

人工智能

教学班级：计科2班

专业：计算机科学

学号：19335262

姓名：张航悦

人工智能

任务1：TF-IDF

一、实验内容

1. 算法原理
2. 伪代码
3. 关键代码展示

二、实验结果及分析

三、思考题

任务2：k-NN分类

一、实验内容

1. 算法原理
2. 伪代码
3. 关键代码展示
4. 创新点

二、实验结果及分析

1. 结果展示及分析
2. 模型性能展示和分析

任务3：k-NN回归

一、实验内容

1. 算法原理
2. 伪代码
3. 关键代码

二、实验结果及分析

1. 结果展示及分析
2. 模型性能展示和分析

三、思考题

任务1：TF-IDF

一、实验内容

1. 算法原理

- 利用编码来处理文档数据
- 编码原因：图像由多个像素点构成，像素值之间是可计算的。与图像不同，文本一般很难直接被进行计算，所以我们需要对文本进行编码。
- 文档编码表示

- one-hot

使用一个V维向量表示一篇文章，向量的长度V为词汇表的大小。1表示存在对应的单词在该文章中存在，0表示不存在。one-hot表示仅考虑存在性，只取0/1两种值。**但这样就忽略了词语出现的频次所带来的影响，所以就进行了改进。**

- TF：词频归一化后的概率

每个文档的词频归一化后的概率。

$$tf_{i,j} = n_{i,j} / (\sum_v n_{v,d})$$

但是还是存在一些问题，**因为TF的方法容易收到一些常见词语的影响**，比如：的，吗，他们数量很多，可能会影响到一些文本关键主题的突出。

- IDF：逆向文档频率

如果这个词语在每篇文章都出现过，说明它能反应文本特性的能力不足，应该缩小它在词语向量中的权重。假设 $|C_i|$ 表示第*i*个词在 $|C|$ 篇文档中出现的次数，则

$$idf_i = \log \frac{|C|}{|C_i|+1}$$

- TF-IDF

$TF-IDF = TF * IDF$ ，可以把IDF理解为TF的一个权重值。用以评估一字词对于一个文件集或一个语料库中的其中一份文件的重要程度。字词的重要性随着它在文件中出现的次数成正比增加，但同时会随着它在语料库中出现的频率成反比下降。

2. 伪代码

进行TF-IDF运算的伪代码如下。

```
1 Procedure tfidf(dict,file_num,dictionary,record)
2  /*
3  input: dictionary是一个字典 键值为单词 用来统计各单词在文档中出现的总数
      和在每个文档中出现的次数
4  input: dict是对dictionary进行键值按字典序排序后的字典
5  input: file_num是输入文件semeval.txt的行数 即要计算的样例的个数
6  input: record是一个列表 顺序记录每个文档所包含的单词数
7
8  output: 返回TFIDF矩阵
9  */
10     word_num:=len(dict) # 得到总单词个数
11     result 创建存储结果的二维列表
12
13     for i:=0 to len(dict): #遍历单词
14         key:=dict[i]
15         itf=math.log(file_num/(dictionary[key][0]+1),10) #
idf计算公式
16         for j:=0 to len(dictionary[key]) # 遍历计算该单词在各个文
      档中的tfidf
17             val:=dictionary[key][j]
18             if val!=0 then
19                 result[j-1][i]:=val/record[j-1]*itf
20             end if
21
22         end for
23
24     end for
25
26     return result
27
```

3. 关键代码展示

该实验代码可分为三部分来实现。

- 读取文件

在读取文档同时计算了每个单词在各个文档中出现的次数和出现总数，只用遍历一次便可完成，效率较高

```

1 def read_csv():
2     file=open("semeval.txt")
3     lines=file.readlines() # 将每行文档内容存入一个列表中
4     file_num=len(lines) # 记录文件数目
5     record=[] # 创建一个列表 记录每个文档包含的单词数
6     dictionary={} # 创建一个空字典 用来统计各单词在文档中出现的总数 和
    在每个文档中出现的次数
7
8     for line in lines:
9         this_line=line.split('\t',2)[2] # 取一行中的单词部分
10        this_line.strip('\n') # 去掉末尾换行符
11        word_list=this_line.split() # 分割出每一个单词
12        record.append(len(word_list)) # 将文档的单词数记录在
    record中
13
14        for word in word_list: # 遍历这一行包含的单词
15            if word not in dictionary: # 如果该单词还未加入字典, 即
    该单词第一次在文档中出现
16                temp=[0 for i in range(file_num + 1)] # 创建一个
    临时列表 插入字典
17                temp[0]=1 # 0号单元用于存储该单词在文档中出现的总次数
18                temp[lines.index(line)+1]=1 # 更新该单词在该文档
    中的次数
19                dictionary[word]=temp # 插入字典
20
21            else: # 如果该单词已存在
22                if dictionary[word][lines.index(line) + 1]==
    0: # 若该单词首次在该文档中出现
23                    dictionary[word][0]= dictionary[word][0]
    + 1 # 更新该单词在文档中出现的总次数
24                    dictionary[word][lines.index(line) + 1]=
    dictionary[word][lines.index(line) + 1] + 1 # 更新该单词在该文档
    中的次数
25
26
27        dict=sorted(dictionary) # 结果要求按字典序排序 对键值进行排序
28        return dict,record,dictionary,file_num
29

```

- 计算TF-IDF矩阵

```

1 def tfidf(dict,file_num,dictionary,record):
2     word_num=len(dict) # 得到总单词个数（去除重复）

```

```

3     result=[[0 for i in range(word_num)] for i in
range(file_num)] # 创建存储结果的二维列表
4
5     for i in range(len(dict)): # 遍历单词
6         key=dict[i] # 当前单词
7         itf=math.log(file_num/(dictionary[key][0]+1),10) #
idf计算公式
8
9         for j in range(1, len(dictionary[key])): # 计算该单词在
各个文档中的tfidf
10             val=dictionary[key][j]
11             if val!=0:
12                 result[j-1][i]=val/record[j-1]*itf # 因为从j=1
开始存储第0个文件的信息 所以要j-1
13
14     return result

```

- 写入csv

```

1 def write_csv(result):
2     file = open("19335262_Zhanghangyue_TFIDF.txt", 'w')
3     for i in range(len(result)):
4         file.write(str(i + 1) + '\t') # str函数将数字转成字符串
5         for j in range(len(result[i])):
6             if result[i][j] != 0: #将结果不为0的写入
7                 file.write(str(result[i][j]) + ' ')
8         file.write('\n')
9     file.close()

```

- 主函数

```

1 def main():
2     dict,record,dictionary,file_num=read_csv()
3     result=tfidf(dict,file_num, dictionary, record)
4     write_csv(result)

```

二、实验结果及分析

- 实验要求只输出非0元素，结果如下所示。

```
1 0.4363994646005813 0.23422032704910614 0.33026244833605223 0.4363994646005813 0.3568792554806376 0.4657480077765282
2 0.6986220116647923 0.2531831680016751 0.623364512748797 0.6986220116647923
3 0.33026244833605223 0.3654046758885344 0.3994246729978553 0.4657480077765282 0.24038425475796782 0.3862277986565844
4 0.9314960155530564 0.5171499993242916 0.8311526836650627
5 0.3492530070538584 0.4657480077765282 0.4657480077765282 0.2921825602501574 0.3994246729978553 0.41557634183253134
6 0.20340946315524178 0.3862277986565844 0.4657480077765282 0.33026244833605223 0.4657480077765282 0.3862277986565844
7 0.6986220116647923 0.36057638213695176 0.6986220116647923 0.654599196900872
8 0.3199044638779115 0.3750700003848156 0.3994246729978553 0.17373993313995517 0.41557634183253134 0.4657480077765282
9 0.9314960155530564 0.8311526836650627 0.7137585109612752
10 0.3104986718510188 0.19938756073990768 0.27705089455502085 0.2909329764003875 0.24360311725902295 0.2574851991043896 0.0938997822238063 0.17670756266
11 0.1561468846994041 0.27705089455502085 0.22823615079610282 0.0850115854415511 0.2574851991043896 0.2574851991043896 0.16244995408261823 0.10088859127
12 0.2618396787603488 0.2493458050995188 0.2317366791939507 0.19493900066449127 0.27944880466591693 0.27944880466591693 0.22504200023088938 0.0907997321
13 0.4657480077765282 0.3027940735617203 0.33026244833605223 0.12751737816232667 0.3492530070538584 0.31084485349081276
14 0.33026244833605223 0.6497966688816375 0.4657480077765282 0.4363994646005813 0.4363994646005813
15 0.1442305528547807 0.2618396787603488 0.2618396787603488 0.1452065365836963 0.1604156348488878 0.15157344457063407 0.27944880466591693 0.152731631825
16 0.4986916101990376 0.3588976093318339 0.39631493800326273 0.4986916101990376 0.16901960800285135
17 0.18434984907593321 0.17310923094683797 0.21496061897378227 0.21496061897378227 0.07790251323128465 0.15242882230894722 0.17825898399534668 0.1124653
18 0.24367375083061407 0.2476968362520392 0.095638033621745 0.2995685047483915 0.34931100583239616 0.11349966518583596 0.2813025002886117 0.233133640118
19 0.20340946315524178 0.4657480077765282 0.4657480077765282 0.4657480077765282 0.4363994646005813 0.3862277986565844
20 0.20801928477052042 0.28967084899243833 0.2300306921524806 0.2520420995344407 0.327299598450436 0.327299598450436 0.2567656696456157 0.11349966518583
21 0.18028819106847588 0.3116822563743985 0.28967084899243833 0.327299598450436 0.327299598450436 0.23642475745840322 0.327299598450436 0.32729959845043
22 0.4657480077765282 0.3994246729978553 0.3750700003848156 0.12751737816232667 0.3750700003848156 0.41557634183253134
23 0.21015533996302507 0.0850115854415511 0.2574851991043896 0.161340596204107 0.3104986718510188 0.2909329764003875 0.2909329764003875 0.21659888962721
24 0.4363994646005813 0.41557634183253134 0.3248983344408187 0.14084967333570944 0.41557634183253134 0.4657480077765282
25 0.5626050005772234 0.599137009496783 0.654599196900872 0.5040841990688814
26 0.3116822563743985 0.2520420995344407 0.21913692018761805 0.22166234344865396 0.2520420995344407 0.27405350691640085 0.2676594416104782 0.27405350691
27 0.17732987475892317 0.27944880466591693 0.2618396787603488 0.2618396787603488 0.2618396787603488 0.2317366791939507 0.2618396787603488 0.169019608002
28 0.5236793575206976 0.47930960759742636 0.4282551065767652 0.47930960759742636 0.45008400046177877
29 0.3199044638779115 0.41557634183253134 0.41557634183253134 0.3750700003848156 0.4363994646005813 0.4657480077765282
30 0.45008400046177877 0.15302085379479202 0.4986916101990376 0.45008400046177877 0.5236793575206976
31 0.28967084899243833 0.22709555517129024 0.34931100583239616 0.34931100583239616 0.3116822563743985 0.10563725500178209 0.28967084899243833 0.26193975
32 0.4634733583789014 0.2025465344013401 0.41910360846463013 0.5236793575206976 0.4986916101990376
33 0.34931100583239616 0.327299598450436 0.095638033621745 0.13586473382015415 0.327299598450436 0.3116822563743985 0.327299598450436 0.2243110058323961
34 0.19938756073990768 0.2379195036537584 0.27705089455502085 0.2909329764003875 0.21015533996302507 0.2909329764003875 0.21326964258527434 0.2909329764
35 0.33605613271258755 0.4363994646005813 0.1811529784268722 0.4363994646005813 0.19266646495075537 0.4363994646005813
36 0.35620829299931256 0.27420382618106703 0.32148857175841333 0.25332839251274736 0.14891994269139014 0.33105239884850096 0.3423640054267331
37 0.4657480077765282 0.4363994646005813 0.4363994646005813 0.12751737816232667 0.24201089430616052 0.2545527197093593
38 0.20076028032780527 0.2880481137536465 0.3740566839433554 0.33105239884850096 0.2784842866635589 0.3740566839433554 0.399212578094167
```

三、思考题

1. IDF 的第二个计算公式中分母多了个1是为什么？

因为可能出现某个单词在所有文本中都没有出现过的情况，这样会导致分母为0的情况出现，加1可以避免这种情况的出现。在文本总量很大的时候，加1对结果影响不大。

2. IDF数值有什么含义？ TF-IDF数值有什么含义？

IDF是指逆向文档频率，如果这个词语在每篇文章都出现过，说明它能反应文本特性的能力不足，它可能只是一个常见词，则应该缩小它在词语向量中的权重。

TF-IDF就是把TF值和IDF值相乘。某一特定文件内的高词语频率，以及该词语在整个文件集中的低文件频率，可以产生出高权重的TF-IDF。因此，TF-IDF倾向于过滤掉常见的词语，保留重要的词语。

任务2：k-NN分类

一、实验内容

1. 算法原理

- k-NN分类问题是预测离散值的问题
- k-NN处理分类问题的步骤
 - 将文本转换成向量形式，可采用one-hot或TF-IDF向量表示
 - 相似度计算：在所有向量中找出和它距离最近的k个向量
计算距离的公式。当k=1的时候是曼哈顿距离，当k=2的时候是最常用的欧式距离。
$$dist(X, Y) = \sqrt[k]{\sum_1^n (x_i - y_i)^k}$$
 - 类别计算：取最相似的k个样本标签的众数
- one-hot和TF-IDF的实现代码虽都给出，但由于TF-IDF预测效果更为准确，后续实验均采用TF-IDF来实现k-NN算法
- 在实验中我利用numpy来实现TF-IDF矩阵的运算。我对训练集、验证集和测试集分别创建了自己的TF矩阵，但只创建了一个IDF向量存储每个单词的IDF值，即三个TF矩阵对应同一个IDF向量。后续计算只用进行IDF向量与TF矩阵的点乘即可。这样做既节省了储存空间，也利用了numpy加快了矩阵的运算。

2. 伪代码

- 整体思路

```
1  Proceuder knn_classification():
2      从训练集、验证集和测试集分别读出单词和对应的情感标签，存在
      train_word,train_feeling,valid_word,valid_feeling,test_word,te
      st_feeling中
3      wordlist:=count_word(train_word,valid_word,test_word) #将
      训练集、验证集、测试集单词合并成一个大的单词表，已去重
4
      trtrain_tfidf,valid_tfidf,test_tfidf:=tf_idf(wordlist,train_
      word,valid_word,test_word) #获得训练集、验证集、测试集的tfidf矩阵
5      /*调参过程*/
6      for k:1 to 10 do
7          predict_tag = knn_predict(valid_tfidf, train_tfidf,
      train_feeling, k)
8          accuracy = cal_accuracy(predict_tag, valid_feeling)
9          print(accuracy)
10     end for
11
12     /*写入结果*/
```

```

13         test_predict=knn_predict(test_tfidf,train_tfidf,train_feeling
14         ,k=5)
15         将预测结果test_predict写入文件

```

- k-NN分类

```

1  Function knn_predict(valid_met,train_met,train_feeling,k):
2  \* input: valid_met为验证集或测试集矩阵
3           train_met为训练集矩阵
4           train_feeling为训练集的情感标签信息
5           k为k-NN算法k值的选定 *
6
7  result初始化为初始值为0，大小为文档数的列表
8  for index:0 to 文档数 do
9      row:=valid_met[index] #row为第index行的文本向量
10     sum_list:=cal(row,train_met,n=1) #计算记录
11     sum_index=np.argsort(sum_list) # 将sum_list中的元素从小到大排列，提取其对应的index
12     dict初始化为空字典
13     for i:0 to k do
14         if dict中已经有了训练集排序后第i个文档的情感 then
15             dict中该情感对应次数+1
16         else
17             dict中添加键值对{第i个情感,1}
18         end if
19     end for
20     result[index]:=dict中情感的众数
21 end for
22 return result

```

- TFIDF矩阵计算

```

1  Function tfidf():
2      word_num:=word_list的大小
3      times:=np.zeros(word_num) #创建一个times向量，用来记录一个单词在文档中出现的总次数，用于计算idf
4      train_tf:=np.zeros((len(train_word), word_num)) #创建训练集单词的tf矩阵
5      valid_tf:=np.zeros((len(valid_word), word_num)) #创建验证集单词的tf矩阵
6      test_tf:=np.zeros((len(test_word), word_num)) #创建测试集单词的tf矩阵
7      filenum:=三个单词集合的大小

```



```

8
9      /*计算训练集的tf矩阵*/
10     for i:0 to len(train_lines) do
11         row:=train_lines[i]
12         temp:=从row中分离出的单词列表
13         num:=len(temp) #num为row中包含的单词数
14         for word in temp:
15             if 该单词在该文档中第一次出现 then
16                 times向量对应单词的位置+1
17             end if
18             对tf矩阵进行更新
19         end for
20     end for
21
22     /*同理计算验证集和测试集的tf矩阵*/
23
24     idf:=np.log10(filenum/(times+1)) # 计算得到idf向量
25     train_tfidf:=idf*train_tf # 计算得到tfidf矩阵
26     valid_tfidf:=idf*valid_tf
27     test_tfidf:=idf*test_tf
28     return train_tfidf,valid_tfidf,test_tfidf

```

- onehot矩阵计算

```

1  Function onehot(givenword,wordlist):
2  \* input: givenword 要求出onehot矩阵的单词列表
3          wordlist 总单词列表          *\
4
5      onehot初始化为 被求单词数x总单词数的二维矩阵
6      for index:0 to len(givenword)) do
7          temp:=该行中的单词
8          for i:0 to len(wordlist) do
9              if wordlist[i] in temp then # 如果单词在该文档内出现
onehot矩阵对应位置为1
10                 onehot[index][i]=1
11             end if
12         end for
13     end for
14     return onehot
15

```

- 距离计算

```

1  Function cal(row,train_tfidf,n):
2      row:=转变为行向量

```

```

3      temp:=将row进行行扩展得到的矩阵
4
5      if n==2 then #计算欧式距离
6          temp1:=temp-train_met
7          sum_list:=temp1行内求和形成的数组
8          sum_list:=sum_list开方后得到的数组
9      end if
10
11     if n==1 then #计算曼哈顿距离
12         temp1:=temp-train_tfidf的矩阵，对应单元取绝对值
13         sum_list=temp1行内求和形成的数组
14     end if
15     return sum_list

```

3. 关键代码展示

- k-NN分类

```

1  def knn_predict(valid_met,train_met,train_feeling,k):
2      """
3
4      :param valid_met: 待测集矩阵
5      :param train_met: 训练集矩阵
6      :param train_feeling: 训练集情感标签
7      :param k: knn k值选定
8      :return: 预测出的情感标签列表
9      """
10     result=[0 for i in range(valid_met.shape[0])] # 创建大小为待
    测集文件数的列表
11     for index in range(0,valid_met.shape[0]): #遍历待测集矩阵的每
    一行
12         row=valid_met[index]
13         sum_list=cal(row,train_met,n=1)
14         sum_index=np.argsort(sum_list) # 将sum_list中的元素从小到
    大排列，提取其对应的index
15         dict = {}
16         for i in range(k):#取k个样本
17             if train_feeling[sum_index[i]] in dict:
18                 dict[train_feeling[sum_index[i]]]+=1
19             else:
20                 dict[train_feeling[sum_index[i]]]=1
21         result[index] = max(dict,key=lambda x: dict[x]) # 求众
    数

```

```
22
23     return result
```

- TFIDF 矩阵计算

```
1  def tf_idf(wordlist,train_word,valid_word,test_word):
2      """
3
4      :param wordlist: 所有单词构成的表
5      :param train_word: 训练集单词
6      :param valid_word: 验证集单词
7      :param test_word: 测试集单词
8      :return: 三个集合的tfidf矩阵
9      """
10     word_num = len(wordlist) #计算出单词个数
11     times=np.zeros(word_num) #创建一个times向量，用来记录一个单词在
    文档中出现的总次数，用于计算idf
12     #times的下标和wordlist中单词的下标对应
13     train_tf=np.zeros((len(train_word), word_num)) #创建训练集单
    词的tf矩阵 规格为文档数x总单词数
14     valid_tf=np.zeros((len(valid_word), word_num)) #创建验证集单
    词的tf矩阵
15     test_tf=np.zeros((len(test_word), word_num)) #创建测试集单词
    的tf矩阵
16     filenum=len(train_word)+len(valid_word)+len(test_word) #计
    算文件总数
17     for i,row in enumerate(train_word):
18         temp=row.split(' ') #将一个文档内的的单词分离出来
19         num=len(temp) #num为一个文档中的单词数
20         for word in temp:
21             index=wordlist.index(word) # 获取该单词在wordlist中的
    下标
22             if train_tf[i][index]==0:
23                 times[index]+=1 # 若该单词在该文档中第一次出现 更新
    times值
24                 train_tf[i][index]=(train_tf[i][index]*num+1)/num
    #对tf值进行更新
25
26     for i,row in enumerate(valid_word):
27         temp=row.split(' ')
28         num=len(temp)
29         for word in temp:
30             index=wordlist.index(word)
31             if valid_tf[i][index]==0:
32                 times[index]+=1
```

```

33         valid_tf[i][index]=(valid_tf[i][index]*num+1)/num
34
35     for i,row in enumerate(test_word):
36         temp=row.split(' ')
37         num=len(temp)
38         for word in temp:
39             index=wordlist.index(word)
40             if test_tf[i][index]==0:
41                 times[index]+=1
42             test_tf[i][index]=(test_tf[i][index]*num+1)/num
43
44     idf=np.log10(filenum/(times+1)) # 计算得到idf向量
45     train_tfidf=idf*train_tf # 计算得到tfidf矩阵
46     valid_tfidf=idf*valid_tf
47     test_tfidf=idf*test_tf
48     return train_tfidf,valid_tfidf,test_tfidf
49

```

- onehot 矩阵计算

```

1  def onehot(givenword, wordlist):
2      """
3
4      :param givenword: 要求onehot矩阵的单词列表
5      :param wordlist: 所有单词构成的表
6      :return: 所求单词表的onehot矩阵
7      """
8      onehot=np.zeros(shape=(len(givenword), len(wordlist))) #
    创建一个二维矩阵
9      for index in range(len(givenword)):
10         temp=givenword[index].split(' ') # 将一行内的单词分割出来
11         for i, word in enumerate(wordlist):
12             if word in temp: # 如果单词在该文档内出现 onehot矩阵对
    应位置为1
13                 onehot[index][i]=1
14     return onehot
15

```

- 距离计算

```

1  def cal(row,train_met,n):
2      """
3
4      :param row: 一行单词
5      :param train_met: 训练集tfidf或onehot矩阵

```

```

6      :param n: 计算方式 n=1为曼哈顿距离 n=2为欧氏距离
7      :return: 记录该行和训练集的距离的列表
8      """
9      row=row.reshape(1,-1)
10     temp=np.repeat(row,train_met.shape[0],axis=0) # 将该行复制
    将temp扩充为一个维度和train_met维度相同的矩阵 便于计算
11
12     if n==2: # 计算欧式距离
13         temp1=np.square(temp-train_met)
14         sum_list=np.sum(temp1,axis=1) #行内求和
15         sum_list=np.sqrt(sum_list)
16
17     if n==1: #计算曼哈顿距离
18         temp1=np.abs(temp-train_met)
19         sum_list=np.sum(temp1,axis=1)
20     return sum_list

```

- 主函数

```

1  def main():
2      train_word,train_feeling=read_csv("train_set.csv")
3      valid_word,valid_feeling=read_csv("validation_set.csv")
4      test_word=read_csv("test_set.csv")
5      wordlist=count_word(train_word,valid_word,test_word)
6      valid_onehot=onehot(valid_word,wordlist)
7      train_onehot=onehot(train_word,wordlist)
8
9
10     #train_tfidf,valid_tfidf,test_tfidf=tf_idf(wordlist,train_word,valid_word,test_word)
11
12     for k in range(3,17):
13
14         predict_tag=knn_predict(valid_onehot,train_onehot,train_feeling,k)
15         accuracy=cal_accuracy(predict_tag,valid_feeling)
16         print('k = ' + str(k) + ', accuracy = ' + str(accuracy))
17
18         #predict_tag = knn_predict(valid_tfidf, train_tfidf, train_feeling, k)
19         #accuracy[k-3] = cal_accuracy(predict_tag, valid_feeling)
20         #print('k = ' + str(k) + ', accuracy = ' + str(accuracy[k-3]))

```

```

19
20 '''
21
22     test_predict=knn_predict(test_tfidf,train_tfidf,train_feeling,
23                               k=5)
24
25     test_output = pd.DataFrame({'words (split by space)':
26                                 test_word, 'label':test_predict})
27
28     test_output.to_csv('19335262_Zhanghangyue_KNN_classification.c
29 sv',index=None,encoding='utf8') # 参数index设为None则输出的文件前
30 面不会再加上行
31
32 '''
33

```

4. 创新点

引入了 `numpy` 矩阵，，加速运算。在计算距离的过程中将 `test` 向量扩展成了 `test` 矩阵，利用 `numpy` 封装好的矩阵操作，一次性计算出1个 `test` 和训练集所有样例之间的距离。在计算 `tfidf` 矩阵的时候利用一个 `idf` 向量和三个 `tf` 矩阵进行点乘运算进而得到 `tfidf` 矩阵。

二、 实验结果及分析

1. 结果展示及分析

通过循环调用KNN分类算法，遍历不同的 `k` 值，就可以找到在这个模型下训练集上表现最好时候的 `k` 值和它在验证集上对应的准确度。采用TFIDF的方法，取 `k=3` 到 `16`，依次输出精确度。

- 取曼哈顿距离时结果为

```

1 | k = 3, accuracy = 0.41479099678456594
2 | k = 4, accuracy = 0.40192926045016075
3 | k = 5, accuracy = 0.40514469453376206
4 | k = 6, accuracy = 0.3987138263665595
5 | k = 7, accuracy = 0.40514469453376206
6 | k = 8, accuracy = 0.3954983922829582
7 | k = 9, accuracy = 0.3762057877813505
8 | k = 10, accuracy = 0.3890675241157556
9 | k = 11, accuracy = 0.40192926045016075
10 | k = 12, accuracy = 0.40836012861736337

```

```
11 k = 13, accuracy = 0.41479099678456594
12 k = 14, accuracy = 0.3987138263665595
13 k = 15, accuracy = 0.40836012861736337
14 k = 16, accuracy = 0.40514469453376206
```

- 取欧氏距离时结果为

```
1 k = 3, accuracy = 0.26366559485530544
2 k = 4, accuracy = 0.2379421221864952
3 k = 5, accuracy = 0.21543408360128619
4 k = 6, accuracy = 0.2090032154340836
5 k = 7, accuracy = 0.19614147909967847
6 k = 8, accuracy = 0.1864951768488746
7 k = 9, accuracy = 0.17363344051446947
8 k = 10, accuracy = 0.17684887459807075
9 k = 11, accuracy = 0.1864951768488746
10 k = 12, accuracy = 0.19614147909967847
11 k = 13, accuracy = 0.2057877813504823
12 k = 14, accuracy = 0.22186495176848875
13 k = 15, accuracy = 0.22508038585209003
14 k = 16, accuracy = 0.2347266881028939
```

- 可能是由于训练及验证样本数量比较少，可以看到准确度随着k值总体在小范围内摆动，变化没有一个很明显的变化趋势，
- 取曼哈顿距离计算时，准确率明显高于取欧式距离。在k=3和13时，都达到了41.479%的准确率。

2. 模型性能展示和分析

对比采用 onehot 和 TFIDF 编码方式，取曼哈顿和欧式距离时k值和准确率。

| | 曼哈顿距离 | 欧式距离 | onehot | TFIDF | 准确率 | k |
|-----|-------|------|--------|-------|---------|------|
| 初始 | 1 | 0 | 0 | 1 | 41.497% | 3或13 |
| 优化1 | 0 | 1 | 0 | 1 | 23.367% | 3 |
| 优化2 | 0 | 1 | 1 | 0 | 40.836% | 13 |
| 优化3 | 1 | 0 | 1 | 0 | 40.836% | 13 |
| 最终 | 1 | 0 | 0 | 1 | 41.497% | 3或13 |

看到这个结果的时候我也觉得很神奇，欧式距离和曼哈顿距离的onehot编码方式在同一k值取得同一准确率（我debug了一下，的确不是bug的问题）。

可以看到在短文本情况下，onehot 编码方式的平均下来，在短文本的数据集上倾向比TFIDF 编码方式的表现更好，采用曼哈顿距离表现比欧式距离效果更好。

任务3：k-NN回归

一、实验内容

1. 算法原理

- k-NN回归问题是预测连续值的问题
- k-NN处理回归问题的步骤
 - 将文本转换成向量形式，可采用one-hot或TF-IDF向量表示
 - 相似度计算：在所有向量中找出和它距离最近的k个向量
 - 根据相似度加权计算：将距离最近的k个向量的距离的倒数作为权重，计算test属于该标签的概率（以happy为例）

$$P(\text{test } y \text{ is happy}) = \sum_{x=1}^n \frac{\text{train } x \text{ probability}}{d(\text{train } x, \text{test } y)}$$

- 在回归问题的应用景下，相关系数用于计算实际概率向量以及预测概率向量之间的相似性

$$COR(X, Y) = \frac{cov(X, Y)}{\sigma_X \sigma_Y} = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2 \sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

- 分别计算6个情感的真实概率值和预测概率值的相关系数，然后对6个维度取平均值计算得到最终的相关系数，评价模型。

2. 伪代码

k-NN回归的代码与k-NN分类的代码总体相同，主要区别在于读入函数和k-NN的标签判断的函数。相同部分不再赘述。

- 读入函数

```
1 Function read_csv(file_name):
2     打开文件，跳过表头开始读取
3     word新建空列表 #存储单词
4     feeling新建一个情感标签数x文档数的二维列表
5     #0-5行分别代表对应的6个情感，每一列代表一个文档的结果
```



```

6
7     if 读取的文件是训练集或验证集:
8         for row in reader:
9             在word中添加该行单词
10            for j:0 to 5 do
11                在feeling的对应行添加信息
12            end for
13        end for
14    end if
15    return word,feeling
16
17    if 读取的文件是测试集:
18        for row in reader:
19            在word中添加该行单词
20        end for
21    end if
22
23    return word

```

- k-NN回归预测函数

为保证六个情感标签概率之和为1，要记得做归一化处理。

```

1  Function knn_predict(valid_met,train_met,train_feeling,k):
2  \* input: valid_met为验证集或测试集矩阵
3            train_met为训练集矩阵
4            train_feeling为训练集的情感标签信息
5            k为k-NN算法k值的选定          *\
6
7      feeling_num:=len(train_feeling)
8      新建result 为便于后续写入csv文件 将result存为情感数x文件数的二维列
9      表
10     for index: 0 to valid_met.shape[0] do
11         row:=valid_met[index]
12         sum_list:=cal(row,train_met,n=1)
13         sum_index:=np.argsort(sum_list)
14         total:=0
15         for i:0 to k do
16             for j:0 to 5 do
17                 if sum_list[sum_index[i]]==0:
18                     sum_list[sum_index[i]]=0.001 # 若这两个标签
19                     完全相同，则将距离改为0.001 避免分母为0的情况出现
20                     result[j][index]+=train_feeling[j]
21                     [sum_index[i]]/float(sum_list[sum_index[i]]) # 根据公式计算概率

```

```

20         total+=train_feeling[j]
    [sum_index[i]]/float(sum_list[sum_index[i]]) # 计算总和 用于做归
    一化处理
21         end for
22     end for
23
24     for i:0 to 5 do
25         result[i][index]/=total #归一化处理
26     end for
27 end for
28 return result

```

3. 关键代码

- 读入函数

```

1 def read_csv(file_name):
2     """
3
4     :param file_name: 要读入的csv文件名
5     :return:返回单词列表 情感标签列表
6     """
7
8     with open(file_name,'r') as f:
9         reader=csv.reader(f)
10        first_row=next(reader)
11        feeling_num=len(first_row)-1
12
13        word=[] #存储单词
14        feeling=[[[] for i in range(feeling_num)] #创建一个情感标
15        签数x文档数的二维列表
16        #第0-5行分别代表对应的5个情感
17
18        if file_name=="train_set1.csv" or
19        file_name=="validation_set1.csv":
20            for row in reader:
21                word.append(row[0])
22                for j in range(feeling_num):
23                    feeling[j].append(float(row[j+1])) #一个文
24                    档的情感标签是一列 注意要转成float 不然无法运算
25
26            return word,feeling
27        else:

```

```

25         for row in reader:
26             word.append(row[1])
27         return word

```

- k-NN回归函数

```

1  def knn_predict(valid_met,train_met,train_feeling,k):
2      feeling_num=len(train_feeling)
3      result=[[0.0 for i in range(valid_met.shape[0])] for j in
4              range(feeling_num)]
5      # 为便于后续写入csv文件 将result列表存为情感数x文件数的规格
6      for index in range(0,valid_met.shape[0]):
7          row=valid_met[index]
8          sum_list=cal(row,train_met,n=1)
9          sum_index=np.argsort(sum_list)
10         total=0
11         for i in range(k):
12             for j in range(feeling_num):
13                 if sum_list[sum_index[i]]==0:
14                     sum_list[sum_index[i]]=0.001 # 若这两个标签
15                     # 完全相同,则将距离改为0.001 避免分母为0的情况出现
16                     result[j][index]+=train_feeling[j]
17                     [sum_index[i]]/float(sum_list[sum_index[i]]) # 根据公式计算概率
18                     total+=train_feeling[j]
19                     [sum_index[i]]/float(sum_list[sum_index[i]]) # 计算总和 用于做归
20                     # 一化处理
21         for i in range(feeling_num):
22             result[i][index]/=total #归一化处理
23     return result

```

二、实验结果及分析

1. 结果展示及分析

通过循环调用KNN回归算法，遍历不同的k值，就可以找到在这个模型下训练集上表现最好时候的k值和它在验证集上对应的相关系数。采用TFIDF的方法，取k=3到16，依次输出相关系数。相关系数利用提供的 validation 相关度评估.xlsx 完成计算

- 取曼哈顿距离时结果为

| | |
|----|---------------------------------------|
| 1 | k=3, coefficient=0.31635738697656807 |
| 2 | k=4, coefficient=0.32525132555084196 |
| 3 | k=5, coefficient=0.3227160899046577 |
| 4 | k=6, coefficient=0.32522254267169204 |
| 5 | k=7, coefficient=0.31211085366060115 |
| 6 | k=8, coefficient=0.316372564257192 |
| 7 | k=9, coefficient=0.30785183237412894 |
| 8 | k=10, coefficient=0.31469633558551763 |
| 9 | k=11, coefficient=0.32209775996910595 |
| 10 | k=12, coefficient=0.32451177938791037 |
| 11 | k=13, coefficient=0.3269034632830868 |
| 12 | k=14, coefficient=0.33245083187909436 |
| 13 | k=15, coefficient=0.32580554764405917 |
| 14 | k=16, coefficient=0.326640413606699 |
| 15 | k=17, coefficient=0.3264130073165443 |
| 16 | k=18, coefficient=0.32156733955725253 |
| 17 | k=19, coefficient=0.32148629511863946 |
| 18 | k=20, coefficient=0.32196083321168806 |

k=13时, 相关程度最高, 相关系数为0.3269034632830868

- 取欧式距离时结果为

| | |
|----|---------------------------------------|
| 1 | k=3, coefficient=0.2584148871442702 |
| 2 | k=4, coefficient=0.25389903943289444 |
| 3 | k=5, coefficient=0.25701020049336293 |
| 4 | k=6, coefficient=0.23673054270264302 |
| 5 | k=7, coefficient=0.24448088405709925 |
| 6 | k=8, coefficient=0.24896973852120216 |
| 7 | k=9, coefficient=0.24264594129196726 |
| 8 | k=10, coefficient=0.2386779279005535 |
| 9 | k=11, coefficient=0.24944531158353658 |
| 10 | k=12, coefficient=0.24757271951825477 |
| 11 | k=13, coefficient=0.24995861533584782 |
| 12 | k=14, coefficient=0.25392654817599386 |
| 13 | k=15, coefficient=0.24920159446999635 |
| 14 | k=16, coefficient=0.24917490031039183 |
| 15 | k=17, coefficient=0.2602504488449418 |
| 16 | k=18, coefficient=0.2541627393655171 |
| 17 | k=19, coefficient=0.25134492256685137 |
| 18 | k=20, coefficient=0.24169062856307888 |

k=17时, 相关程度最高, 相关系数为0.2602504488449418

2. 模型性能展示和分析

| | 曼哈顿距离 | 欧式距离 | onehot | TFIDF | 相关系数 | k |
|-----|-------|------|--------|-------|--------|----|
| 初始 | 1 | 0 | 0 | 1 | 0.3269 | 13 |
| 优化1 | 0 | 1 | 0 | 1 | 0.2603 | 17 |
| 优化2 | 0 | 1 | 1 | 0 | 0.2699 | 10 |
| 优化3 | 1 | 0 | 1 | 0 | 0.2798 | 10 |
| 最终 | 1 | 0 | 0 | 1 | 0.3269 | 13 |

在分类和回归问题上，都是TFIDF编码方式结合曼哈顿距离的效果最好。结果与分类问题类似。

三、 思考题

为什么是倒数？如果要求同一测试样本的各个情感概率总和为1，应该如何处理？

1. 因为距离越远理论上和当前文本的相似程度就越小，则它作为权重削减它的情感对当前文本情感的影响力，所以需要取倒数，分母越大，数越小。

2. 归一化处理：将所有情感概率加起来的到一个total，然后每个情感概率除以这个total作为归一化后的情感概率。