

NLP----关键词提取算法 (TextRank,TF/IDF)

河北凝讯科技订阅号 2018-11-24

TF/IDF

基本思想：TF是计算一个词在一篇文档中出现的频率，IDF是一个词在多少篇文档中出现过，显然TF越高证明这个词在这篇文章中的代表性就越强，而IDF越低则证明这个词在具有越强的区分能力。因此中和这两个数，就能较好地算出文档的关键词。

关键公式

$$tf * idf(i, j) = tf_{ij} * idf_i = \frac{n_{ij}}{\sum_k n_{kj}} * \log(\frac{|D|}{1 + |D_i|})$$

$|D_i|$ 是文档中出现词i的文档数量， $|D|$ 是文档数

附上书上抄来的代码

```
1 import jieba
2 import jieba.posseg as psg
3 import math
4 import functools
5
6 # 停用词表加载方法
7
8
9 def get_stopword_list():
10     # 停用词表存储路径，每一行为一个词，按行读取进行加载
11     # 进行编码转换确保匹配准确率
12     stop_word_path = './data/stopword.txt'
13     stopwords_list = [sw.replace('\n', '')
14                       for sw in open(stop_word_path, encoding = 'utf-8').readlines()]
15     return stopwords_list
16
17 # 分词方法，调用结巴接口
18
19
20 def seg_to_list(sentence, pos=False):
21     if not pos:
22         # 不进行词性标注的分词方法
23         seg_list = jieba.cut(sentence)
24     else:
25         # 进行词性标注的分词方法
26         seg_list = psg.cut(sentence)
27     return seg_list
28
29 # 去干扰词
30
31
32 def word_filter(seg_list, pos=False):
33     stopwords_list = get_stopword_list()
34     filter_list = []
35     # 根据POS参数选择是否词性过滤
36     # 不进行词性过滤，则将词性都标记为n，表示全部保留
37     for seg in seg_list:
38         if not pos:
```

```
39     word = seg
40     flag = 'n'
41     else:
42         word = seg.word
43         flag = seg.flag
44         if not flag.startswith('n'):
45             continue
46         # 过滤高停用词表中的词, 以及长度<2的词
47         if not word in stopword_list and len(word) > 1:
48             filter_list.append(word)
49
50     return filter_list
51
52 # 数据加载, pos为是否词性标注的参数, corpus_path为数据集路径
53
54
55 def load_data(pos=False, corpus_path='./data/corpus.txt'):
56     # 调用上面方法对数据集进行处理, 处理后的每条数据仅保留非干扰词
57     doc_list = []
58     for line in open(corpus_path, 'r', encoding = 'utf-8'):
59         content = line.strip()
60         seg_list = seg_to_list(content, pos)
61         filter_list = word_filter(seg_list, pos)
62         doc_list.append(filter_list)
63     return doc_list
64
65
66 def train_idf(doc_list):
67     idf_dic = {}
68     #总文档数
69     tt_count = len(doc_list)
70
71     #每个词出现的文档数
72     for doc in doc_list:
73         for word in set(doc):
74             idf_dic[word] = idf_dic.get(word,0.0)+1.0
```

```
75
76     #按公示转换为idf值，分母加一进行平滑处理
77     for k,v in idf_dic.items():
78         idf_dic[k]=math.log(tt_count/(1.0+v))
79
80     #对于没有在字典中的词，默认其仅在一个文档中出现，得到默认idf值
81     default_idf = math.log(tt_count/1.0)
82     return idf_dic,default_idf
83
84 def cmp(e1,e2):
85     import numpy as np
86     res = np.sign(e1[1]-e2[1])
87     if res != 0:
88         return res
89     else:
90         a = e1[0]+e2[0]
91         b = e2[0]+e1[0]
92         if a>b:
93             return 1
94         elif a == b:
95             return 0
96         else:
97             return -1
98
99 class TfIdf(object):
100     #统计tf值
101     def get_tf_dic(self):
102         tf_dic = {}
103         for word in self.word_list:
104             tf_dic[word] = tf_dic.get(word,0.0)+1.0
105
106         tt_count = len(self.word_list)
107         for k,v in tf_dic.items():
108             tf_dic[k] = float(v)/tt_count
109
110         return tf_dic
111
112     #四个参数分别是：训练好的idf字典，默认idf值，处理后的待提取文本，关键词数量
```

```

113 def __init__(self,idf_dic,default_idf,word_list,keyword_num):
114     self.word_list = word_list
115     self.idf_dic,self.default_idf = idf_dic,default_idf
116     self.tf_dic = self.get_tf_dic()
117     self.keyword_num = keyword_num
118
119     #按公式计算tf_idf
120     def get_tfidf(self):
121         tfidf_dic = {}
122         for word in self.word_list:
123             idf = self.idf_dic.get(word,self.default_idf)
124             tf = self.tf_dic.get(word,0)
125
126             tfidf = tf*idf
127             tfidf_dic[word] = tfidf
128
129             #根据tf_idf排序, 取排名前keyword_num的词作为关键词
130         for k,v in sorted(tfidf_dic.items(),key=functools.cmp_to_key(cmp),reverse=True)[:self.keyword_num]:
131             print(k+"/",end=' ')
132         print()
133
134     def tfidf_extract(word_list, pos=False, keyword_num=10):
135         doc_list = load_data(pos)
136         idf_dic, default_idf = train_idf(doc_list)
137         tfidf_model = TfIdf(idf_dic, default_idf, word_list, keyword_num)
138         tfidf_model.get_tfidf()
139
140     if __name__ == '__main__':
141         text = '6月19日,《2012年度“中国爱心城市”公益活动新闻发布会》在京举行。' + \
142             '中华社会救助基金会理事长许嘉璐到会讲话。基金会高级顾问朱发忠,全国老龄' + \
143             '办副主任朱勇,民政部社会救助司助理巡视员周萍,中华社会救助基金会副理事长耿志远,' + \
144             '重庆市民政局巡视员谭明政。晋江市人大常委会主任陈健倩,以及10余个省、市、自治区民政局' + \
145             '领导及四十多家媒体参加了发布会。□中华社会救助基金会秘书长时正新介绍本年度“中国爱心城' + \
146             '市”公益活动将以“爱心城市宣传、孤老关爱救助项目及第二届中国爱心城市大会”为主要内容,重庆市' + \
147             '、呼和浩特市、长沙市、太原市、蚌埠市、南昌市、汕头市、沧州市、晋江市及遵化市将会积极参加' + \
148             '这一公益活动。□中国雅虎副总编张银生和凤凰网城市频道总监赵耀分别以各自媒体优势介绍了活动' + \
149             '的宣传方案。□会上,中华社会救助基金会与“第二届中国爱心城市大会”承办方晋江市签约,许嘉璐理' + \
150             '事长接受晋江市参与“百万孤老关爱行动”向国家重点扶贫地区捐赠的价值400万元的款物。晋江市人大' + \
151             '常委会主任陈健倩介绍了大会的筹备情况。'
152
153         pos = True
154         seg_list = seg_to_list(text, pos)
155         filter_list = word_filter(seg_list, pos)
156
157         print('TF-IDF模型结果:')
158         tfidf_extract(filter_list)
159

```

TextRank

基本思路：每个词将自己的分数平均投给附近的词，迭代至收敛或指定次数即可，初始分可以打1

附上代码


```

1 def get_stopword_list():
2     path = './data/stop_words.utf8'
3     stopwords_list = [sw.replace('\n','') for sw in open(path,'r',encoding='utf8').readlines()]
4     return stopwords_list
5
6 def seg2list(text):
7     import jieba
8     return jieba.cut(text)
9
10 def word_filter(seg_list):
11     stopwords_list = get_stopword_list()
12     filter_list = []
13     for w in seg_list:
14         if not w in stopwords_list and len(w)>1:
15             filter_list.append(w)
16     return filter_list
17
18 str = '6月19日,《2012年度“中国爱心城市”公益活动新闻发布会》在京举行。' + \
19     '中华社会救助基金会理事长许嘉璐到会讲话。基金会高级顾问朱发忠,全国老龄' + \
20     '办副主任朱勇,民政部社会救助司助理巡视员周萍,中华社会救助基金会副理事长耿志远,' + \
21     '重庆市民政局巡视员谭明政。晋江市人大常委会主任陈健倩,以及10余个省、市、自治区民政局' + \
22     '领导及四十多家媒体参加了发布会。□中华社会救助基金会秘书长时正新介绍本年度“中国爱心城' + \
23     '市”公益活动将以“爱心城市宣传、孤老关爱救助项目及第二届中国爱心城市大会”为主要内容,重庆市' + \
24     '、呼和浩特市、长沙市、太原市、蚌埠市、南昌市、汕头市、沧州市、晋江市及遵化市将会积极参加' + \
25     '这一公益活动。□中国雅虎副总编张银生和凤凰网城市频道总监赵耀耀分别以各自媒体优势介绍了活动' + \
26     '的宣传方案。□会上,中华社会救助基金会与“第二届中国爱心城市大会”承办方晋江市签约,许嘉璐理' + \
27     '事长接受晋江市参与“百万孤老关爱行动”向国家重点扶贫地区捐赠的价值400万元的款物。晋江市人大' + \
28     '常委会主任陈健倩介绍了大会的筹备情况。'
29
30 win={}
31 seg_list = seg2list(str)
32 filter_list = word_filter(seg_list)
33 #构建投分表,根据窗口
34 for i in range(len(filter_list)):
35     if filter_list[i] not in win.keys():
36         win[filter_list[i]]=set()
37     if i-5 < 0:
38         lindex = 0
39     else:
40         lindex = i-5
41     for j in filter_list[lindex:i+5]:
42         win[filter_list[i]].add(j)
43
44 # 投票
45 time = 0
46 score = {w:1.0 for w in filter_list}
47 while(time<50):
48     for k,v in win.items():
49         s = score[k]/len(v)
50         score[k] = 0
51         for i in v:
52             score[i]+=s
53     time+=1
54
55 l = sorted(score.items(), key=lambda score:score[1],reverse=True)
56 print(l)
57
58

```