下：

读入物品索引数据集并迭代 V 矩阵

```

1 train_trans = open(uv.train_trans_path, mode='r')
2 train_trans.seek(0, 0)
3 while True:
4     line = train_trans.readline()
5     if line == '':
6         break
7     try:
8         v_col = int(line.split('|')[0])
9     except ValueError:
10        print(line)
11        print(v_col)
12        exit(0)
13    print("current v_col:{}".format(v_col))
14    user_num = int(line.split('|')[1])
15    up = [0 for _ in range(uv.d)]
16    down = [0 for _ in range(uv.d)]
17    user = 0
18    for _ in range(user_num):
19        line = train_trans.readline()
20        user = int(line.split(' ')[0])
21        score = int(line.split(' ')[1])
22        for v_line in range(uv.d):
23            up[v_line] += uv.u[user, v_line] * (score - (
24                np.dot(uv.u[user, :], uv.v[:, v_col]) - uv.u[user, v_line]
25                * uv.v[v_line, v_col]))
26            down[v_line] += uv.u[user, v_line] ** 2
27        for v_line in range(uv.d):
28            if down[v_line] == 0:
29                continue
30            z_point = up[v_line] / down[v_line]
31            if z_point > uv.v[v_line, v_col]:
32                uv.v[v_line, v_col] += abs(z_point - uv.v[v_line, v_col]) * 1
33            else:
34                uv.v[v_line, v_col] -= abs(z_point - uv.v[v_line, v_col]) * 1
35    train_trans.close()

```



### 3. 计算 RMSE

该部分使用 `calc_rmse` 函数进行计算均方根误差。

导入之前划分的测试集文件，每读入测试集文件中的一行评分项，就进行误差的计算。

误差计算完毕之后，返回对应的误差。

`calc_rmse` 函数实现如下：

计算 RMSE

```
1 def calc_rmse(self):
2     train = open(self.train_path, mode='r')
3     se = 0
4     count = 0
5     while True:
6         line = train.readline()
7         if line == '':
8             break
9         user = int(line.split('|')[0])
10        # print(user)
11        rate_num = int(line.split('|')[1])
12        for _ in range(rate_num):
13            line = train.readline()
14            item = int(line.split('|')[0])
15            score = int(line.split('|')[1])
16            error = (score - np.dot(self.u[user, :], self.v[:, item])) ** 2
17            se += error
18            # print("error:{}".format(error))
19        count += rate_num
20    rmse = (se / count) ** 0.5
21    train.close()
22    return rmse
```

#### 4. 计算 test.txt 中所需项的得分

该部分使用 display 函数进行计算得分。遍历 test.txt 测试文件，每遍历到对应用户对应物品项的一行，就进行一次估计值计算，写入结果文件中展示。

将 U 矩阵中所求用户对应的行与 V 矩阵对应所求物品项对应的列进行矩阵点乘，即可得到对应用户对应物品的预测得分，将该得分保存在最终文件当中。

display 函数实现如下：

计算预测得分

```

1  def display(self):
2      valid = open(self.test_path, mode='r')
3      res = open(self.res_path, mode='w')
4      while True:
5          line = valid.readline()
6          if line == '':
7              break
8          user = int(line.split('|')[0])
9          rate_num = int(line.split('|')[1])
10         res.write(line)
11         for _ in range(rate_num):
12             line = valid.readline()
13             item = int(line.split(' ')[0])
14             score = np.dot(self.u[user, :], self.v[:, item])
15             line = "{ } {} \n".format(item, score)
16             res.write(line)
17         valid.close()
18         res.close()

```

#### 5. 保存矩阵数据

为了保存矩阵训练后的结果，避免之后计算需要重新计算，就将矩阵保存在本地。

使用 save\_mat 的函数进行矩阵数据保存。

save\_mat 的函数实现如下：

保存矩阵信息

```

1  def save_mat(self):
2      save = open(self.save_path, mode='w')
3      for u_row in range(self.u_dim):
4          for u_col in range(self.d):
5              save.write("%f|" % self.u[u_row, u_col])
6          save.write("\n")
7      save.write("\n")
8      for v_row in range(self.d):
9          for v_col in range(self.v_dim):
10             save.write("%f|" % self.v[v_row, v_col])
11         save.write("\n")
12     save.write("\n")
13     save.close()

```

## (二) 基于物品属性的线性回归

相比上述使用 epoch 循环训练, 该方法使用线性回归来确定矩阵的数值。

有效利用了数据集中 itemAttribute.txt 的数据。

同样利用 \_\_train 数据集进行回归训练, 再利用 \_\_test 数据集测试估计分数计算与正确分数的误差。

### 1. 导入 itemAttribute.txt 数据集并处理物品属性值矩阵

打开训练集 itemAttribute.txt 文件, 遍历该文件的过程中设置物品属性值矩阵的值。

若 itemAttribute.txt 文件中的 Attribute 为 None, 则视为该物品的该属性的数值为 0。

由于视为物品属性值的数值的绝对值大小变化幅度较大, 故数据预处理过程中先对数据进行归一化处理, 方便后续训练, 并转置获得正确行列结构的矩阵。

导入 itemAttribute.txt 并处理物品属性矩阵, 实现如下:

导入 itemAttribute.txt 数据集并生成物品属性矩阵

```
1 LR = LinearRegressionForUsers()
2 feature = open(LR.feature_path, mode='r')
3 feature_mat = np.zeros((624961,2))
4 for line in feature:
5     item = int(line.split('|')[0])
6     attr1 = line.split('|')[1]
7     attr2 = line.split('|')[2].strip()
8     attr1 = int(attr1) if attr1 != 'None' else 0
9     attr2 = int(attr2) if attr2 != 'None' else 0
10    feature_mat[item,0]=attr1
11    feature_mat[item,1]=attr2
12    feature.close()
13    min = [np.min(feature_mat[... , i]) for i in range(feature_mat.shape[1])]
14    max = [np.max(feature_mat[... , i]) for i in range(feature_mat.shape[1])]
15    #print(feature_mat.shape[1])
16    feature_mat = np.array([(ele - min[i]) / (max[i] - min[i]) for ele in
17        feature_mat[... , i]] for i in range(feature_mat.shape[1]))
18    feature_mat=feature_mat.transpose()
19    #print(feature_mat)
```

## 2. 导入训练数据集并生成线性回归模型

处理完毕特征矩阵后，导入训练数据集。

开始进行线性回归计算，结合之前处理得到的物品属性矩阵，对每一个用户训练一个具有两个参数的线性回归模型，将生成好的模型保存入列表，用于后续对每个用户的未评分物品进行分数估计。

导入后，成功生成基于物品属性的二元线性回归模型。

导入训练数据集，实现如下：

导入训练数据集

```

1  train = open(LR.train_path, mode='r')
2  LR_set=[]
3  while True:
4      line = train.readline()
5      if line == '':
6          break
7      user = int(line.split('|')[0])
8      print("current user:{}".format(user))
9      rate_num = int(line.split('|')[1])
10     current_x=np.ones((rate_num,2))
11     current_y=np.ones(rate_num)
12     for i in range(rate_num):
13         line = train.readline()
14         item = int(line.split(' ')[0])
15         score = int(line.split(' ')[1])
16         current_y[i]=score
17         current_x[i,0]=feature_mat[item,0]
18         current_x[i, 1] = feature_mat[item, 1]
19     LR.fit_gd(current_x, current_y, eta=0.005)
20     LR_set.append(LR)
21     LR=LinearRegressionForUsers()
22     # if user==0:
23     #     break
24 train.close()

```

线性回归方法的实现如下，关键部分为梯度下降的实现，每次通过导数公式计算好梯度，将参数  $\theta$  参数减去梯度，重复迭代；当迭代次数超过预设次数或前后 MSE 的差值小于预设值时停止迭代，该用户的回归模型训练完成：

线性回归方法

```

1  def fit_normal(self, X_train, y_train):
2      assert X_train.shape[0] == y_train.shape[0], \
3          "the size of X_train must be equal to the size of y_train"
4      X_b = np.hstack([np.ones((len(X_train), 1)), X_train])
5      self._theta = np.linalg.inv(X_b.T.dot(X_b)).dot(X_b.T).dot(y_train)
6      self.intercept_ = self._theta[0]
7      self.coef_ = self._theta[1:]

```

```

8         return self
9
10    def fit_gd(self, X_train, y_train, eta=0.01, n_iters=1e4):
11        assert X_train.shape[0] == y_train.shape[0], \
12            "the size of X_train must be equal to the size of y_train"
13
14    def J(theta, X_b, y):
15        try:
16            return np.sum((y - X_b.dot(theta)) ** 2) / len(y)
17        except:
18            return float('inf')
19
20    def dJ(theta, X_b, y):
21        return X_b.T.dot(X_b.dot(theta) - y) * 2. / len(X_b)
22
23    def gradient_descent(X_b, y, initial_theta, eta, n_iters=1e4, epsilon
24                          =1e-8):
25        theta = initial_theta
26        cur_iter = 0
27        while cur_iter < n_iters:
28            gradient = dJ(theta, X_b, y)
29            last_theta = theta
30            theta = theta - eta * gradient
31            yPred = X_b.dot(theta)
32            num = (yPred - y).dot(yPred - y)
33            d = y.size
34            self.MSE.append(num / d)
35            self.theta_list.append((cur_iter + 1))
36            if (abs(J(theta, X_b, y) - J(last_theta, X_b, y)) < epsilon):
37                break
38            cur_iter += 1
39        return theta
40
41    X_b = np.hstack([np.ones((len(X_train), 1)), X_train])
42    initial_theta = np.zeros(X_b.shape[1])
43    self._theta = gradient_descent(X_b, y_train, initial_theta, eta,
44                                  n_iters)
45    self.intercept_ = self._theta[0]
46    self.coef_ = self._theta[1:]
47    return self

```

### 3. 计算 test.txt 中所需项的得分

该部分使用 predict 函数进行计算得分。

遍历 test.txt 文件，每遍历到对应用户对应物品项的一行，就进行一次对应的计算。

使用线性回归模型和特征矩阵计算，由于 itemAttribute 文件中存在不少物品的属性值条目缺失，故训练过程中相应的物品属性条目数值为 0 未对模型训练做出贡献，且在估分过程中直接将自变量为 0 将参数影响消除，故在最终估分过程中若分数结果为模型参数的截距值则将分数值直接标记为-1，后续综合多方法时去除该分数条目采用其它算法；将得分保存在最终文件当中。

导入 test.txt 文件并进行预测得分，实现如下：

导入 test.txt 文件并预测

```

1  valid = open(LR.validate_path, mode='r')
2  res=open(LR.res_path,mode='a')
3  # se=0
4  # count=0
5  while True:
6      line = valid.readline()
7      if line == '':
8          break
9      user = int(line.split('|')[0])
10     print("current test_user:{}".format(user))
11     rate_num = int(line.split('|')[1])
12     test_x = np.ones((rate_num,2))
13     # test_y = np.ones(rate_num)
14     LR=LR_set[user]
15     items=[]
16     for i in range(rate_num):
17         line = valid.readline()
18         item = int(line.split(' ')[0])
19         # score = int(line.split(' ')[1])
20         items.append(item)
21         # test_y[i]=score
22         test_x[i, 0] = feature_mat[item, 0]
23         test_x[i, 1] = feature_mat[item, 1]
24     y_predT = LR.predict(test_x)
25     # se += (y_predT - test_y).dot(y_predT - test_y)
26     # count+=rate_num
27     res.write("{}|{}\n".format(user,rate_num))
28     for ele0,ele1 in zip(y_predT,items):
29         if ele0==LR.intercept_:
30             ele0=-1
31             res.write("{} {}{}\n".format(ele1,ele0))
32     # print(LR.intercept_)
33     # if user==0:
34     #     break
35 valid.close()
36 res.close()

```

predict 函数实现如下:

predict 函数

```

1  def predict(self, X_predict):
2      assert self.intercept_ is not None and self.coef_ is not None, \
3          "must fit before predict!"
4      assert X_predict.shape[1] == len(self.coef_), \
5          "the feature number of X_predict must be equal to X_train"
6
7      X_b = np.hstack([np.ones((len(X_predict), 1)), X_predict])
8
9      return X_b.dot(self._theta)

```

### (三) UV 分解与线性回归的结果合并

将 UV 分解的结果和基于物品属性线性回归的结果按照规则: 若线性回归估计的分数值为-1, 则该分值采用 UV 分解模型得到的分值, 否则采用线性回归模型得到的分值; 即可得到使用 UV 分解并结合使用 itemAttribute.txt 文件的结果。

合并结果的过程中, 可以测算组合方案在验证集上的 RMSE 值以评估最终算法效果。

合并 UV 分解和线性回归的实现如下:

合并结果

```

1  uv_path = r'.\res-valid-100run.txt'
2  regre_path = r'.\res-valid-regre-0.txt'
3  test_path = r'.\_test.txt'
4  res_path = r'.\res-combine.txt'
5  res1 = open(uv_path, mode='r')
6  res2 = open(regre_path, mode='r')
7  res0=open(res_path,mode='w')
8  res3=open(test_path,mode='r')
9  se=0
10 count=0
11 while True:
12     line1=res1.readline()
13     line2=res2.readline()
14     line3=res3.readline()
15     if line1=='':
16         break
17     if line1.find('|')!=-1:
18         res0.write(line1)
19     else:
20         item = int(line1.split(' ')[0])
21         score1 = float(line1.split(' ')[1].strip())
22         score2=float(line2.split(' ')[1].strip())
23         if score2==-1.0:

```

```
24         if score1 < 0:
25             score1 = 0.0
26         if score1 > 100:
27             score1 = 100.0
28         score = score1
29     else:
30         if score2 < 0:
31             score2 = 0.0
32         if score2 > 100:
33             score2 = 100.0
34         score = score2
35     res0.write("{} {} \n".format(item, score))
36     real_score = int(line3.split(' ')[1])
37     se += (real_score - score) ** 2
38     count += 1
39 rmse = (se / count) ** 0.5
40 print(rmse)
41 res0.close()
42 res1.close()
43 res2.close()
44 res3.close()
```



## (四) 基于用户的协同过滤 (UserCF)

由于用户数为 19835, 而项目数为 624961。项目的数量级远大于用户的数量级, 所以采取基于用户的协同过滤能有效减少计算量。

基于用户的协同过滤需要构造用户相似度矩阵。要预测某个用户对某项物品的打分, 可以根据用户相似度矩阵找出数个与该用户相似的其他用户。根据相似的其他用户对该物品的打分来预测所求用户的打分。

使用余弦相似度来评估项之间的相似性。

### 1. 计算每个用户评分均值

由于用户打分的习惯问题, 打分可能存在统一偏高或偏低的问题, 先计算平均值, 再将评分减去平均值, 更能反映用户真实的喜好情况。

将本地的训练数据集导入, 计算每个用户的评分均值。

该过程由 calcMean 函数实现。

calcMean 函数, 实现如下:

计算每个用户评分均值

```
1  def calcMean(self):
2      '计算所有打分的均值'
3      dataFile = open(self.dataFilePath, 'r')
4      lines = dataFile.readlines()
5      lineNum = 0
6      while lineNum < len(lines):
7          self.userNum += 1
8          #print(lines[lineNum].strip().split())
9          uid, count = lines[lineNum].strip().split('|')
10         temp = dict()
11         for i in range(int(count)):
12             item = lines[i + lineNum + 1].strip().split()
13             temp[item[0]] = int(item[1])
14         lineNum += (int(count) + 1)
15         self.rateNum += len(temp)
16         self.rateSum += sum(list(temp.values()))
17         self.userRateVectors.append(temp)
18         self.rateMean = self.rateSum / self.rateNum
19         dataFile.close()
20         #print(self.userNum)
```

### 2. 预测与构建相似度矩阵

如果需要预测某个用户 (预测用户) 对某个物品, 首先寻找拥有对该物品打分过的用户 (参考用户), 若预测用户与参考用户之间没有计算过相似度, 则使用余弦相似度公式进行计算, 并保存到相似度矩阵中。如果已经计算过相似度则不重复计算。

由于用户与用户之间的相似度可以用对角矩阵保存, 从而实现进一步节省空间消耗。

寻找出设定的 N 个最相似用户, 然后根据这些参考用户的评分进行预测。

预测过程中，遍历 test.txt 文件，每遍历到对应用户对应物品项的一行，就进行一次对应的预测计算。

该过程由 predict 函数实现。

predict 函数，实现如下：

#### 预测用户评分

```

1  def predict(self):
2      '对_test.txt中item进行打分预测'
3      testFile = open(self.testFilePath, 'r')
4      resFile = open(self.resultFilePath, 'w+')
5      data = testFile.readline()
6      time = 0
7      while(data!=''):
8          time+=1
9          resFile.write(data)
10         uid,count = data.strip().split('|')
11         userVector = self.getUserVector(uid)
12         for i in range(int(count)):
13             item = testFile.readline().strip().split(' ')[0]
14             #output itemid to resuslt file
15             print('to predict:',uid,item)
16             simList = self.getSimUsers(uid,userVector,item)
17             if(self.allSame):
18                 rate = self.sameRate
19                 self.allSame=False
20             else:
21                 weightsum,Sum=0,0
22                 for i in range(len(simList)):
23                     Sum+=simList[i][0]*simList[i][1]
24                     weightsum+=simList[i][0]
25                     #rate =self.rateMean+self.userDeviation+self.
26                     itemDeviation+ Sum/weightsum
27                     rate = Sum/weightsum
28                 if rate>100:
29                     rate = 100
30                 if rate<=0:
31                     rate = 0
32                 resFile.write(item+' '+str(rate)+'\n')
33             data = testFile.readline()
34         testFile.close()
35         resFile.close()

```

获取与预测用户最为相似的参考用户集合的过程由 `getSimUsers` 函数实现。  
`getSimUsers` 函数，实现如下：

#### 获取参考用户集合

```

1  def getSimUsers(self, userid, userVector, itemID):
2      '获取与uID最为相似的user集合N'
3      res = list()
4      Sum, num = 0, 0
5      for i in range(len(self.userRateVectors)):
6          if itemID in self.userRateVectors[i].keys():
7              Sum += self.userRateVectors[i][itemID]
8              num += 1
9              itemTemp = self.userRateVectors[i].copy()
10             #归一化
11             m = sum(list(itemTemp.values())) / len(itemTemp)
12             for j in itemTemp.keys():
13                 itemTemp[j] -= m
14             #print(temp)
15             userTemp = userVector.copy()
16             #归一化
17             m = sum(list(userTemp.values())) / len(userTemp)
18             self.userDeviation = m - self.rateMean
19             for j in userTemp.keys():
20                 userTemp[j] -= m
21             if sum(itemTemp.values()) != 0 and sum(userTemp.values()) != 0:
22                 s = self.sim(userTemp, itemTemp, int(userid), int(i))
23                 res.append((s, self.userRateVectors[i][itemID]))
24                 res.sort(key=takeFirst, reverse=True)
25                 if len(res) > self.N:
26                     res.pop()
27             elif sum(userTemp.values()) == 0:
28                 self.allSame = True
29                 self.sameRate = sum(userVector.values()) / len(userVector)
30         try:
31             self.itemDeviation = Sum / num - self.rateMean
32         except ZeroDivisionError:
33             print('Sum:', Sum)
34             print('Num', num)
35
36     return res

```

### 3. 计算 RMSE

该部分使用 calcRMSE 函数进行计算均方根误差。

导入之前划分的测试集文件，每读入测试集文件中的一行评分项，就进行误差的计算。

误差计算完毕之后，返回对应的误差。

calcRMSE 函数实现如下：

计算 RMSE

```
1  def calcRMSE(self):
2      Sum =0
3      num =0
4      resFile = open(self.resultFilePath, 'r')
5      testFile = open(self.testFilePath, 'r')
6      data1 = testFile.readline()
7      data2 = resFile.readline()
8      while(data1!='' and data2!=''):
9          count = data1.strip().split('|')[1]
10         num+=int(count)
11         for i in range(int(count)):
12             item1=testFile.readline().strip().split()[1]
13             item2=resFile.readline().strip().split()[1]
14             Sum+=(float(item2)-float(item1))**2
15             data1=testFile.readline()
16             data2 = resFile.readline()
17         rmse = math.sqrt(Sum/num)
18         resFile.close()
19         testFile.close()
20         return rmse
```

## 四、实验结果

### (一) RMSE

算法	RMSE(验证集上)
奇异值分解	31.28
线性回归	29.89
协同过滤	43.02
奇异值分解 + 线性回归	28.97

表 1: 不同算法的 RMSE 对比表

### (二) 训练时间

算法	运行时间
奇异值分解 (epoch=100)	2.5
线性回归	1.3
协同过滤	7.2
奇异值分解 + 线性回归	3.8

表 2: 不同算法在训练集上的训练时间对比表 (单位: 小时)

### (三) 空间消耗

算法	消耗空间
奇异值分解	$O(N^2 + M^2)$
线性回归	$O(2N)$
协同过滤	$O(N^2)$
奇异值分解 + 线性回归	$O(N^2 + M^2)$

表 3: 不同算法的空间消耗对比表

## 五、 算法分析

### (一) 理论分析

设用户数为  $N$ ，物品数为  $M$ 。

使用 SVD 算法，需要保存  $U$ 、 $V$  两个矩阵，两个矩阵的空间复杂度分别为  $O(N^2)$  和  $O(M^2)$ 。所以总的空间复杂度为  $O(N^2 + M^2)$ 。

而使用基于物品属性的线性回归算法，需要保存每个用户的两个属性，空间复杂度为  $O(2N)$ 。

使用下三角矩阵保存协同过滤的相似度矩阵，若是基于用户，则需要保存的矩阵的空间复杂度为  $O(N^2)$ ；若是基于物品，则需要保存的矩阵的空间复杂度为  $O(M^2)$ 。显然，由于课程提供的数据集用户数  $N$  远小于物品数  $M$ ，使用基于用户的协同过滤，空间复杂度为  $O(N^2)$  远小于基于物品的协同过滤的空间复杂度  $O(M^2)$ 。

基于用户的协同过滤算法计算用户相似度矩阵需要  $O((N-1)M)$  的时间复杂度。而其他两个算法涉及机器学习，训练所需要的时间与参数的设置以及收敛的速度有关。

### (二) 实验结果分析

由于协同过滤在验证集上的 RMSE 误差值较大，本次实验最终采取的是奇异值分解与线性回归相结合的推荐系统算法。

消耗时间为 25 次迭代的 UV 分解模型的 45mins，加上学习率为 0.005 的线性回归模型 1h 20mins 共为 2h 5mins。

空间消耗为  $O(N^2 + M^2)$ 。

表现较好。将其对 test.txt 中预测评分项的预测结果作为最终答案。

## 六、 总结

通过本次实验，小组成员深刻了解了奇异值分解、线性回归、协同过滤应用于推荐系统的原理，学会了计算余弦相似度。

尝试使用了 itemAttribute.txt 中的 Attribute 提高算法的性能。

小组成员学习了基于用户和基于物品的推荐系统两种推荐方式的特点和适用场景。小组成员学会了实现基于物品属性的线性回归和基于用户的协同过滤的方法。

本次实验为小组成员将来的学习奠定了坚实的基础。