



南开大学  
Nankai University

南 开 大 学

计 算 机 学 院

大数据计算及应用期末报告

---

推荐系统实验报告

---

姜奕兵 1710218

皮春莹 1711436

周辰霏 1712991

年级：2017 级

专业：计算机科学与技术

2020 年 6 月 13 日

## 摘要

基于 C++ 实现了协同过滤算法，将给定的数据集 train.txt 划分为训练集和验证集，用训练集对模型进行训练，用验证集计算 RMSE，验证模型的预测结果，最后将 test.txt 作为模型输入，得到预测的分数。

关键字：推荐系统；协同过滤算法

## 目录

一、 问题描述	1
二、 算法细节	1
(一) 数据预处理	1
(二) 划分验证集	1
1. 提取 txt 文本信息	4
(三) 计算用户的相似度	5
(四) 预测最终分数	9
1. 计算 baseline estimate	10
2. 计算修正后的 $r_{xi}$	11
3. 最终预测分数 $End_{xi}$	12
4. 对所有数据进行测试	13
三、 实验结果	15
四、 实验分析	16

## 一、 问题描述

使用学过的推荐算法预测 test.txt 文件中用户对电影的评分。在这里我们使用了基于用户的协同过滤算法。

数据集描述:

- train.txt 文件存储了用户 ID、物品 ID 及对应的评分，用于训练和验证模型，包括 19835 个用户，共 5001513 个评分。
- test.txt 文件存储了用户 ID、物品 ID，用于测试模型，包括与 train.txt 中相同的 19835 个用户，每个用户需要计算对应 6 种物品的评分。
- itemAttribute.txt 用于优化模型，包括 624961 件物品的属性描述。

## 二、 算法细节

基于用户的协同过滤通过不同用户对物品的评分来评测用户之间的相似性。基于用户的相似性做推荐，即给用户推荐和他兴趣相投的其他用户喜欢的物品。算法实现流程在下面分别进行介绍。

### (一) 数据预处理

### (二) 划分验证集

定义一些全局变量 user\_num、item\_num，以及结构体 itemScore，用于存放物品 ID 和分数:

定义

```
1 // 一共有0~19834个用户 (19835)
2 const int user_num = 19835;
3 // 存放每个用户的物品数目
4 int item_num[user_num];
5
6 struct itemScore
7 {
8     int item_id;
9     double item_sco;
10 };
11
12 // 开19835行 (用户)、不固定列的矩阵score_matrix, 存放txt读入的数据
13 itemScore **score_matrix = new itemScore *[user_num];
```

划分

```
1 // 分界线
2 int FineLine[user_num];
3 void div_read2buffer()
4 {
5     // 打开train.txt
6     ifstream infile1;
```

```
7   infile1.open("train.txt");
8   // 若失败,则输出错误消息,并终止程序运行
9   assert(infile1.is_open());
10
11   string s; // 存放读取行数据
12   string blank = " ";
13   string line = "|";
14   int UserID = -1; // user递增, 可以直接先赋值为-1
15
16   int NumOfMovie = 0; // 目前的用户所看过的物品数目
17   int clock = 0;
18
19   while (getline(infile1, s)) // 逐行读取数据
20   {
21       int LocLine = s.find(line); // 查找该行中"|"所在位置
22       if (s == "")
23           break;
24
25       if (LocLine != -1) // "|"存在, 新的用户, 读取他看过的物品数
26       {
27           UserID++;
28           clock = 0;
29
30           int len = s.length();
31           string ss = s.substr(LocLine + 1, len - LocLine);
32           NumOfMovie = stoi(ss); // 该用户看过的物品数目
33
34           item_num[UserID] = NumOfMovie;
35
36           // 动态初始化data数组中该UserID行
37           score_matrix[UserID] = new itemScore[NumOfMovie];
38       }
39
40       else // "|"不存在, 还是上一个用户, 现在s存的是物品号和评分
41       {
42           int LocBlack = s.find(blank); // 空格所在位置
43           string item_id = s.substr(0, LocBlack);
44           int leng = s.length();
45           string sco = s.substr(LocBlack + 2, leng - LocBlack);
46
47           double score = stof(sco); // 当前用户对目前行的物品的评分
48           int MovieID = stoi(item_id); // 当前物品编号
49
50           struct itemScore MS = { MovieID, score };
51           score_matrix[UserID][clock] = MS;
52           clock++;
53       }
54   }
```

```

55     infile1.close();
56     cout << "traindata.txt已读取完毕" << endl;
57 }

```

然后调用 `div_write2file()` 根据上面统计的数量及存储的信息，将数据按 15% 的比例划分为两部分。

### 统计信息

```

1  void div_write2file()
2  {
3
4      for (int i = 0; i < user_num; i++)
5          FineLine[i] = int(item_num[i] * 0.85);
6      cout << "Calculate FineLine." << endl;
7
8      fstream file1;
9      file1.open("train1.txt", std::ios::out | std::ios::app);
10     cout << "opened train1.txt" << endl;
11     for (int i = 0; i < user_num; i++)
12     {
13         for (int j = 0; j < FineLine[i]; j++)
14         {
15             if (!j)
16             {
17                 file1 << i << "|" << FineLine[i] << endl;
18                 file1 << score_matrix[i][j].item_id << " "
19                     << score_matrix[i][j].item_sco << endl;
20             }
21             else
22                 file1 << score_matrix[i][j].item_id << " "
23                     << score_matrix[i][j].item_sco << endl;
24         }
25     }
26     file1.close();
27
28     file1.open("verify.txt", std::ios::out | std::ios::app);
29     cout << "opened verify.txt" << endl;
30     for (int i = 0; i < user_num; i++)
31     {
32         int k = item_num[i] - FineLine[i];
33         for (int j = FineLine[i]; j < item_num[i]; j++)
34         {
35             if (j == FineLine[i])
36             {
37                 file1 << i << "|" << k << endl;
38                 file1 << score_matrix[i][j].item_id << " "
39                     << score_matrix[i][j].item_sco << endl;
40             }
41             else

```

```

42         file1 << score_matrix[i][j].item_id << " "
43         << score_matrix[i][j].item_sco << endl;
44     }
45 }
46 file1.close();
47 }

```

### 1. 提取 txt 文本信息

为了便于计算，增加全局变量（如下所示），item\_num 记录了 train1.txt 每个用户的物品数目，buffer 记录了 test.txt 中的用户号及对应物品号，ver\_item\_num 记录了 verify.txt 每个用户的物品数目，ver\_buffer 记录了 verify.txt 中的用户号及对应物品号，itemScore 的定义在上面已经提过，score\_matrix 用于记录 train1.txt 中的用户号、物品号及分数。

#### 变量设定

```

1 // 存放 train1.txt 每个用户的物品数目
2 int item_num[user_num];
3 int buffer[user_num][6];
4 // 存放 verify.txt 每个用户的物品数目
5 int ver_item_num[user_num];
6 vector<vector<int>>>ver_buffer;
7 itemScore **score_matrix = new itemScore *[user_num];

```

score\_matrix 计算过程由 load\_score() 函数实现，代码如下：

#### 分数矩阵

```

1 // 读取数据到 score_matrix 矩阵
2 void load_score()
3 {
4     // 打开 train1.txt
5     ifstream infile1;
6     infile1.open("train.txt");
7     assert(infile1.is_open()); // 若失败, 则输出错误消息, 并终止程序运行
8
9
10    string s; // 存放读取行数据
11    string blank = " ";
12    string line = "|";
13    int user_id = -1; // user 递增, 可以直接先赋值为 -1
14
15    int item_count = 0; // 目前的用户的物品数目
16    int clock = 0;
17
18    while (getline(infile1, s)) // 逐行读取数据
19    {
20        int line_pos = s.find(line); // 查找该行中 "|" 所在位置
21        if (s == "")
22            break;

```

```

23
24         if (line_pos != -1) // "|" 存在, 新的用户, 读取他的物品数
25         {
26             user_id++;
27             clock = 0;
28
29             int len = s.length();
30             string ss = s.substr(line_pos + 1, len - line_pos);
31             item_count = stoi(ss); // 该用户的物品数目
32
33             item_num[user_id] = item_count;
34
35             // 动态初始化data数组中该user_id行
36             score_matrix[user_id] = new itemScore[item_count];
37         }
38
39         else // "|" 不存在, 还是上一个用户, 现在s存的是物品号和评分
40         {
41             int blank_pos = s.find(blank); // 空格所在位置
42             string item_id = s.substr(0, blank_pos);
43             int leng = s.length();
44             string sco = s.substr(blank_pos + 2, leng - blank_pos);
45
46             double score = stof(sco); // 当前用户对目前行的物品的评分
47             int itemID = stoi(item_id); // 当前物品编号
48
49             struct itemScore MS = { itemID, score };
50
51             // 将该行数据存入效用矩阵对应位置
52
53             score_matrix[user_id][clock] = MS;
54             clock++;
55         }
56     }
57     infile1.close();
58     cout << "train1.txt 已读取完毕" << endl;
59 }

```

### (三) 计算用户的相似度

这里采用的是皮尔逊相关系数 (Pearson Correlation Coefficient), 它是余弦相似度在维度值缺失情况下的一种改进。

用户 x 和用户 y 之间的 Pearson 相似度的计算公式如图1所示。

Pearson 系数的取值范围和意义如图2所示。

由于 Pearson 相似度的分子部分需要多次地找两个用户共同的物品 ID, 为了更快速地进行二

### □ Pearson correlation coefficient

■  $S_{xy}$  = items rated by both users  $x$  and  $y$

$r_x, r_y \dots$  avg.  
rating of  $x, y$

$$sim(x, y) = \frac{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)(r_{ys} - \bar{r}_y)}{\sqrt{\sum_{s \in S_{xy}} (r_{xs} - \bar{r}_x)^2} \sqrt{\sum_{s \in S_{xy}} (r_{ys} - \bar{r}_y)^2}}$$

图 1: Pearson 相似度的计算公式

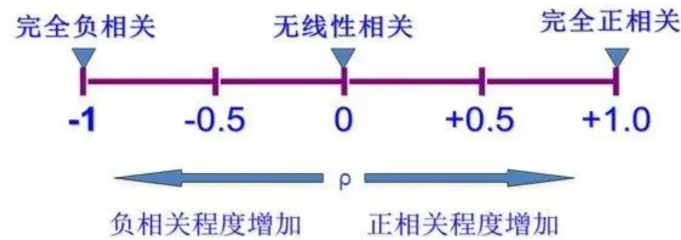


图 2: Pearson 系数的取值范围和意义

分查找，先将 score\_matrix 按照物品按编号从小到大升序排列：

#### 二分查找共同 ID

```

1 // 将每个用户的物品按编号从小到大升序排列
2 void sort_item_id()
3 {
4     for (int i = 0; i < user_num; i++)
5     {
6         sort(score_matrix[i], score_matrix[i] + item_num[i], compare)
7         ;
8     }
9 }
10 // 二分查找
11 int binary_search(itemScore **R, int uid, int n, int K)
12 {
13     // 在有序表R[0..n-1]中进行二分查找，成功时返回结点的位置，失败时返回-1
14     int left = 0, right = n - 1, mid; // 置当前查找区间上、下界的初值
15     while (left <= right)
16     {
17         if (R[uid][left].item_id == K)
18             return left;
19         if (R[uid][right].item_id == K)
20             return right; // 当前查找区间R[left..right]非
21                             // 空
22         mid = left + ((right - left) / 2);
23         // 使用(left+right)/2会有整数溢出的问题
24         if (R[uid][mid].item_id == K)
25             return mid; // 查找成功返回
26         if (R[uid][mid].item_id < K)

```



```

26         left = mid + 1;           // 继续在R[mid+1..right]
           中查找
27     else
28         right = mid - 1;          // 继续在R[left..mid-1]中
           查找
29     }
30     if (left > right)
31         return -1; // 当left > right时表示所查找区间内没有结果，查找失败
32 }

```

编程实现 Pearson 相似度的计算公式：

根据 Pearson 相似度的计算公式, 里面需要计算某一用户的所有评分的平均值, 编写 get\_avg\_user(int user\_id) 来实现这一功能：

#### 计算分数

```

1 // 求用户x对所有产品打分的平均值
2 double get_avg_user(int user_id)
3 {
4     double avg = 0;
5     double Sum = 0;
6     for (int i = 0; i < item_num[user_id]; i++)
7     {
8         Sum = Sum + score_matrix[user_id][i].item_sco;
9     }
10    avg = Sum / item_num[user_id];
11    return avg;
12 }

```

为了便于计算任意两个用户之间的 Pearson 相似度, 声明一个全局变量 pearson\_matrix 系数矩阵, pearson\_matrix[i][j] 表示第 i 个用户的第 j 个评分的物品 ID 以及评分减去用户 i 打分的平均值。由 get\_pearson\_matrix() 函数来计算 pearson\_matrix。

#### pearson 矩阵

```

1 itemScore **pearson_matrix = new itemScore *[user_num];
2 // 将score_matrix矩阵的数据转成pearson_matrix系数矩阵
3 void get_pearson_matrix()
4 {
5     for (int i = 0; i < user_num; i++)
6     { // 外层遍历所有用户
7         double avguser = get_avg_user(i);
8         pearson_matrix[i] = new itemScore[item_num[i]];
9
10        for (int j = 0; j < item_num[i]; j++)
11        { // 内层遍历当前用户看过的所有物品
12            pearson_matrix[i][j].item_id = score_matrix[i][j].
              item_id;
13            pearson_matrix[i][j].item_sco =
14            score_matrix[i][j].item_sco - avguser;
15        }
16    }
17 }

```

```

16     }
17 }

```

求用户 A 和用户 B 的 pearson 相似度时，分子部分需要对 AB 都有的物品所有物品编号 s 进行遍历，求出用户 A 对物品 s 的评分与用户 B 对物品 s 的评分的乘积之和，存放在 pear\_sum 变量中。分母部分需要用户 A、B 对所有物品的评分各自平方求和开根号在相乘。最后得出用户 A 和用户 B 的 pearson 相似度，代码如下：

#### 求解用户相似度

```

1 // 求用户A和用户B的pearson相似度
2 double pear_AB(int UserA, int UserB)
3 {
4     // 求出AB都有的物品a的编号
5     // 计算：用户A对物品a的评分*用户B对物品a的评分，存在multi_res
6     double multi_res[10000];
7     int clock = 0;
8
9     if (item_num[UserA] <= item_num[UserB])
10    {
11        for (int i = 0; i < item_num[UserA]; i++) // 外层遍历A用户的物
            品
12        {
13            int j = binary_search(pearson_matrix, UserB,
14            item_num[UserB], pearson_matrix[UserA][i].item_id);
15            if (j != -1)
16            {
17                multi_res[clock] = pearson_matrix[UserA][i].item_sco
18                * pearson_matrix[UserB][j].item_sco;
19                clock++;
20            }
21        }
22    }
23
24    if (item_num[UserA] > item_num[UserB])
25    {
26        for (int i = 0; i < item_num[UserB]; i++) // 外层遍历B用户的物
            品
27        {
28            int j = binary_search(pearson_matrix, UserA,
29            item_num[UserA], pearson_matrix[UserB][i].item_id);
30            if (j != -1)
31            {
32                multi_res[clock] = pearson_matrix[UserB][i].item_sco
33                * pearson_matrix[UserA][j].item_sco;
34                clock++;
35            }
36        }
37    }

```

```

38
39 // 将multi_res数组里的各数字求和，作为pearson系数的分子
40 double pear_sum = 0;
41 for (int i = 0; i < clock; i++)
42     pear_sum = pear_sum + multi_res[i];
43
44 // 求pearson系数的分母
45
46 // 求用户A对所有物品的评分各自平方求和开根号
47 double sum_userA = 0;
48 for (int i = 0; i < item_num[UserA]; i++)
49     sum_userA = sum_userA + pearson_matrix[UserA][i].item_sco
50     * pearson_matrix[UserA][i].item_sco;
51 double sqrt_userA = sqrt(sum_userA);
52
53 // 求用户B对所有物品的评分各自平方求和开根号
54 double sum_userB = 0;
55 for (int i = 0; i < item_num[UserB]; i++)
56     sum_userB = sum_userB + pearson_matrix[UserB][i].item_sco
57     * pearson_matrix[UserB][i].item_sco;
58 double sqrt_userB = sqrt(sum_userB);
59
60
61 // 求pearson系数
62 double pearson = pear_sum / (sqrt_userA * sqrt_userB);
63
64 return pearson;
65 }

```

#### (四) 预测最终分数

需要预测用户  $x$  对物品  $i$  的打分，如果仅从用户之间的相关性考虑，这个分数可以记为  $r_{ij}$ ，计算公式如图3所示。

$$r_{xi} = \frac{\sum_{j \in N(i;x)} s_{ij} \cdot r_{xj}}{\sum_{j \in N(i;x)} s_{ij}}$$

$s_{ij}$ ... similarity of items  $i$  and  $j$   
 $r_{xj}$ ... rating of user  $u$  on item  $j$   
 $N(i;x)$ ... set items rated by  $x$  similar to  $i$

图 3: 计算  $r_{ij}$  的公式

仅仅依靠用户相关性预测出来的分数，没有考虑到物品的特点以及特定用户打分的习惯，因此需要作出调整，加入基准，调节预测分数，如图4所示，将公式的加法第二个操作数记为修正后的  $r_{ij}$ ，把这一最终的分数记为  $End_{xi}$ 。下面分  $b_{ij}$  和  $r_{ij}$  两部分来计算  $End_{xi}$ 。

$$r_{xi} = b_{xi} + \frac{\sum_{j \in N(i;x)} s_{ij} \cdot (r_{xj} - b_{xj})}{\sum_{j \in N(i;x)} s_{ij}}$$

baseline estimate for  $r_{xi}$

$$b_{xi} = \mu + b_x + b_i$$

- $\mu$  = overall mean movie rating
- $b_x$  = rating deviation of user  $x$   
= (avg. rating of user  $x$ ) -  $\mu$
- $b_i$  = rating deviation of movie  $i$

图 4: 计算最终分数的公式

### 1. 计算 baseline estimate

为了计算总体平均值  $\mu$ ，编写函数 `get_avg_all()`，将计算结果保存在变量 `u` 中：

计算均值

```

1 // 求总体平均值
2 double get_avg_all()
3 {
4     double avg = 0;
5     double Sum = 0;
6     int count = 0;
7
8     for (int i = 0; i < user_num; i++)
9     {
10         for (int j = 0; j < item_num[i]; j++)
11         {
12             Sum = Sum + score_matrix[i][j].item_sco;
13             count++;
14         }
15     }
16
17     avg = Sum / count;
18     return avg;
19 }
20 double u = get_avg_all();

```

公式中需要求某一用户评分的平均值，编写函数 `get_avg_item(int it_id)` 来计算：

计算产品全部均值

```

1 // 求产品i的所有评分的平均值
2 double get_avg_item(int it_id)
3 {
4     double avg = 0;
5     double Sum = 0;
6     int count = 0;

```

```

7
8     for (int i = 0; i < user_num; i++) // 外层遍历所有用户
9     {
10         int j = binary_search(score_matrix, i, item_num[i], it_id);
11         // for (int j = 0; j < item_num[i]; j++) // 内层遍历每一用户看
            过的所有物品
12         // {
13         if (j != -1)
14         {
15             Sum = Sum + score_matrix[i][j].item_sco;
16             count++;
17         }
18         // }
19     }
20
21     avg = Sum / count;
22     return avg;
23 }

```

计算  $b_{xi}$ :

计算 bxi 均值

```

1 // 最后产品得分中的 bxi 的值
2 double bxi(int user_id, int it_id, double u)
3 {
4     double bx = get_avg_user(user_id) - u;
5     double bi = get_avg_item(it_id) - u;
6
7     double bxi = u + bx + bi;
8
9     return bxi;
10 }

```

## 2. 计算修正后的 $r_{xi}$

现在的  $r_{xi}$ ，分母是所有与目标用户有关、且看过目标物品的其他用户，与目标用户的相似度之和，分子是看过目标物品的其他用户与目标用户的相似度，与看过目标物品的其他用户对目标物品的打分与相应 baseline 之差的乘积之和。计算过程中需要对无关联或关联很小的用户进行特殊处理，避免出现分母为 0 的情况。代码如下：

rxix

```

1 // 最后产品得分中的 rxix 的值
2 double rxix(int user_id, int it_id, double u)
3 {
4     // 求所有与目标用户有关、且看过目标物品的其他用户，与目标用户的相似度之和
5     // 求(看过目标物品的其他用户与目标用户的相似度)*
6     (看过目标物品的其他用户对目标物品的打分)之和
7     double sum_sij = 0;

```

```

8      double sum_sij_rjx = 0;
9
10     for (int i = 0; i < user_num; i++) // 外层遍历用户
11     {
12         int j = binary_search(score_matrix, i, item_num[i], it_id);
13
14         if (j != -1) // 如果该用户也看过目标物品
15         {
16             double a = pear_AB(user_id, i);
17             if (a > 0) // 如果两用户是相似用户
18             {
19                 sum_sij += a;
20                 sum_sij_rjx += a * (score_matrix[i][j].item_sco
21                     - bxi(i, score_matrix[i][j].item_id, u));
22             }
23         }
24     }
25
26     if (sum_sij != 0 && sum_sij_rjx != 0)
27     {
28         double rix = sum_sij_rjx / sum_sij;
29         return rix;
30     }
31
32     else
33     {
34         double rix = 1000;
35         return rix;
36     }
37 }

```

### 3. 最终预测分数 $End_{xi}$

最终得到的预测分数即为  $b_{xi}$  与  $r_{xi}$  求和的结果，但要对上面  $r_{xi}$  说到的特殊情况进行判断和处理。

#### 最终分数

```

1 // 最终产品的终极得分
2 double End(int user_id, int it_id, double u)
3 {
4     double End = 0;
5     double rxi1 = rxi(user_id, it_id, u);
6
7     if (rxi1 == 1000) // 如果目标用户的相似用户都没有看过目标物品
8     {
9         // 以目标用户打分平均值作为结果
10        End = get_avg_user(user_id);
11    }

```

```

12     else
13     {
14         End = rxil + bxi(user_id , it_id , u);
15         if (End > 100)
16             End = 100;
17         if (End < 0)
18             End = 0;
19     }
20     cout << "user_id = " << user_id << " , it_id = " << it_id << " , End = " <<
        End << endl;
21     return End;
22 }

```

#### 4. 对所有数据进行测试

需要进行预测的是验证集 verify.txt 和测试集 test.txt，调用 read2buffer() 将需要测试的数据读入，分别存在 vec\_buffer 和 buffer 数组，并在 ver\_item\_num 数组中存放 verify.txt 每个用户的物品数目：

##### 获取数据

```

1 //把需要测试的数据（test.txt和verify）读入
2 void read2buffer()
3 {
4     ifstream infile;
5     infile.open("test.txt");
6     assert(infile.is_open());
7
8     string s; //存放读取行数据
9     string line = "|";
10    int user_id = -1; //user递增，可以直接先赋值为-1
11    int Row = -1;
12    int i = 0;
13
14    while (getline(infile , s))
15    {
16        int line_pos = s.find(line); //查找该行中"|"所在位置
17        if (s == "")
18            break;
19        if (line_pos != -1)
20        {
21            Row++;
22            i = 0;
23        }
24        else
25        {
26            buffer[Row][i] = stoi(s);
27            i++;
28        }

```

```

29     }
30     infile.close();
31
32     infile.open("verify.txt");
33     assert(infile.is_open());
34     string blank = " ";
35     user_id = -1; // user 递增, 可以直接先赋值为-1
36     Row = -1;
37
38     while (getline(infile, s))
39     {
40         int line_pos = s.find(line); // 查找该行中 "|" 所在位置
41         if (s == "")
42             break;
43         if (line_pos != -1)
44         {
45             ver_buffer.push_back({});
46             Row++;
47             string ss = s.substr(line_pos + 1, s.length() - line_pos);
48             // 该用户的物品数目
49             ver_item_num[Row] = stoi(ss);
50         }
51         else
52         {
53             int blank_pos = s.find(blank); // 空格所在位置
54             string item_id = s.substr(0, blank_pos);
55             ver_buffer[Row].push_back(stoi(item_id));
56         }
57     }
58     infile.close();
59 }

```

最后便是对每条数据调用 End() 函数计算分数, 再按照格式规定写入文本文件, 由 write2file(double u) 函数实现:

#### 计算全部分数

```

1 // 把读入的数据进行测试, 并写入结果文档 (result.txt)
2 void write2file(double u)
3 {
4     ofstream file1;
5     file1.open("test_result.txt", std::ios::out | std::ios::app);
6     for (int i = 0; i < user_num; i++)
7     {
8         for (int j = 0; j < 6; j++)
9         {
10             if (!j)
11             {
12                 file1 << i << "|" << "6" << endl;

```



```

13         file1 << buffer[i][j] << " " << int(End(i, buffer[i][j], u)) <<
14         endl;
15     }
16     else
17         file1 << buffer[i][j] << " " << int(End(i, buffer[i][j], u)) <<
18         endl;
19     }
20 }
21 file1.close();
22
23 // 计算验证集分数
24 file1.open("verify_result.txt", std::ios::out | std::ios::app);
25 for (int i = 0; i < user_num; i++)
26 {
27     for (int j = 0; j < ver_item_num[i]; j++)
28     {
29         if (!j)
30         {
31             file1 << i << "|" << ver_item_num[i] << endl;
32             file1 << ver_buffer[i][j] << " " << int(End(i, ver_buffer[i][j], u
33             )) << endl;
34         }
35         else
36             file1 << ver_buffer[i][j] << " " << int(End(i, ver_buffer[i][j], u
37             )) << endl;
38     }
39 }
40 file1.close();
41 }

```

## 主函数

```

1 int main()
2 {
3     load_score();
4     read2buffer();
5     double u = get_avg_all();
6     sort_item_id();
7     get_pearson_matrix();
8     write2file(u);
9
10    return 0;
11 }

```

## 三、 实验结果

在实验过程中我们对验证集测试结果与原结果对比最终得到 rmse 结果为 28.591，说明整体结果较为可靠。训练时间整体较长，计算验证集用时 6 天，由于我们没有渠道使用大型云服务器

只能使用小型云主机故耗时较长，而测试集结果用时2天，整体上和测试数据的大小成正比。内存占用比例约为30%-40%。

User ID	Original Rating	Predicted Rating
017		90
9101	88	90
112040	65	90
450906	75	90
453488	79	90
464229	85	90
393064	93	90
438303	80	0
1135		90
514730	92	90
516754	91	90
517718	90	90
521369	90	90
528144	93	90
532453	91	90
532569	91	90
534596	89	90
539126	93	90
540934	93	90
541988	91	90
543790	93	90
546722	89	90
546992	92	90
548098	92	90
550452	93	90
553890	92	90
556366	91	90
556458	96	90
558154	92	90
566644	91	90

图 5: 验证结果与原数据对比

```
dc2-user@117:~/pro$ g++ rmse.cpp
dc2-user@117:~/pro$ ./a.out
Score:28.591
```

图 6: rmse 测试结果

## 四、实验分析

### 计算 RMSE

```
1 // 默认验证集的原有分数和预测分数的位置是对应的
2 vector<int> predict, rating;
3
4 void getScore()
5 {
6     ifstream infile1;
7     infile1.open("verify.txt");
8     assert(infile1.is_open()); // 若失败,则输出错误消息,并终止程序运行
9     string s; // 存放读取行数据
10    string blank = " ";
11    string line = "|";
12    while (getline(infile1, s)) // 逐行读取数据
13    {
14        int LocLine = s.find(line); // 查找该行中"|"所在位置
15        if (s == "")
16            break;
17
18        if (LocLine == -1)
19        {
20            int LocBlack = s.find(blank); // 空格所在位置
21            string sco = s.substr(LocBlack + 2, s.length() - LocBlack);
```

```
22         double score = stof(sco);
23         rating.push_back(score);
24     }
25 }
26 infile1.close();
27
28 infile1.open("verify_result.txt");
29 assert(infile1.is_open()); // 若失败,则输出错误消息,并终止程序运行
30 while (getline(infile1, s)) // 逐行读取数据
31 {
32     int LocLine = s.find(line); // 查找该行中"|"所在位置
33     if (s == "")
34         break;
35     if (LocLine == -1)
36     {
37         int LocBlack = s.find(blank); // 空格所在位置
38         string sco = s.substr(LocBlack + 2, s.length() - LocBlack);
39         double score = stof(sco);
40         predict.push_back(score);
41     }
42 }
43 infile1.close();
44 }
45 double calcuRMSE()
46 {
47     getScore();
48     double res = 0;
49     int N = predict.size();
50     for (int i = 0; i < N; i++)
51     {
52         res += (pow((predict[i] - rating[i]), 2));
53     }
54     res = pow(res / N, 0.5);
55     return res;
56 }
```

## 参考文献

[1] 推荐系统 1-2