

# Experiment-1.2

## Ranked retrieval model

### 一、实验要求

#### 1.输入输出要求：

在Homework1.1的基础上实现最基本的Ranked retrieval model

- **Input** : a query (like Ron Weasley birthday)
- **Output**: Return the top K (e.g.,K = 10) relevant tweets.
- **Query** : 支持and, or ,not ; 查询优化可以选做；

#### 2.Use SMART notation: Inc.Itc

- **Document**: logarithmic tf (l as first character), no idf and cosine normalization
- **Query**: logarithmic tf (l in leftmost column), idf (t in second column), no normalization

#### 3.改进Inverted index

- 在Dictionary中存储每个term的DF
- 在posting list中存储term在每个doc中的TF with pairs (docID, tf)

#### 4.选做

- 支持所有的SMART Notations。

### 二、实验步骤

#### 1.创建新的倒排索引记录表(二元组)

(1)首先，计算每个doc中每个term出现的频率，并计算如下：

$$t,d = 1 + \log(tf_{t,d})$$

然后，对得到的数值进行归一化(normalization)。

(2)计算df和idf，即文档频率以及逆文档频率。其中idf=log(N/df)

(3)将原来的postings以及现在的postings\_for\_topk，以及document(存放逆文档频率的字典)保存成文件，使用numpy中的save函数保存为npz文件。

修正后的get\_postings代码如下：

```
1 def get_postings():
2     global postings
3     global all
4     global postings
```

```

5     f = open("tweets1.txt", 'r')
6     lines=f.readlines()
7     i=1
8     for line in lines:
9         all.append(i)
10        line=preprocess(line)
11        #print(line)
12        #计算TF
13        unique_terms=set(line)
14        t1=0
15        for every_term in unique_terms:
16            #统计TF
17            term_frequency =1+math.log10(line.count(every_term))
18            t1+=term_frequency*term_frequency
19            #print(term_frequency)
20            if every_term in postings.keys():
21                postings[every_term].append(i)
22            else:
23                postings[every_term]=[i]
24
25        ans1=math.sqrt(t1)
26        for every_term in unique_terms:
27            term_frequency = 1 + math.log10(line.count(every_term))
28            #重新计算一次，保存到二元组中的数值为归一化之后的数值
29            if every_term in postings_for_topk.keys():
30                postings_for_topk[every_term].append((i,
term_frequency/ans1))
31            else:
32                postings_for_topk[every_term] = [(i, term_frequency/ans1)]
33
34        #计算文档频率
35        for every_term in unique_terms:
36            if every_term in document.keys():
37                document[every_term]+=1
38            else:
39                document[every_term]=1
40        i=i+1
41
42    number=i-1
43    #计算idf
44    for every_key in document.keys():
45        document[every_key]=math.log10(float(number)/document[every_key])
46
47    np.save("idf.npy",document)
48    np.save("postings.npy",postings)
49    np.save("postings_for_topk.npy",postings_for_topk)
50    # print(all)
51    print("预处理完成")
52    print("文本数量共计: ",number,"项")
53    print("倒排索引记录表建立完成")

```

**注意：** postings , postings\_for\_topk,document都是词典 是在全局声明的。

## 2.Use SMART notation: Inc.Itc

参考计算方式如下图：

## tf-idf weighting has many variants

Term frequency		Document frequency		Normalization	
n (natural)	$tf_{t,d}$	n (no)	1	n (none)	1
l (logarithm)	$1 + \log(tf_{t,d})$	t (idf)	$\log \frac{N}{df_t}$	c (cosine)	$\frac{1}{\sqrt{w_1^2 + w_2^2 + \dots + w_M^2}}$
a (augmented)	$0.5 + \frac{0.5 \times tf_{t,d}}{\max_t(tf_{t,d})}$	p (prob idf)	$\max\{0, \log \frac{N-df_t}{df_t}\}$	u (pivoted unique)	$1/u$
b (boolean)	$\begin{cases} 1 & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$			b (byte size)	$1/CharLength^\alpha$ , $\alpha < 1$
L (log ave)	$\frac{1 + \log(tf_{t,d})}{1 + \log(\text{ave}_{t \in d}(tf_{t,d}))}$				

**SMART notation** for tf-idf variants.

Here CharLength is the number of characters in the document.

## tf-idf example: Inc.ltc

Query: "best car insurance". Document: "car insurance auto insurance".

word	query					document				product
	tf-raw	tf-wght	df	idf	weight	tf-raw	tf-wght	weight	n'lized	
auto	0	0	5000	2.3	0	1	1	1	0.52	0
best	1	1	50000	1.3	1.3	0	0	0	0	0
car	1	1	10000	2.0	2.0	1	1	1	0.52	1.04
insurance	1	1	1000	3.0	3.0	2	1.3	1.3	0.68	2.04

Key to columns: **tf-raw**: raw (unweighted) term frequency, **tf-wght**: logarithmically weighted term frequency, **df**: document frequency, **idf**: inverse document frequency, **weight**: the final weight of the term in the query or document, **n'lized**: document weights after cosine normalization, **product**: the product of final query weight and final document weight

$$\sqrt{1^2 + 0^2 + 1^2 + 1.3^2} \approx 1.92$$

$$1/1.92 \approx 0.52$$

$$1.3/1.92 \approx 0.68$$

$$\text{Final similarity score between query and document: } \sum_i w_{qi} \cdot w_{di} = 0 + 0 + 1.04 + 2.04 = 3.08$$

(1) 将query分成多个term，在postings\_for\_topk中可以返回一串二元组，每个二元组包含了docID以及tf，根据公式计算并累加求和得到针对相关句子的score。

(2) 根据score对涉及到的句子分数从高到低进行排序，选择topk个相关的句子，如果涉及到的句子不足topk个，那么直接把所有的句子由相似度从高到低输出。

相关代码如下

```

1 def IncIltc(terms):
2     # 计算tf
3     query_dict = {}
4     unique_terms = set(terms)
5     score = {}
6
7     # 计算涉及到的doc的分数
8     for every_term in unique_terms:
9         query_dict[every_term] = 1 + math.log10(terms.count(every_term))
10        for every_tuple in postings_for_topk[every_term]:
11            # 对应的tf
12            if every_tuple[0] in score.keys():

```

```

13         score[every_tuple[0]]\
14             += query_dict[every_term] * document[every_term] *
every_tuple[1]
15     else:
16         score.update({every_tuple[0]:
17             query_dict[every_term] *
document[every_term] * every_tuple[1]})
18     # 根据分数对score字典进行排序
19     ans = sorted(score.items(), key=lambda item: item[1], reverse=True)
20     #print(ans)
21     answer = []
22     if len(ans) <= topk:
23         for every_doc in ans:
24             answer.append(every_doc[0])
25     else:
26         temp = ans[:topk]
27         for every_doc in temp:
28             answer.append(every_doc[0])
29     return answer
30
31

```

### 3.已保存变量载入

载入保存的文件，而不是像第一个实验中动态建立，可以加快运行速度。

```

1 def load_npy():
2     print("文件载入中.....")
3     document=np.load("idf.npy").item()
4     postings=np.load("postings.npy").item()
5     postings_for_topk=np.load("postings_for_topk.npy").item()
6     tweet_text = open("tweet_text.txt", 'r+').readlines()
7     print("文件载入完成")
8     return document,postings,postings_for_topk,tweet_text

```

### 4.设计简单交互和运行提示

```

1 def search2():
2     global postings
3     global postings_for_topk
4     global tweet_text
5     input_str = input("Please input your query:")
6     if input_str=="Exit":
7         exit(0)
8     if input_str == "Back":
9         return False
10    terms = split_input1(input_str)
11    answer=Incltc(terms)
12
13    for docid in answer:
14        print(tweet_text[docid-1])
15
16    return True
17
18
19

```

```

20 def choose_model():
21     print("Tips:\n"
22           "You can choose the following two models:\n"
23           "Model 1(Boolean Retrieval Model):\n"
24           "    (1)you can input only one term;\n"
25           "    (2)your input can include 'and','or' and 'not',each operator
is a binary operator;\n"
26           "Model 2:(Ranked retrieval model)\n"
27           "    (1)You can enter the query text freely just like web
search\n"
28           "You can choose the mode by entering Arabic numerals like '1' or
'2'\n"
29           "You can input 'Exit' to exit the program\n"
30           "You can input 'Back' to go back to choose option\n")
31     while True:
32         opt=int(input("input your option here:"))
33         if opt==1:
34             while search1():
35                 pass
36         elif opt==2:
37             while search2():
38                 pass
39         else:
40             print("No such option")
41

```

### 三、实验结果

#### 1.界面与返回功能

```

文件载入中.....
文件载入完成
Tips:
You can choose the following two models:
Model 1(Boolean Retrieval Model):
    (1)you can input only one term;
    (2)your input can include 'and','or' and 'not',each operator is a binary operator;
Model 2:(Ranked retrieval model)
    (1)You can enter the query text freely just like web search
You can choose the mode by entering Arabic numerals like '1' or '2'
You can input 'Exit' to exit the program
You can input 'Back' to go back to choose option

input your option here:1
Please input your query:Back
input your option here:2
Please input your query:

```

#### 2.Ranked retrieval model

此处取topk=10

Test1 :

```
Please input your query:i love you
"\u266a I, I love you like John loves Sherlock! I, I love you like John loves Sherlock! And I keep hitting
re-peat-peat-peat-peat-peat-peat! \u266b"

"@KhloeKardashian you are my favorite kardashian i love you khloe =)"

"i love Khloe Kardashian"

"I love Kourtney Kardashian!!!"

"I love the middle"

"@KUWTKardashians No prob love you!!"

"I love you Justin\u266d\u266c97 @justinbieber #KnockKnockLive @KnockKnockFOX @RyanSeacrest http://t.co/24r118DjZk"

"RT @HammerFist3: I love u dad I hate you Outback kids meal"

"I love the middle \u266d\u266cfa"

"oh mans, christina. i love you, but really? http://tinyurl.com/5uph7so #superbowl #fail"
```

## Test2

```
Please input your query:i want to fly
"I want mcdonalds breakfast."

"@Kelliekidd30 i want McDonalds .."

"RT @MarketWatch: Amazon describes how it wants delivery drones to fly http://t.co/6byHFI0yob"

"I just want my snow blower to work"

"I totally want to cruise in this storm!"

"I want some Gilligans hush puppies !"

"I don't want an agreement, I just want something real"

"I want my dad to come home quicker because i want to buy some clothes online"

"When I grow up I want to be Tony Mendez #argo @FenellaHH"

"I really want some breakfast from McDonalds!"
```

## 四、实验改进与不足

(1) 实验中将倒排索引记录表存储在文件中，方便了每次直接从文件中读取变量。

(2) 对于tweet数据集只保留了text信息，同时，对每一行文本使用集合set处理，同时保留了词频的信息。