

数据科学与计算机学院

移动信息工程专业

本科生实验报告

(2017-2018学年秋季学期)

课程名称：人工智能

专业方向	班级	学号	姓名
移动互联网	1501	15352005	蔡景韬

## 一、实验题目

- 文本数据集简单处理

## 二、实验内容

### 1. 算法原理

- 得到总词条和词汇
  - 总词条
    - 总词条是指数据集中所有有用的词汇分行分列储存。
    - `vector<vector<string>>` 存储总表。
    - 采用边读边储存的形式对总表进行存储。
    - 我们需要的信息在1个数字与7个字符串之后。
    - 每次都让文件流输出到一个int对象 `file>>num`
      1. 输出成功, `file.good() == true`, 此时标志着一行读取完毕, 并读入第二行的数字部分, (将已储存的信息插入总表中), 连续读入7个字符串并舍弃, 为下一个判断分支准备。
      2. 输出失败, 则此时文件流中的是我们需要的信息, 首先 `file.clear()` 恢复文件流可用, 然后读取字符串存储在 `vector<string>` 中。
  - 词汇
    - 词汇是指数据集中所有有用的单词不重复的储存。
    - `vector<string>` 存储词汇。
    - 读入单个词时, 就开始判断是否插入到词汇中。
    - 为了方便得出矩阵, 增加 `map<string,int>` 的变量 `wordsMap`, 记录词在词汇中的id。
    - 读入一个词时, `wordsMap.find()`, 若不在词汇表中, 则插入词汇表, 并在 `wordsMap` 中映射该词的id `wordsMap[str]=words.size()-1`。
  - 记录文章数 `docNum=all.size()` 与词汇数 `wordSize=words.size()`。
- OneHot矩阵
  - OneHot矩阵是使用词汇的大小作为列数, 文章数作为行数, 当第i篇文章中出现第j个单词, 则 `OneHot[i][j]=1`。

- `OneHot[docNumMax][wordSizeMax]` 二维数组存储矩阵。
- 整个OneHot矩阵置为0。
- 遍历总词单，将每一行的单词对应的id作为下标，索引OneHot矩阵，将其置为1  
`OneHot[i][ wordsMap[str] ] = 1。`
- 由于TF矩阵的需要，增加cntOneHot数组存储每个词在一行中出现的次数。

## • TF矩阵

- TF (*Term Frequency*)：向量的每一个值标志对应的词语出现的次数归一化后的概率（一行中逐个单词出现的概率）
- `TF[docNumMax][wordSizeMax]` 二维数组存储矩阵。
- 整个TF矩阵置为0。
- 遍历总词单，将每一行的单词对应的id作为下标，索引TF矩阵与cntOneHot，将每个词在一行中出现的概率赋值给TF矩阵 `cntOneHot[i][id] / all[i].size()`。

## • TF\_IDF矩阵

- IDF (*逆向文件频率*)：表示单词在所有文章中出现的概率的倒数，并求对数。
- `IDF[wordSizeMax]` 二维数组存储矩阵。
- 遍历总词条每一列，如果OneHot矩阵值为1，则该词出现次数+1，之后用文章总数除以次数+1  
`log2((double)docNum / (double)(++cnt))。`
- TF\_IDF：表示单词在某篇文章中的标识度（可以用来区分这篇文章），等于词频TF \* 反文件频率IDF。
- 遍历TF矩阵，将TF矩阵值与IDF对应列值相乘即可得到TF\_IDF矩阵 `TF_IDF[i][j] = TF[i][j]*IDF[j]`

## • 三元表

- OneHot矩阵为稀疏矩阵，使用三元表可以大大减少存储空间。
- 使用TriTable类储存三元表

```

1 struct triTable{
2     int row,col,num;
3     vector<triPoint> table;
4 }TriTable;

```

■ triPoint是一个含有三个int类型成员类

- 遍历总词单，将每一行的单词对应的id与行数i成对压入集合（排序与去同），每行完成后将set中所有对(i,id)组成三元组(i,j,1)，插入三元表。

## • 三元表加和

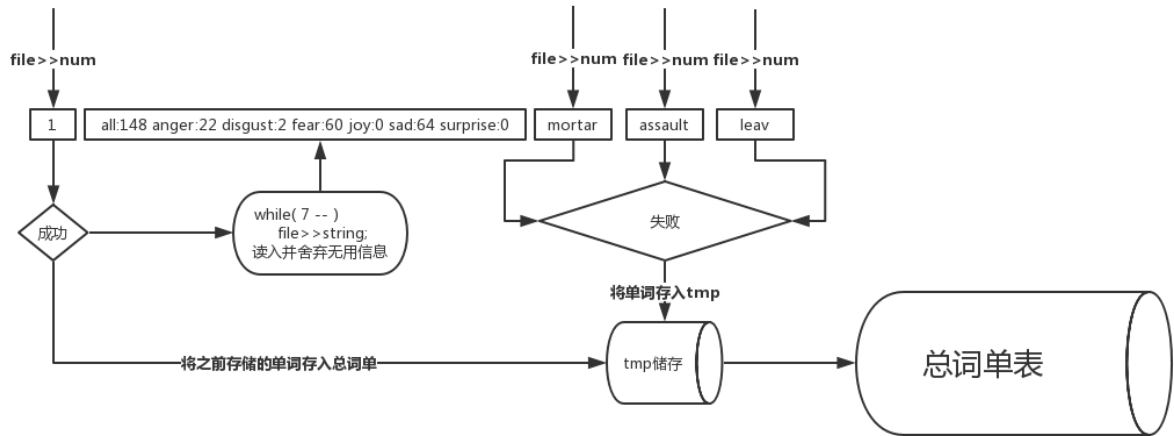
- 三元表结构体同上文。
- 使用权重值表示三元组在表中的前后位置，`num=i*10000+j`，num小的三元组排在表的靠前位置。
- 读入A，B两个三元表，新建ans变量存储结果。两个迭代器同时遍历A、B两表，考虑三种情况：
  1. 权重值 `aNum < bNum` :  
压入A当前位置的三元组；(A表迭代器推进)
  2. 权重值 `aNum == bNum` :  
压入A当前的三元组，其k值加上B表当前三元组k值；(A、B表迭代器同时推进)
  3. 权重值 `aNum > bNum` :

压入B当前位置的三元组；(B表迭代器推进)

- 如果A、B两个迭代器未到结尾，则把剩下三元组全部压入。

## 2. 伪代码

- 处理数据集



- 词汇表

```
1 > 词汇表: words
2 > 词汇映射Map: wordsMap(单词映射到词汇表下标id)
3 > 输入一个单词str
4 > if( str在wordsMap中存在 ){
5     将 str 插入 words;
6     将 (str,words.size()-1) 插入 wordsMap;
7 }
```

- OneHot矩阵

```
1 > 增加计数OneHot矩阵(单词在一行中出现的次数): cntOneHot
2 > 两个矩阵初始化置零
3 > 遍历总词单:
4     for(i: 遍历所有行)
5         for(j: 遍历所有列){
6             OneHot[i][单词_ij对应的词汇表下标] = 1;
7             cntOneHot[i][单词_ij对应的词汇表下标]++;
8         }
```

- TF矩阵

```
1 > TF矩阵初始化置零
2 > 遍历总词单
3     for(i: 遍历所有行)
4         for(j: 遍历所有列)
5             TF[i][单词_ij对应的词汇表下标]=这个单词出现的次数 / 该行词数;
```

- IDF矩阵

```

1 > 遍历OneHot矩阵
2   for(j: 遍历所有列){
3       单词计数: cnt = 0
4       for(i: 遍历所有行)
5           if(OneHot[i][j], 单词出现一次) cnt++;
6       IDF[j] = log2(总文本数/++cnt);
7   }

```

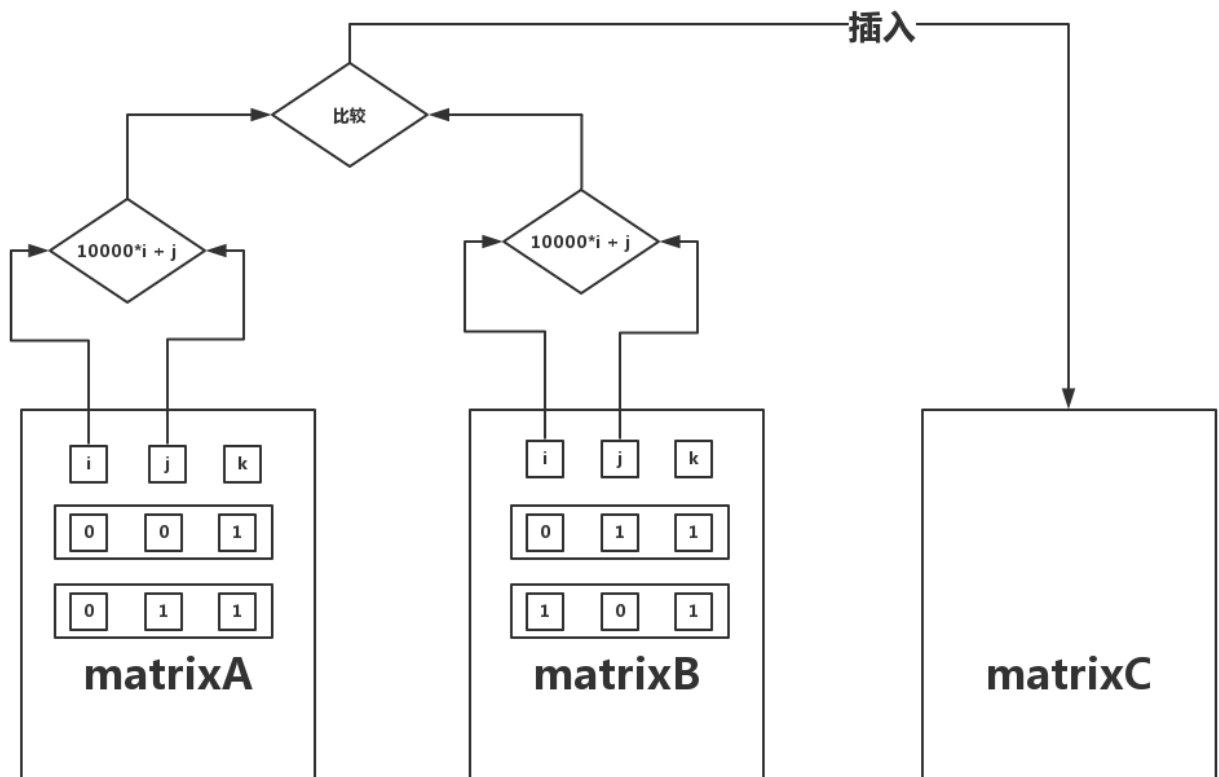
- 三元表

```

1 > 三元表数据结构:
2   int row:行数,col:列数,num:三元组数
3   vector<tripPoint(i,j,k)> table:三元表
4 > 遍历总词单
5   for(i: 遍历所有行){
6       声明集合: set<pair> idSet;
7       for(j: 遍历所有列)
8           将 单词_ij对应词汇表的下标 与 行数i 成对, 插入idSet
9       for(遍历idSet)
10          将 idSet中每对pair 与 1 组成三元组,插入table
11   }

```

- 三元表加和



### 3. 关键代码截图

- 处理数据集

```

1 while(!file.eof()){
2     if(file>>num){ // return file.good()
3         if(num!=1) all.push_back(tmp); // 把已储存的一行单词存入总词单（第一行不用压入）
4         tmp.clear(); // 清空tmp内容
5         int cnt=7;
6         while(cnt-->0) file>>word; // 读取并舍弃无用内容
7     }
8     else{
9         file.clear(); // 恢复文件流的状态
10        file>>word;
11        if(file.fail()) break;
12        tmp.push_back(word);
13        addWord(word); // 调用函数，判断并将单词插入词汇表
14    }
15 }all.push_back(tmp); // 将最后一行压入

```

o void addWord(string) 函数

```

1 if(wordsMap.find(str)==wordsMap.end()){ // 如果词汇不在词汇表中
2     words.push_back(str);
3     // 单词映射到词汇表对应的下标
4     wordsMap.insert(pair<string,int>(str,words.size()-1));
5 }

```

## • OneHot矩阵

```

1 memset(OneHot,0,sizeof(OneHot));
2 memset(cntOneHot,0,sizeof(cntOneHot)); // 计算单词在一行中出现的次数
3 for(int i=0; i<docNum; i++){
4     for(int j=0; j<all[i].size(); j++){
5         int id=wordsMap[all[i][j]]; // 单词在词汇表的位置下标
6         OneHot[i][id]=1; // 下标转化为OneHot矩阵的列数
7         cntOneHot[i][id]++;
8     }

```

## • TF矩阵

```

1 memset(TF,0,sizeof(TF));
2 for(int i=0; i<docNum; i++){
3     for(int j=0; j<all[i].size(); j++){
4         int id=wordsMap[all[i][j]]; // 单词在词汇表的位置下标
5         TF[i][id]=(double)cntOneHot[i][id]/(double)all[i].size(); // 单词在一行的出现概率
6     }

```

## • TF\_IDF

```

1  for(int j=0; j<wordSize; j++){
2      int cnt=0;
3      for(int i=0; i<docNum; i++){
4          if(OneHot[i][j]) cnt++; // 记录每个单词在所有文本中的出现次数
5          IDF[j]=log2((double)docNum/(double)(++cnt)); // 计算IDF矩阵
6      }
7      for(int i=0; i<docNum; i++){
8          for(int j=0; j<wordSize; j++){
9              TF_IDF[i][j] = TF[i][j]*IDF[j]; // 计算TF_IDF矩阵

```

- 三元表

```

1  TriTable.row = docNum;
2  TriTable.col = wordSize;
3  int cnt=0;
4  for(int i=0; i<docNum; i++){
5      set<pair<int,int> > idSet; // 记录三元组的ij对
6      for(int j=0; j<all[i].size(); j++){
7          int id=wordsMap[all[i][j]];
8          idSet.insert(pair<int,int>(i,id));
9      }
10     for(set<pair<int,int> >::iterator it=idSet.begin(); it!=idSet.end(); it++){
11         // ij对和1组成三元组，压入表
12         TriTable.table.push_back(triPoint(it->first,it->second,1));
13         cnt++ ;
14     }
15 }
16 TriTable.num=cnt; // 记录三元组数

```

- 三元表加和

```

1  while ( ait!=aEnd && bit!=bEnd ) { // A且B表都没到表尾
2      // 对行数加权，并与列数求和，大小关系表示其在三元表中的前后关系
3      aNum = ait->i*10000 + ait->j;
4      bNum = bit->i*10000 + bit->j;
5      // A当前三元组在前，先压入A当前三元组，并推进A
6      if(aNum<bNum) ans.table.push_back(*ait++);
7      // A、B两个三元组同行同列
8      else if(aNum==bNum) {
9          ans.table.push_back(triPoint(ait->i,ait->j,ait->k + bit->k));
10         ait++,bit++;
11     }
12     // B当前三元组在前，先压入B当前三元组，并推进B
13     else ans.table.push_back(*bit++);
14 }
15 // 如果A、B表任一个未到表位，将剩下的全部压入
16 while(ait!=aEnd) ans.table.push_back(*ait++);
17 while(bit!=bEnd) ans.table.push_back(*bit++);
18 ans.num = ans.table.size();

```

## 4. 创新点&优化

- 每次都让文件流输入到int类型，边读边处理，避免了读入一行再切割字符串的繁琐。
- 多加一个单词到词汇表下标的映射map，使得在求取OneHot与TF时，将时间复杂度从 $O(nm)$ 降为 $O(n \log m)$ 。
- 发现使用freopen比fstream所需的时间要长，于是改用fstream。猜测是因为前者重定向标准输入输出，而后者是一个比标准输入输出更为简便的文件流，所以后者速度就快（就像cin、cout与scanf、printf?）。
- 在计算三元表求和时，将行数i乘以10000的权重，再加上列数j，将三元组在三元表的前后关系转化为数值的大小关系。

### 三、实验结果及分析

#### 1. 实验结果展示示例

- OneHot矩阵

Row	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8	Col 9	Col 10	Col 11	Col 12	Col 13	Col 14	Col 15	Col 16	Col 17	Col 18	Col 19	Col 20
1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	1	0	0	0	0	1	1	1	1	1	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	0
5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1
6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

- TF矩阵

Row	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8	Col 9	Col 10	Col 11	Col 12
1	0.166667	0.166667	0.166667	0.166667	0.166667	0.166667	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0.25	0.25	0.25	0.25	0
3	0	0	0	0	0	0	0.166667	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0	0	0	0
6	0	0	0	0	0	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0
10	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0.111111	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	0	0	0	0
13	0	0	0	0	0	0	0	0	0	0	0	0
14	0	0	0	0	0	0	0	0	0	0	0	0
15	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0.0769231	0	0	0

- TF\_IDF矩阵（与同学结果的对拍）

