

Branch: master Course-Design / T6 / 6.cpp

Find file Copy path

lttzz 词频统计(bug再次修复)

7e903d4 6 hours ago

1 contributor

381 lines (355 sloc) 10 KB

```
1  /*****
2      > File Name:      6.cpp
3      > File Category: Course Design
4      > Author:      lttzz
5      > Mail:      3344517687@qq.com
6      *****/
7
8  //缺陷: 1)不规则动词复数,动词五态,s结尾的单词等可能会出现问
9  //      2)C++等末尾为两个符号的会被误删符号
10
11 /*思路:
12  *从文件中读取字符串,将其统一成小写字母后送到divideString()函数里,
13  *字符串末尾的标点符号删除(如果有的话),处理完的结果以文件流的形式保存到本地备用
14  *接着读取上一步处理后所得结果,进行二次处理.二次处理包括: 1)含数字单词变字母操作 2)复数变单数操作
15  *上一步结果以键值对的形式存入到map中, 由于map内部本身就是按序存储的(key),故直接输出map即可实现按字典序输出每个单词及其出现次数
16  *然后将键值对从map中取出来, 插入到vector中, 用sort()按词频降序排列后保存至本地并输出
17  */
18 #include <iostream>
19 #include <fstream>
20 #include <map>
21 #include <vector>
22 #include <algorithm>
23 #include <iomanip>
24 #include <iterator>
25 using namespace std;
26
27 ifstream pin("./bbe.txt");
28 ofstream pout("./tmp.txt");
29 ifstream fin("./tmp.txt");
30 ofstream fout("./out.txt");
31
32 map<string, int> wordlist;
33 vector<pair<string,int> > Sort;
34 string temp_input, second_process;
35
36 void init();
37 void sortWordTime(void);
38 void selectMode(int op);
39 void divideString(string initial);
40 string processString(string simple);
41 bool checkAbnormal(string str);
42 void inputString(string after_process);
43 int cmp(const pair<string, int>& x, const pair<string, int>& y);
44
45
46 int main(void)
47 {
48     init();
49     int op;
50     cin >> op;
51     selectMode(op);
52     return 0;
53 }
54
55
56 /**
57  * 若文本出现单符号或者单数字, 返回false, 直接跳过此次读入的内容
58  */
59 bool checkAbnormal(string str)
60 {
61     if (str.length() == 1)
62     {
```

```

63     if (!(str[0] >= 'a' && str[0] <= 'z') || !(str[0] >= 'A' && str[0] <= 'Z'))
64     {
65         return false;
66     }
67 }
68 return true;
69 }
70
71 void init()
72 {
73     if(!fin || !fout || !pin || !pout)
74     {
75         cout << "Error: can't input or output file" << endl;
76     }
77     while(pin >> temp_input)
78     {
79         transform(temp_input.begin(), temp_input.end(), temp_input.begin(), ::tolower);    //转化为小写字母
80         if (checkAbnormal(temp_input))
81         {
82             divideString(temp_input);
83         }
84     }
85
86     pin.close();
87     pout.close();
88
89     while(fin >> second_process)
90     {
91         string last_word;
92         last_word = processString(second_process);
93         inputString(last_word);    //存入键值表
94     }
95     sortWordTime();
96 };
97
98 /**
99  * 首先删去标点符号，增加空格，输出到pout文件中
100  * 其次，将一个字符串拆分成两个字符串
101  * @param initial 待处理的原始字符串
102  */
103 void divideString(string initial)
104 {
105     if (initial[0] == '\\')    //删除字符串开头的引号(单引号和双引号)
106     {
107         initial.erase(0, initial.find_first_not_of('\\'));
108     }
109     if (initial[0] == '\"')
110     {
111         initial.erase(0, initial.find_first_not_of '\"');
112     }
113
114     //删除末尾符号，显然当末尾为符号时，若末尾前一位仍为符号，则末尾必为引号，都删除之
115     int len = (int)initial.length();
116     if (initial[len - 1] > 'z' || initial[len - 1] < 'a') {    //删除字符串后面的标点符号(如果有的话)
117         initial.erase(len-1, 1);
118         len--;
119     }
120     if (initial[len - 1] > 'z' || initial[len - 1] < 'a') {    //删除字符串后面的标点符号(如果有的话)
121         initial.erase(len-1, 1);
122         len--;
123     }
124
125     if(initial == "i'm")    //处理 ' 号分词问题
126     {
127         pout << "i am ";
128     }
129     else if(initial == "i'd")
130     {
131         pout << "i would ";
132     }
133     else if(initial == "i've")
134     {
135         pout << "i have ";
136     }

```

```
137     else if(initial == "wasn't")
138     {
139         pout << "was not ";
140     }
141     else if(initial == "isn't")
142     {
143         pout << "is not ";
144     }
145     else if(initial == "don't")
146     {
147         pout << "do not ";
148     }
149     else if(initial == "he's")
150     {
151         pout << "he is ";
152     }
153     else if(initial == "she's")
154     {
155         pout << "she is ";
156     }
157     else if(initial == "can't")
158     {
159         pout << "can not ";
160     }
161     else if(initial == "wouldn't")
162     {
163         pout << "would not ";
164     }
165     else if(initial == "you're")
166     {
167         pout << "you are";
168     }
169     else if(initial == "it's")
170     {
171         pout << "it is";
172     }
173     else if(initial == "i'll")
174     {
175         pout << "i will ";
176     }
177     else if(initial == "you'll")
178     {
179         pout << "you will ";
180     }
181     else if(initial == "he'll")
182     {
183         pout << "he will ";
184     }
185     else if(initial == "he'd")
186     {
187         pout << "he would ";
188     }
189     else if(initial == "she'd")
190     {
191         pout << "she would ";
192     }
193     else if(initial == "--") //过滤破折号
194     {
195         return;
196     }
197     else if((initial[0] > 'z' || initial[0] < 'a') && (len == 1)) //过滤单符号
198     {
199         return;
200     }
201     else //原样输出
202     {
203         pout << initial << ' ';
204     }
205 }
206
207 /**
208  * 含数字单词变字母操作, 复数变单数操作(以's结尾的单词删除's)
209  * ss结尾了解一下, was, is等等了解一下, 2333
210 */
```

```

211 string processString(string simple)
212 {
213     string ss[] = {"s", "dos", "is", "was", "as", "bus", "gas", "has", "his", "news", "odds", "plus", "status", "thank
214
215     string last;
216     int simple_len = (int)simple.length();
217     if(simple == "1st")
218     {
219         last = "first";
220     }
221     else if(simple == "2nd")
222     {
223         last = "second";
224     }
225     else if(simple == "3rd")
226     {
227         last = "third";
228     }
229     else if(simple == "4th")
230     {
231         last = "forth";
232     }
233     else if(simple[simple_len-1] == 's')
234     {
235         if (simple[simple_len-2] == 's')           //双s结尾单词不是复数,故原样输出
236         {
237             return simple;
238         }
239
240         for (it : ss)                             //没辙了,手动输入白名单遍历特判吧
241         {
242             if (it == simple)
243             {
244                 return simple;
245             }
246         }
247
248         if (simple[simple_len-2] == '\\')
249         {
250             simple.erase(simple_len-2, 1);
251             simple_len--;
252         }
253
254         simple.erase(simple_len-1, 1);
255         simple_len--;
256         if(simple[simple_len-1] == 'e')
257         {
258             simple.erase(simple_len-1, 1);
259             simple_len--;
260             if(simple[simple_len-1] == 'i')
261             {
262                 simple.erase(simple_len-1, 1);
263                 simple += 'y';
264             }
265         }
266         last = simple;
267     }
268     else
269     {
270         last = simple;
271     }
272     return last;
273 }
274
275 /**
276  * string->map
277  * @param after_process 第一次处理后的字符串
278  */
279 void inputString(string after_process) {
280     auto it = wordlist.find(after_process);
281     if(it == wordlist.end())
282     {
283         wordlist.insert(pair<string, int>(after_process, 1));           //如果该单词没有出现,打包插入map,并将值设为1
284     }

```

```

285     else
286     {
287         wordlist[after_process]++;           //统计单词出现次数
288     }
289 }
290
291 /**
292  * sort()函数使用的cmp函数 使pair按单词出现次数多少排序
293  */
294 int cmp(const pair<string, int>& x, const pair<string, int>& y) {
295     return x.second > y.second;
296 }
297
298 /**
299  * map->vector 并调用sort()排序
300  */
301 void sortWordTime(void)
302 {
303     for (auto it = wordlist.begin(); it != wordlist.end(); it++)
304     {
305         if (it->first == " ")
306         {
307             continue;
308         }
309         Sort.push_back(make_pair(it->first, it->second));
310     }
311     sort(Sort.begin(), Sort.end(), cmp);
312 }
313
314 /**
315  * 选择模式 0--仅做处理后的文章,1-3对应于题目所做要求1-3
316  * @param op [description]
317  */
318 void selectMode(int op)
319 {
320     if (0 == op)
321     {
322
323     }
324     else if (1 == op)
325     {
326         cout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
327         fout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
328         for (auto it = wordlist.begin(); it != wordlist.end(); it++)
329         {
330             if(it->first[0] >= 'a' && it->first[0] <= 'z')           //如果键是单词 同时标准和文件输出键和值
331             {
332                 cout << left << setw(20) << it->first << setw(20) << it->second << '\n';
333                 fout << left << setw(20) << it->first << setw(20) << it->second << '\n';
334             }
335         }
336     }
337     else if (2 == op)
338     {
339         int MAXN;
340         cin >> MAXN;
341         map<string, int> wd;
342         cout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
343         fout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
344         for (vector<pair<string,int> >::iterator it = Sort.begin(); it != Sort.end(); it++)
345         {
346             if (MAXN == 0)
347             {
348                 break;
349             }
350             if(it->first[0] >= 'a' && it->first[0] <= 'z')
351             {
352                 wd.insert(pair<string, int>((*it).first, (*it).second));
353                 MAXN--;
354             }
355         }
356         for (auto it = wd.begin(); it != wd.end(); it++)
357         {
358             if(it->first[0] >= 'a' && it->first[0] <= 'z')

```

```
359         {
360             cout << left << setw(20) << it->first << setw(20) << it->second << '\n';
361             fout << left << setw(20) << it->first << setw(20) << it->second << '\n';
362         }
363     }
364 }
365 else if (3 == op)
366 {
367     cout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
368     fout << left << setw(20) << "Wordlist" << setw(20) << "Times" << "\n\n";
369     for (vector<pair<string,int> >::iterator it = Sort.begin(); it != Sort.end(); it++)
370     {
371         if(it->first[0] >= 'a' && it->first[0] <= 'z')
372         {
373             cout << left << setw(20) << it->first << setw(20) << it->second << '\n';
374             fout << left << setw(20) << it->first << setw(20) << it->second << '\n';
375         }
376     }
377 }
378
379 fin.close();
380 fout.close();
381 }
```