1. 引言

假如你经营着一家网店,里面卖各种商品(Items),有很多用户在你的店里面买过东西,并对买过的 Items进行了评分,我们称之为历史信息,现在为了提高销售量,必须主动向用户推销产品,所以关键 是要判断出用户除了已经买过的商品之外还会喜欢哪些商品,这就需要利用用户购买商品过程产生的 历史信息。协同过滤通常分为基于用户的协同过滤和基于商品的协同过滤。

- 基于用户的协同过滤: 利用用户之间的相似度进行推荐
- 基于物品的协同过滤: 利用物品之间的相似度进行推荐

2. 原理

关于协同过滤的原理网上到处都有,思想很简单,这里就不赘述,下面举一个简单的实例来说明基于 用户的协同过滤:

```
1 5 3 0 4 0 0 1 2 4 4
2 3 1 2 0 0 2 0 0 1 2
3 4 0 1 2 1 0 1 2 4 0
4 0 2 3 0 0 4 2 1 1 2
5 1 3 4 1 2 0 0 0 0 2
6 2 0 0 0 2 3 5 1 4 0
7 4 1 1 0 2 0 1 0 0 3
8 1 2 0 3 0 0 2 1 0 3
9 3 1 0 3 0 4 0 0 1 0
10 2 1 3 1 0 1 0 0 2 1
```

上面每一行代表一个用户,每一列代表一个商品,比如第2行第一列的3表示用户2对商品1的评分为3,0代表对应的用户还没有购买过该商品,现在想预测用户2对商品4的评分:

- 找出对商品4评过分的用户: 用户1, 3, 5, 8, 9, 10, 评分分别为: 4, 2, 1, 3, 3, 1
- 分别计算用户2与用户1, 3, 5, 8, 9, 10之间的相似度,相似度的计算方法有很多,常用的分为3类:欧氏距离,余弦相似度,皮尔逊相关系数,网上很容易查到,这里以常用的余弦相关系数说明:

要计算用户2与用户1之间的相似度,首先找到二者都评过分的商品为:商品1,2,9,10,用户1对这4个商品的评分向量为r1=[5344],用户2对这4个商品评分向量为r2=[3112];所谓余弦相似度就是利用两个向量之间夹角的余弦值来衡量两个向量之间的相似度,显然夹角越小,余弦值就越大,两个向量就越靠近,即二者越相似,于是用户2和用户1之间的相似度就为sim2_1=(5*3+3*1+4*1+4*2)/(||r1||*||r2||)=0.953,其中||r||代表向量r的模长或者2范数,类似地分别计算出用户2与用户358910之间的sim23,sim25,sim28,sim29,sim210

- 最后利用相似度加权得到用户2对商品4的预测评分: predict = 4*sim2_1 + 2*sim2_3 + 1*sim2_5 + 3*sim2_8 + 3*sim2_9 + 1*sim2_10
- 基于物品相似度就是与上面计算过程几乎相似,只是计算的是物品之间的相似度

3. 实现

关于Matlab的实现可以参考: http://blog.csdn.net/google19890102/article/details/28112091, 这里我用C++实现,并用movielens.rar进行测试,这个数据集是包括训练集和测试集,已经处理成矩阵形式。

• 首先给出读取训练数据和保存预测结果的头文件

```
#ifndef LOAD H
#define LOAD_H
#include <iostream>
#include <fstream>
#include <vector>
#include <string>
using namespace std;
template <typename T>
vector<vector<T> > txtRead(string FilePath, int row, int col)
   ifstream input(FilePath);
   if (!input.is_open())
        cerr << "File is not existing, check the path: \n" << FilePath << endl;
        exit(1);
    vector<vector<T> > data(row, vector<T>(col, 0));
    for (int i = 0; i < row; ++i)
        for (int j = 0; j < col; ++j)
            input >> data[i][j];
    return data;
template<typename T>
void txtWrite(vector<vector<T> > Matrix, string dest)
    ofstream output (dest);
   vector<vector<T> >::size_type row = Matrix.size();
   vector(T)::size_type col = Matrix[0].size();
    for (vector<vector<T> >::size_type i = 0; i < row; ++i)</pre>
        for (\text{vector}(T)::\text{size\_type } j = 0; j < \text{col}; ++j)
```

```
output << Matrix[i][j];
}
output << endl;
}
#endif</pre>
```

• 再给出评价预测好坏的计算RMSE的头文件

```
1 #ifndef EVALUATE H
 2 #define EVALUATE_H
 3 #include <cmath>
 4 #include <vector>
 6 double ComputeRMSE(vector<vector<double> > predict, vector<vector<double> > test)
7 {
       int Counter = 0;
       double sum = 0;
       for (vector<vector<double> >::size_type i = 0; i < test.size(); ++i)</pre>
11
           for (vector \( \double \)::size_type j = 0; j \( \text[0].size(); ++j) \)
12
13
                if (predict[i][j] && test[i][j])
14
               {
15
                    ++Counter;
                    sum += pow((test[i][j] - predict[i][j]), 2);
17
18
               }
           }
19
20
       return sqrt(sum / Counter);
21
22 }
23
24 #endif
```

最后给出主函数:

```
1 #include "load.h"
2 #include "evaluate.h"
3 #include <vector>
4 #include <string>
```

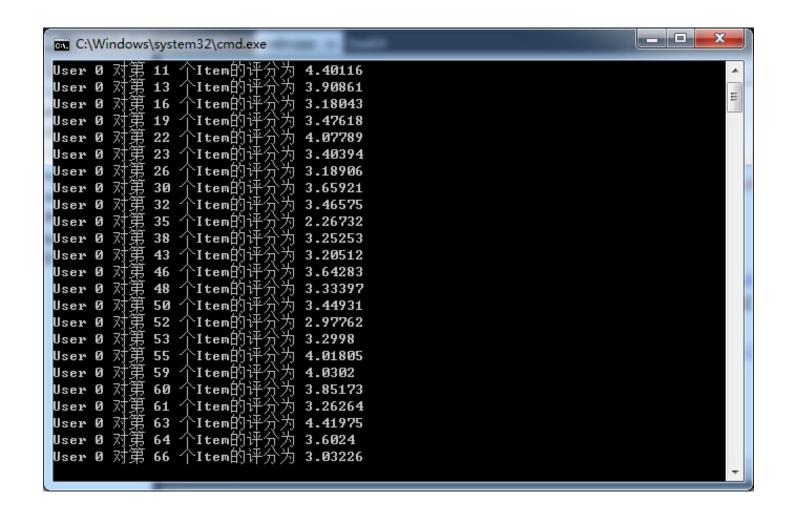
```
5 #include <cmath>
 6 #include <assert.h>
 7 using namespace std;
 9 double norm(vector \( \double \rangle \) A)
10 {
        double res = 0;
11
        for(vector \( \double \) :: size_type i = 0; i < A. size(); ++i)</pre>
12
13
            res += pow(A[i], 2);
14
15
       return sqrt(res);
16
17 }
18
19 double InnerProduct(vector(double) A, vector(double) B)
20 {
21
        double res = 0;
        for(vector \( \double \) :: size_type i = 0; i < A. size(); ++i)</pre>
22
23
            res += A[i] * B[i];
24
25
        return res;
26
27 }
28
29 double ComputeSim(vector double A, vector double B, int method)
30 {
31
        switch (method)
32
        case 0://欧氏距离
33
            {
34
                 vector <double > C;
35
                 for(vector \( \double \) :: size_type i = 0; i \( A. \size(); ++i) \)
36
37
                      C. push_back((A[i] - B[i]));
38
39
                 return 1 / (1 + norm(C));
40
                 break;
41
42
            }
        case 1://皮尔逊相关系数
43
44
                 double A_mean = 0;
45
                 double B_{mean} = 0;
46
                 for(vector(double)::size_type i = 0; i < A. size(); ++i)</pre>
47
                {
48
                      A_mean += A[i];
49
```

```
B_mean += B[i];
 50
 51
                  A_mean /= A.size();
 52
                  B_mean /= B.size();
 53
                  vector ⟨double⟩ C(A);
 54
                  vector < double > D(B);
 55
                  for(vector<double>::size_type i = 0; i < A.size(); ++i)</pre>
 56
 57
                       C[i] = A[i] - A_mean;
 58
                       D[i] = B[i] - B mean;
 59
 60
                  assert(norm(C) * norm(D));
 61
                  return InnerProduct(C, D) / (norm(C) * norm(D));
 62
                  break;
 63
 64
         case 2:
 65
             {
 66
                  assert(norm(A) * norm(B));
 67
                  return InnerProduct(A, B) / (norm(A) * norm(B));
 68
                  break;
 69
             }
 70
         default:
 71
 72
             {
                  cout << " Choose method:" << endl;</pre>
 73
                  cout << "0:欧氏距离\n1:皮尔逊相关系数\n2:余弦相似度\n";
 74
                  return -1;
 75
 76
 77
 78
 79 }
 80
 81 void FindCommon(vector \double \rangle A, vector \double \rangle B, vector \double \rangle &C, vector \double \rangle
&D)
 82 {
 83
         for(vector \( \double \)::size_type i = 0; i < A. size(); ++i)</pre>
 84
             if (A[i] && B[i])
 85
 86
                 C.push back(A[i]);
 87
                 D.push_back(B[i]);
 88
 89
             }
 90
         }
 91 }
 92
 93
```

```
94 vector \( vector \( double \) \rightarrow \( UserBasedCF \) (vector \( vector \( double \) \rightarrow \tan int users \( Num, int \)
itemsNum)
95 {
        vector<vector<double> > predict(usersNum, vector<double>(itemsNum, 0));
96
        for (int i = 0; i < usersNum; ++i) //对每个用户进行预测
97
        {
98
            //找出user i未评分的item j, 预测user i 对item j的评分
99
            for (int j = 0; j < itemsNum; ++j)
100
101
102
103
104
                if (train[i][j])
105
                     continue;
                //如果item j没有被user i评过分,找出对 item j评过分的用户
106
107
                else
108
                    vector \( \double \rangle \ \sim; \)
109
                    vector \( \double \range \) historyScores;
110
                    for (int k = 0; k < usersNum; ++k)
111
                   {
112
                         //如果user k对item j 评过分,计算user k与user i的相似度
113
                         if (train[k][j]) //找出对item j 评过分的user k
115
                       {
116
                             // 为了计算user k与user i的相似度,必须找出二者共同评过分的items
117
                             // 把二者对共同评过分的items的评分分别存储在两个vector中
118
                             vector \( \double \rangle \) commonA, commonB;
119
                           FindCommon(train[i], train[k], commonA, commonB);
120
                             //如果二者存在共同评过分的items,计算相似度
121
                             if (!commonA.empty())
122
123
                                 sim. push_back (ComputeSim (commonA, commonB, 2));
124
                                 // 把user k对item j 的历史评分记录下来
125
                               historyScores.push back(train[k][j]);
126
127
                           }
128
                       }
129
130
                   }
                    // 计算出所有与user i存在共同评过分的items的users与user i之间的相似度,
131
                    // 保存在sim中,这些users对目标items j(即user i没有评过分)的历史评分记
132
133
                    // 录在historyScores中。利用这两个vector,计算出相似度加权平均分作为预
                     // 测user i对item j的评分
134
                    double SimSum = 0;
135
                     if (!sim.empty())
136
137
```

```
138
                          for(vector \( \double \)::size_type m = 0; m \( \sim. \size(); ++m)
139
                              SimSum += sim[m];
140
141
                     predict[i][j] = InnerProduct(sim, historyScores) / (SimSum);
142
                     cout << "User "<< i << " 对第 " << j << " 个Item的评分为 " <<
143
predict[i][j] << endl;</pre>
144
145
            }
146
147
       return predict;
148
149 }
150
151 int main()
152 {
153
        string FilePath1("E:\\Matlab code\\recommendation
system\\data\\movielens\\train.txt");
        string FilePath2("E:\\Matlab code\\recommendation
system\\data\\movielens\\test.txt");
155
        int row = 943;
156
        int col = 1682;
157
        vector<vector<double> > train = txtRead<double>(FilePath1, row, col);
158
        vector<vector<double> > predict = UserBasedCF(train, row, col);
159
        txtWrite(predict, "predict.txt");
160
        vector<vector<double> > test = txtRead<double>(FilePath2, 462, 1591);
161
        double rmse = ComputeRMSE(predict, test);
162
        cout << "RMSE is " << rmse <<endl;</pre>
163
        return 0;
164
165 }
```

4. 运行



由于程序没有优化,循环比较多,时间比较长,程序没写好,如果读者有兴趣帮我优化,请联系我,多谢,欢迎有兴趣的可以自己构造一个小点的数据集试一试,以前我用这个数据在Matlab中运行的 RMSE是1左右,所以如果读者运行结果得到测试集上的RMSE是0.9-1.3之间问题应该不大,如果偏离太多,程序设计可能就有问题。