

# Web信息处理与应用实验报告

## lab3-推荐系统实践

曾勇程 PB18000268

张弛 PB18000269

### Web信息处理与应用实验报告

#### lab3-推荐系统实践

实验目的

实验内容

实验条件

关键代码讲解

基于用户的协同过滤

SVD算法

社交关系的思考

实验结果分析

所作优化

## 实验目的

本实验要求以给定的豆瓣电影打分数据集为基础，实现一个电影推荐系统。

该系统的主要功能是从给定的2173个用户对58431个电影所打出的227万多个评分中学习或训练出每个用户的观影偏好打分偏好、每部电影的隐含特征，从而对于测试集中给定的目标用户和电影给出预测的评分值，便于我们向用户推荐影片。

## 实验内容

数据来源为豆瓣电影的评分记录，包含用户ID，电影ID，打分，打分时间，电影标签。需要运用个性化推荐技术。

本次实验中我们小组实现了基于**内容**和基于**模型**的协同过滤推荐，即**用户协同过滤**和**SVD**算法。我们的工作可以分为以下几个步骤：

1. 对文档数据集进行**预处理**。这一步骤包括逐行读取文档内容，获得用户和电影信息，并进行重新编码处理，将打分数据**存入矩阵**数组中。注意，为了区分评分为0的和用户未打分两种情况，录入评分时我们**对所有评分加1**，得以和矩阵中未评分项0区分。
2. 实现**基于用户的协同过滤算法**。首先根据矩阵计算余弦相似度。具体计算方式即每个评分减去该用户的平均评分后，计算余弦相似度。需要注意的是，用户和电影的平均得分都**只考虑矩阵中的非0项**。然后对于每一个用户，搜索其**K个最近邻**的其他对目标电影评分的用户，以相似度**加权偏置**打分再加上用户偏置评分计算预测评分。注意，这里我们只采用相似度为正的近邻用户。为了用对协同过滤方法中**矩阵稀疏**的问题，我们尝试了先通过电影过滤进行**数据补足**再用户过滤的方法，即**混合协同过滤**。
3. 实现**SVD矩阵分解算法**。我们通过特定公式来预测打分，并通过不断地**梯度下降**学习，得到分解的两个用于打分的**特征矩阵**。
4. 思考并使用更多的**特征信息**。我们对于给定的**社交网络**，做出了各种尝试，以提升系统的准确性，如**冷启动**和**参与加权评分**。

## 实验条件

1. 硬件条件：一台PC机。
2. 软件条件：
  - 系统Windows10
  - 语言python3.6
  - 编辑器pycharm2018.3.4

## 关键代码讲解

本次实验用到的库有：

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.metrics.pairwise import cosine_similarity
```

### 1. 读取数据并对用户和电影编码

```
1 trainDF = pd.read_csv("data/training.dat", sep=",", header=None,
2     names=["UserID", "MovieID", "Rating", "Timestamp", "Tag1"])
3 # 可以继续加上Tag2, ...
4 testDF = pd.read_csv("data/testing.dat", sep=",", header=None,
5     names=["UserID", "MovieID", "Timestamp", "Rating"])
6
7 Num2UserID = dict(enumerate(list(trainDF["UserID"].unique())))
8 UserID2Num = {key: value for value, key in Num2UserID.items()}
9 UserNum = len(UserID2Num) # 2173
10
11 Num2MovieID = dict(enumerate(list(trainDF["MovieID"].unique())))
12 MovieID2Num = {key: value for value, key in Num2MovieID.items()}
13 MovieNum = len(MovieID2Num) # 58431
```

### 2. 建立评分矩阵

注意，在矩阵数组中，我们把所有用户对电影的评分+1，这样得以把评分为0的项和未评分项区分开。

```
1 DataMatrix = np.zeros((UserNum, MovieNum))
2 for line in trainDF.itertuples():
3     DataMatrix[UserID2Num[line[1]], MovieID2Num[line[2]]] = line[3]+1
4 # 将评分为0项与矩阵未评分项区分
```

同时，我们也计算出每个电影和用户的平均得分，这里只统计矩阵中非零项的均值。

```
1 DataExist = (DataMatrix != 0)
2 MovieMeanScore = DataMatrix.sum(axis=0) / DataExist.sum(axis=0)
3 # print(MovieMeanScore)
4 UserBias = DataMatrix.sum(axis=1) / DataExist.sum(axis=1)
```

### 3. 获取社交关系

我们将社交关系存储在字典中，每个用户ID对应的关注账号以集合存储。

```

1 SocialRelation = {}
2 readrel = open("data/relation.txt", "r")
3 attention = readrel.readline()
4 while attention:
5     attention = attention.strip("\n")
6     userid = int(attention.split(":", 1)[0])
7     SocialRelation[userid] = set(map(int, attention.split(":", 1)
8     [1].split(",")))
9     attention = readrel.readline()
9 readrel.close()
10 # print(SocialRelation)

```

## 基于用户的协同过滤

### 4. 计算用户间的余弦相似度

先对每个评分减去用户偏置分数，同样，代表未评分的0项不做处理。

```

1 BiasMatrix = DataExist * UserBias[:, np.newaxis]
2 # print(BiasMatrix)
3 NormMatrix = DataMatrix - BiasMatrix
4 UserSimilarity = cosine_similarity(NormMatrix)
5 # print(UserSimilarity)

```

### 5. 建立用户协同过滤评分预测系统

这一步中，我们不断调整近邻数目K的值，以找到预测效果最佳的模型。

以下是基础预测模型函数，我们对所有相似度为正的已评分用户加权平均。

```

1 def UserCFPredict(userid, movieid):
2     useri = UserID2Num[userid]
3     movie = MovieID2Num[movieid]
4     weight = 0
5     rating = 0
6     # no k-nearest
7     for userj in range(UserNum):
8         if (DataMatrix[userj, movie] != 0) and (UserSimilarity[useri, userj]
9         > 0):
10             weight += UserSimilarity[useri, userj]
11             rating += UserSimilarity[useri, userj]*NormMatrix[userj, movie]
12     if weight == 0:
13         return UserBias[useri]
14     return UserBias[useri] + rating/weight

```

以下是K近邻改进后的系统。经过尝试，我们最终认为K=50时效果较好。

```

1 def UserCFPredict(userid, movieid):
2     useri = UserID2Num[userid]
3     movie = MovieID2Num[movieid]
4     weight = 0
5     rating = 0
6     # k-nearest
7     K = 25
8     TopK = []
9     for index in range(K):

```

```

10     TopK.append([0, 0])      # userid, similarity
11     for userj in range(UserNum):
12         if DataMatrix[userj, movie] != 0:
13             ptr = K-1
14             if UserSimilarity[useri, userj] > TopK[ptr][1]:
15                 ptr = ptr-1
16                 while (ptr >= 0) and (UserSimilarity[useri, userj] >
TopK[ptr][1]):
17                     TopK[ptr+1] = TopK[ptr]
18                     ptr = ptr-1
19                     TopK[ptr+1] = [userj, UserSimilarity[useri, userj]]
20     for index in range(K):
21         if TopK[index][1] > 0:
22             weight += TopK[index][1]
23             rating += TopK[index][1]*NormMatrix[TopK[index][0], movie]
24     if weight == 0:
25         return UserBias[useri]
26     return UserBias[useri] + rating/weight

```

## 6. 生成预测数据

对于预测分数，我们四舍五入化为整数后还要减1，同时将过大和过小的数据修正回[0, 5]范围内。

```

1  def PrintResult(filename):
2      writeres = open(filename, "w")
3      for line in testDF.itertuples():
4          userid = line[1]
5          movieid = line[2]
6          predscore = round(UserCFPredict(userid, movieid)-1)
7          if predscore < 0:
8              predscore = 0
9          elif predscore > 5:
10             predscore = 5
11             writeres.write(str(predscore) + "\n")
12     writeres.close()
13     filename = "submission21.txt"
14     # PrintResult(filename)

```

## 7. 对于用户协同过滤的改进

用户协同过滤方法效果往往并不十分准确的主要原因在于评分矩阵过于稀疏，因此我们根据找到的论文中提到的方案，试图先通过协同过滤进行数据补足，以提升预测准确率。

```

1  def DataSupplement(DataMatrix):
2      nDataMatrix = np.zeros((UserNum, MovieNum))
3      for user in range(UserNum):
4          for movie in range(MovieNum):
5              if DataMatrix[user, movie] != 0:
6                  nDataMatrix[user, movie] = DataMatrix[user, movie]
7              else:
8                  userid = Num2UserID[user]
9                  movieid = Num2MovieID[movie]
10                 nDataMatrix[user, movie] = UserCFPredict(userid, movieid)
11     return nDataMatrix
12     # nDataMatrix = DataSupplement(DataMatrix)

```

不过可惜的是，由于用户和电影数量较多，该方案的计算量过大，在我们的电脑中跑不太动，因此我们放弃了这一思路。

## SVD算法

### 8. 数据处理

将训练集和测试集转为numpy数组，同样，所有打分加上1。

```
1 def getData():    #获取训练集和测试集的函数
2     train_data = []
3     test_data = []
4     for line in trainDF.itertuples():
5         train_data.append([UserIDNum[line[1]], MovieID2Num[line[2]],
6                             line[3]+1])
7     for line in testDF.itertuples():
8         test_data.append([UserIDNum[line[1]], MovieID2Num[line[2]]])
9     train_data=np.array(train_data)
10    test_data=np.array(test_data)
11    print('load data finished')
12    print('train data ',len(train_data))
13    print('test data ',len(test_data))
14    return train_data,test_data
15 train_data,test_data=getData()
```

### 9. SVD迭代计算预测两个向量

具体详解见注释：

```
1 class SVD:
2     def __init__(self,mat,K=30):
3         self.mat=np.array(mat)
4         self.K=K
5         self.bi={}
6         self.bu={}
7         self.qi={}
8         self.pu={}
9         self.avg=np.mean(self.mat[:,2])
10        for i in range(self.mat.shape[0]):
11            uid=self.mat[i,0]
12            iid=self.mat[i,1]
13            self.bi.setdefault(iid,0)
14            self.bu.setdefault(uid,0)
15
16        self.qi.setdefault(iid,np.random.random((self.K,1))/10*np.sqrt(self.K))
17
18        self.pu.setdefault(uid,np.random.random((self.K,1))/10*np.sqrt(self.K))
19
20    def predict(self,uid,iid):    #预测评分的函数
21        #setdefault的作用是当该用户或者物品未出现过时，新建它的bi,bu,qi,pu，并设置初
22        始值为0
23        self.bi.setdefault(iid,0)
24        self.bu.setdefault(uid,0)
25        self.qi.setdefault(iid,np.zeros((self.K,1)))
26        self.pu.setdefault(uid,np.zeros((self.K,1)))
```

```

24     rating=self.avg+self.bi[iid]+self.bu[uid]+np.sum(self.qi[iid]*self.pu[uid])
    #预测评分公式
25     #由于评分范围在1到6，所以当分数大于6或小于1时，返回6,1.
26     if rating>6:
27         rating=6
28     if rating<1:
29         rating=1
30     return rating
31
32     def train(self,steps=200,gamma=0.035,Lambda=0.15):    #训练函数，step为迭代
    次数。
33         print('train data size',self.mat.shape)
34         for step in range(steps):
35             print('step',step+1,'is running')
36             KK=np.random.permutation(self.mat.shape[0]) #随机梯度下降算法，kk为
    对矩阵进行随机洗牌
37             rmse=0.0
38             for i in range(self.mat.shape[0]):
39                 j=KK[i]
40                 uid=self.mat[j,0]
41                 iid=self.mat[j,1]
42                 rating=self.mat[j,2]
43                 eui=rating-self.predict(uid, iid)
44                 rmse+=eui**2
45                 self.bu[uid]+=gamma*(eui-Lambda*self.bu[uid])
46                 self.bi[iid]+=gamma*(eui-Lambda*self.bi[iid])
47                 tmp=self.qi[iid]
48                 self.qi[iid]+=gamma*(eui*self.pu[uid]-Lambda*self.qi[iid])
49                 self.pu[uid]+=gamma*(eui*tmp-Lambda*self.pu[uid])
50             gamma=0.95*gamma
51             print('rmse is',np.sqrt(rmse/self.mat.shape[0]))
52
53     def test(self,test_data):    #gamma以0.95的学习率递减
54         print('test data size',test_data.shape)
55         writeres = open(filename, "w")
56         for i in range(test_data.shape[0]):
57             uid=test_data[i,0]
58             iid=test_data[i,1]
59             predsore = round(self.predict(uid, iid)-1)
60             writeres.write(str(predsore) + "\n")
61         writeres.close()
62
63 a=SVD(train_data,40)
64 a.train()
65 a.test(test_data)

```

## 社交关系的思考

我们对社交关系的使用做了以下两种尝试：

冷启动，对于训练集中未参与评分的用户，预测其评分结果时我们引入社交关系，以其关注的账号的平均打分预测。

```

1 def ColdStart(userid, movieid):
2     if userid in SocialRelation:    # userid 在 社会关系中
3         count = 0

```

```
4 ratingsum = 0
5 for leaderid in SocialRelation[user_id]:
6     if DataMatrix[UserID2Num[leaderid], MovieID2Num[movieid]] != 0:
7         ratingsum += DataMatrix[UserID2Num[leaderid],
MovieID2Num[movieid]]
8         count += 1
9     if count != 0:
10        return ratingsum / count
11    else:
12        return MovieMeanScore[MovieID2Num[movieid]]
13 else:
14    return MovieMeanScore[MovieID2Num[movieid]]
```

以及用于**加权预测**，即不再遍历所有其他用户，而在所关注账号中使用相似度加权。不过该方案效果并不好，我们认为这是由于社交关系并不代表用户具有相同的观影偏好，因此不能默认其具有高相似性。

## 实验结果分析

### 用户协同过滤结果（RMSE）：

submission4.txt	1.4895220072226922
submission5.txt	1.5955736001816962
submission6.txt	1.5955736001816962
submission7.txt	1.5180811391737694
submission8.txt	1.497461122796132
submission9.txt	1.489788225339001
submission10.txt	1.4872986671105892
submission11.txt	1.4879806044832304
submission12.txt	1.5180811391737694

基于用户协同过滤的最佳结果未RMSE=1.4872，这一方案中主要通过改变近邻数量K的值来调试，最终取K=50。

### SVD预测结果（RMSE）：

submission16.txt	2.1228008717843587
submission17.txt	1.4359913667046194
submission18.txt	1.4351962339423872
submission19.txt	1.4352410418639572
submission20.txt	1.4328575589011638
submission21.txt	1.4320135332906845

SVD中最佳效果RMSE=1.4320，主要通过改变特征向量长度、学习率、迭代次数等改进，最终我们尝试的取值为长度40，迭代200次。

注：每次的RMSE都不大相同，该结果是我们经过很多次训练以及调参以后出现的某次最优结果，因为我们采用的是**随机梯度下降**，且迭代次数由于运行时间问题并没有达到饱和。

## 所作优化

1. 只考虑相似度大于0的项，并采用K近邻计算相似度；
2. 采用两种方法预测（基于**内容**和基于**模型**的协同过滤推荐，即**用户协同过滤**和**SVD**算法）；
3. 调整**SVD**算法的参数，使其达到近似最优。