程序设计大作业: 单词查询

禹顺尧 2019202290

一.程序说明

本节简要说明应用程序的框架。源程序除主程序文件 main.cpp 外由 5 个主要部分构成,分别编码在不同的头文件和源代码文件中(IDE: CodeBlocks)。

● datadef.h/cpp 负责定义全局变量

● WordVectors.h/cpp 定义 Word 类,声明并实现成员变量与成员函数

● mode.h/cpp (1)构建词典 (2)两种单词查询的模式

● sml_algorithm.h/cpp 负责实现若干相似度的算法

● Word_Methods.h/cpp 实现程序中所需的函数

二.算法效果说明

本节通过截图展示各相似度算法的效果对比。单词由随机算法通过增加, 删除, 替换, 交换字母生成。

余弦相似度	编辑距离相似度	调整余弦相似度	扩展向量后的余弦
			相似度
请输入单词>>isable 最相似的5个单词是: 1. disable 0. 9258 2. sizable 0. 9258 3. lesbian 0. 9258 4. aisle 0. 9129 5. feasible 0. 9037	请输入单词>>isable 最相似的5个单词是: 1. disable 0. 8571 2. disabled 0. 7500 3. sizable 0. 7143 4. visible 0. 6667	请输入单词>>isable 最相似的5个单词是: 1. sizable 0. 8770 2. disable 0. 8652 3. lesbian 0. 8529 4. aisle 0. 8391 5. baseball 0. 8260	请输入单词>>isable 最相似的5个单词是: 1. disable 0. 9199 2. advisable 0. 8301 3. disabled 0. 8044 4. able 0. 7977 5. miserable 0. 7609
请输入单词>>boklet 最相似的5个单词是: 1. booklet 0. 9526 2. bottle 0. 8660 3. textbook 0. 8250 4. belt 0. 8165 5. bolt 0. 8165	请输入单词>>boklet 最相似的5个单词是: 1. booklet	请输入单词>>boklet 最相似的5个单词是: 1. booklet	请输入单词>>boklet 最相似的5个单词是: 1.booklet 0.9342 2.let 0.6742 3.bottle 0.6690 4.letter-box 0.6580 5.kettle 0.6228
请输入单词>>pmanist 最相似的5个单词是: 1.pianist 0.8819 2.assumption 0.8729 3.pastime 0.8571 4.postman 0.8571 5.implant 0.8571	请输入单词〉〉pmanist 最相似的5个单词是: 1. pianist 0.8571 2. humanist 0.7500 3. feminist 0.6250 4. humanism 0.6250 5. manifest 0.6250	请输入单词>>pmanist 最相似的5个单词是: 1. pianist 0. 7980 2. postman 0. 7654 3. implant 0. 7612 4. stamp 0. 7606 5. tsunami 0. 7567	请输入单词>>pmanist 最相似的5个单词是: 1.pianist 0.7877 2.humanist 0.7877 3.manifest 0.7161 4.postman 0.6923 5.humanism 0.6727
请输入单词>>acpoutability 最相似的5个单词是: 1. accountability 0. 9293 2. publication 0. 8908 3. actuality 0. 8908 4. typical 0. 8671 5. topical 0. 8671	请输入单词>>acpoutability 最相似的5个单词是: 1. accountability 0.8571 2. availability 0.6154 3. possibility 0.6154 4. probability 0.6154 5. actuality 0.6154	请输入单词>>acpoutability 最相似的5个单词是: 1. accountability 0. 8828 2. actuality 0. 8109 3. publication 0. 8104 4. availability 0. 7287 5. application 0. 7286	请输入单词>>acpoutabilit 最相似的5个单词是: 1.accountability 0.8500 2.ability 0.7884 3.stability 0.7839 4.liability 0.7259 5.availability 0.7190

请输入单词>>substaniii 最相似的5个单词是: 1. insist 0. 8944 2. distinguishable 0. 8847 3. antimissile 0. 8652 4. distinguish 0. 8652 5. sustain 0. 8642	2. substantive 0.7273 3. substance 0.7000 4. sustain 0.6000	请输入单词>>substaniii 最相似的5个单词是: 1. insist 0. 8211 2. distinguishable 0. 8084 3. distinguish 0. 8063 4. antimissile 0. 7860 5. acquisition 0. 7824	请输入单词>>substaniii 最相似的5个单词是: 1. substantial 0.7147 2. substantive 0.7056 3. sustain 0.6712 4. pianist 0.6390
最相似的5个单词是: 1. consulate 1. 0000 2. speculation 0. 9045 3. consultative 0. 8909 4. consultant 0. 8909 5. counsel 0. 8819	最相似的5个单词是: 1. speculate 0.7778 2. consulate 0.7778 3. insulate 0.6667 4. calculate 0.6667 5. circulate 0.6667	请输入单词>>sonculate 最相似的5个单词是: 1. consulate 1. 0000 2. consultant 0. 8131 3. consultancy 0. 7762 4. consult 0. 7564 5. speculation 0. 7417	请输入单词>>sonculate 最相似的5个单词是: 1.consulate 0.8235 2.speculation 0.7410 3.speculate 0.7233 4.congratulate 0.7206 5.calculate 0.7080

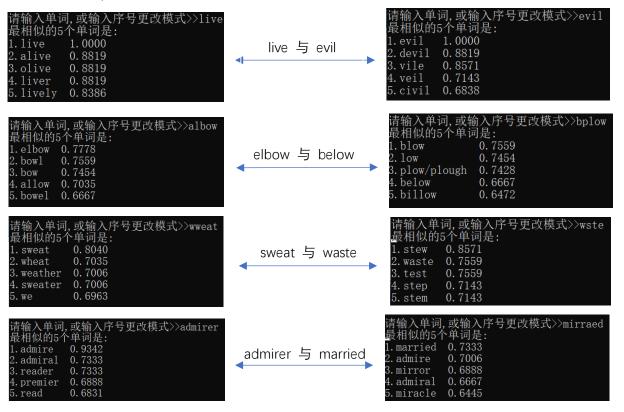
通过对比,发现算法效果:编辑距离相似度>调整余弦相似度>余弦相似度。余弦相似度只关注字母的组合,而编辑距离相似度不仅考虑了字母的组合,而且考虑了字母的排列,因此编辑距离相似度的算法效果更好。

程序中还定义了欧几里得距离相似度、算法效果与余弦相似度类似。

另外,在测试过程中发现,如果输入的单词拼写有误,且单词长度较短时,算法结果可能与预期相悖。推测原因在于拼写有误的单词可能与除原单词以外的某些单词更加相似。

三.程序亮点

1.为体现相同字母组合单词的差别,将向量扩充到 702 维,其中 1-26 维表示 a-z 26 个字母, 27-702 维表示任取两个字母的排列(aa,ab,ac,·····) (详见 sml_algrithm.cpp 中的函数 extended_vec_sml())。函数返回向量的余弦相似度。下面以截图体现其算法效果:



由以上四例可看出,扩充向量后算法能够对相同字母的不同单词进行区分,而普通的余弦相似度算法不能做到这一点。

2.参考网上的文献,对余弦相似度算法进行调整。(均值由另外的程序从词典中统计得来)由上一节的结果对比看出算法效果略有优化。

正因为余弦相似度在数值上的不敏感,会导致这样一种情况存在:

用户对内容评分,按5分制,X和Y两个用户对两个内容的评分分别为(1,2)和(4,5),使用余弦相似度得到的结果是0.98,两者极为相似。但从评分上看X似乎不喜欢2这个内容,而Y则比较喜欢,余弦相似度对数值的不敏感导致了结果的误差,需要修正这种不合理性就出现了调整余弦相似度,即所有维度上的数值都减去一个均值,比如X和Y的评分均值都是3,那么调整后为(-2,-1)和(1,2),再用余弦相似度计算,得到-0.8,相似度为负值并且差异不小,但显然更加符合现实。

那么是否可以在(用户-商品-行为数值)矩阵的基础上使用调整余弦相似度计算呢?从算法原理分析,复杂度虽然增加了,但是应该比普通余弦夹角算法要强。

```
double adj_cos_sml(Word a, Word b)

{
    int t=0;
    double al=0,bl=0;
    for (int i=0;i<=25;i++)
    {
        t+=sq(a.vec[i]-b.vec[i]);
        al=al+sq(a.vec[i]-tj[i]);
        bl=bl+sq(b.vec[i]-tj[i]);
    }
    return (al+bl-t)/(2*sqrt(al)*sqrt(bl));
}</pre>
```

3.编辑距离算法。由动态规划得出的答案表示修改的次数。通过 1- f[1][12] 归一化,表示出相似度,方便比较大小。

```
double ed_dstc_sml(Word a, Word b)

{
    int ll=a.len;
    int f[51][51]={0};
    for (int i=0;i<=l2-1;i++) f[0][i+1]=i+1;
    for (int i=0;i<=l1-1;i++) f[i+1][0]=i+1;
    for (int i=1;i<=l1;i++)
        for (int j=1;j<=l2;j++)
        {
          int k;
          if (a.proto[i-1]==b.proto[j-1]) k=0;else k=1;
          f[i][j]=min(min(f[i-1][j]+1,f[i][j-1]+1),f[i-1][j-1]+k);
    }
    double p=double (f[11][12])/max(11,12);
    return 1-p;
}</pre>
```

4. 函数 get_a_word(): 配合全局变量 n(代表模式的数字), a_word(储存待查询的单词), 通过递归调用, 实现了模式的即时切换, 并保证了程序的鲁棒性。

```
void get_a_word(int k)

{
    char s[51];
    if (k==0) printf("遺輸入单词,或輸入序号更改模式>>");
    fflush(stdin);
    gets(s);
    if (strcmp(s,"1")==0||strcmp(s,"2")==0)

{
        if (strcmp(s,"1")==0) n=1;else n=2;
            get_a_word(0);
            return;
        }
        if (strcmp(s,"3")==0) exit(0);
        while (judge(s)==false)
        {
            printf("輸入有误,...请输入单词或序号>>");
            get_a_word(1);
            return;
        }
        a_word.set(trans(s));
}
```

其中函数 judge 简单判断输入的字符串是否为一个英文单词:

```
bool judge(char *s)//判監法入的字符集是不是一个英文单词

{
    int ll=strlen(s);bool flag=false;
    if (ll==0) return false;
    for (int i=0;i<=ll-l;i++)
        if (s[i]<'A'||s[i]>'z'||(s[i]>'Z'&&s[i]<'a'))
        {
            if (s[i]!='\''&&s[i]!='-') return false;
        }
        else flag=true;
    return flag;
}
```

参考: https://blog.csdn.net/weixin_39050022/article/details/80732249
https://www.jianshu.com/p/a617d20162cf