QINMIN HU 2018年5月29日

LAB 5 检索模型

- 一. 实验目的
 - ~ 在unigram假设下进行检索
- 二. 实验步骤

if s1.text==None:
 doc_list.append('')
else:

for word in s1.text.split():

while ' 'in word or ''' in word or '.' in

word = word.replace("_",

word or '-' in word or '!' in word or ';' in word or ',' in word or '\$' in word or '/" in word or '\" in word o

```
"").replace('`', "").replace(".", "").replace("-', "").replace("!","").replace(";",
"").replace(",", "").replace("$", "").replace("/", "").replace('/', "").replace('\'', "").replace('\'', "").replace("\'', "").replace("\'', "").replace("\'', "").replace('\'', "").
```

 $if \ re.findall(r'\d+', word) \ or \ word ==$

continue
doc_list.append(word)
for s2 in doc.findall('DOCNO'):

```
docno=s2.text.strip()
                                                                    doc_filaname[docno]=len(doc_list)
                    def loadDocument(dirname):
                                  print('loaddocument')
                                  for parent, dirnames, filenames in os. walk (dirname):
                                                   for e in filenames:
                                                                    print(e)
                                                                    getDocLength(dirname,e)
                    def loadDic():
                                  print('loaddic')
                                  f_dic=open("inverted_index.json",encoding='utf-8')
                                  word_list=json.load(f_dic)
                                  for word in word list:
                                                   word_dic[word[0]]=word[1]
                    def loadQuery():
                                  print('loadquery')
                                  xmldoc = ET.parse('topics_new.xml')
                                  for doc in xmldoc.findall('top'):
                                                   test="
                                                   for s1 in doc.findall('num'):
                                                                    query_id.append(s1.text.replace("Number:","").strip())
                                                   for s2 in doc.findall('desc'):
                                                                    test=test+' '+s2.text.replace("Descript:","").strip()
                                                   for s3 in doc.findall('narr'):
                                                                    test=test+' '+s3.text.replace("Narr:","").strip()
                                                   for s4 in doc.findall('title'):
                                                                    test=test+' '+s4.text.replace("Topic:","").strip()
                                                   while '_' in test or ''' in test or '.' in test or '-' in test or '!' in test or ';'
in test or ',' in test or '$' in test or '/" in test or '/' in test or '}' in test or '}' in test or '*' in
test or '#' in test or '^' in test or '|' in test or '~' in test or '=' in test or '\'' in test or '+' in
test or ':' in test or '?' in test:
                                                                    test = test.replace("_", "").replace("`', "").replace(".",
"").replace(''-', "").replace("!","").replace(";","").replace(",", "").replace("$",
"").replace("//", "").replace('\forall ', "").replace(
"").replace("|", "").replace('~', "").replace('=', "").replace('\'',
"").replace('+',"").replace(':',"").replace('?', "")
                                                   if re.findall(r'\d+', test) or test == "":
                                                                    continue
                                                   query[s1.text.replace("Number:","").strip()]=test.split()
```

2. 根据查询在unigram假设下进行检索,打分 $P(D \mid Q) \propto P(D) \prod_{t \in Q} \left[\alpha P_{MLE}(t \mid M_D) + (1-\alpha) P_{MLE}(t \mid M_C)\right]$ $P_{MLE}(t \mid M_D) = (tf_(t,D))/L_D$ $P_MLE(t \mid M_C) = (cf_t)/T$ 平滑方式: Add-one、线性插值、Dirichlet... 10152130138_丁婉宁_unigram.py 中: def getUnigram(title,words): docDic={} for doc,doc_len in doc_filaname.items(): p=1for word in words: if doc in word_dic[word]: pwd1=word_dic[word][doc]*1.0/doc_len else: pwd1=1wDoc=0for d,freq in word_dic[word].items(): wDoc=wDoc+freq $p=p*(0.9*pwd1+(1-0.9)*(wDoc/len(doc_filaname)))$ docDic[doc]=p score[title]=docDic 10152130138_丁婉宁_unigram_dic.py 中: def getUnigram(title,words,T): print('getunigram') docDic={} for doc,doc_len in doc_filename.items(): p=1for word in words: if word in word_dic.keys(): if doc in word_dic[word]: pwd1=word_dic[word][doc]*1.0/doc_len else: pwd1=0else: pwd1=0 $p=p*(0.9*pwd1+(1-0.9)*(word_freq[word]/T))$ docDic[doc]=p score[title]=docDic

评测:

test_10152130138_Unigram_0.5.txt:

```
est1.txt >
runid
                             all
                                       10152130138_dingwanning_unigram
num_q
                             all
                                       50
num_ret
                             all
                                       25000
                                       9805
                             all
num_rel
num_rel_ret
                             all
                                       1641
                             all
                                       0.0417
map
                                       0.0177
gm_map
                             all
                                       0.1167
Rprec
                             all
bpref
                             all
                                       0.0996
recip_rank
                             all
                                       0.2018
iprec at recall 0.00 iprec at recall 0.10
                                       0.3033
                             all
                             all
                                       0.1731
 iprec_at_recall_0.20
                             all
                                       0.1156
iprec at recall 0.30 iprec at recall 0.40
                             all
                                       0.0580
                                       0.0008
                             all
                                       0.0000
 iprec_at_recall_0.50
                             all
iprec_at_recall 0.60
iprec_at_recall 0.70
                             all
                                       0.0000
                                       0.0000
                             all
                                       0.0000
 iprec_at_recall_0.80
                             all
 iprec_at_recall_0.90
                             all
                                       0.0000
 iprec_at_recall_1.00
                             all
                                       0.0000
P_5
P_10
                                       0.1080
                             all
                             all
                                       0.1320
P_15
                             all
                                       0.1507
P_20
                             all
                                       0.1440
P_30
                                       0.1533
                             all
                                       0.1448
P_100
                             all
P_200
                             all
                                       0.1195
P_500
                                       0.0656
                             all
P_1000
                                       0.0328
                             all
```

test_10152130138_Unigram_500.txt:

```
test1.txt ~
runid
                            all
                                      10152130138_dingwanning_unigram
num_q
                            all
                                      50
                            all
                                      50000
num_ret
                                      9805
num_rel
                            all
                                      1667
num_rel_ret
                            all
                            all
                                      0.0202
map
                                      0.0070
gm_map
                            all
                                      0.0664
Rprec
                            all
                                      0.0764
bpref
                            all
recip_rank
                            all
                                      0.1060
 iprec_at_recall_0.00
                            all
                                      0.1813
iprec at recall 0.10
iprec at recall 0.20
                                      0.0866
                            all
                            all
                                      0.0571
iprec_at_recall_0.30
                            all
                                      0.0173
 iprec_at_recall_0.40
                            all
                                      0.0044
 iprec_at_recall 0.50
                                     0.0005
                            all
 iprec_at_recall_0.60
                            all
                                      0.0000
iprec_at_recall_0.70
iprec_at_recall_0.80
                                      0.0000
                            all
                                      0.0000
                            all
                                      0.0000
 iprec_at_recall_0.90
                            all
 iprec_at_recall_1.00
                            all
                                      0.0000
P_5
                            all
                                      0.0400
P_10
P_15
                                      0.0340
                            all
                            all
                                      0.0493
P_20
                            all
                                      0.0510
P_30
                            all
                                      0.0520
P_100
                                      0.0790
                            all
P_200
                            all
                                      0.0802
P_500
                                      0.0549
                            all
P_1000
                                      0.0333
                            all
```

test_10152130138_Unigram.txt

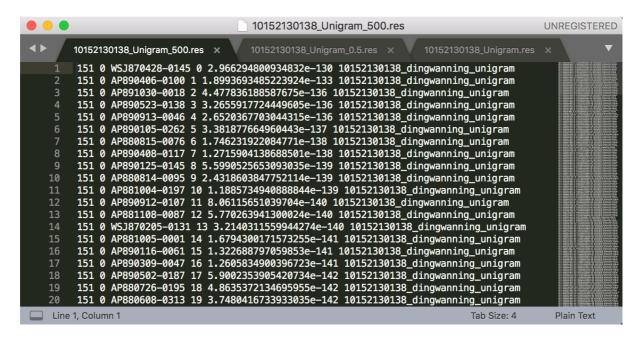
```
est1.txt ~
 runid
                            all
                                     10152130138 dingwanning unigram
                                     50
                            all
num_q
                                     25000
num_ret
                            all
num_rel
                            all
                                     9805
num_rel_ret
                           all
                                     1286
                                     0.0252
map
                            all
 gm_map
                            all
                                     0.0082
                                     0.0869
 Rprec
                           all
                                     0.0772
 bpref
                            all
 recip_rank
                            all
                                     0.1388
 iprec_at_recall_0.00
                            all
                                     0.2462
 iprec_at_recall_0.10
iprec_at_recall_0.20
                           all
                                     0.1110
                                     0.0617
                           all
 iprec_at_recall_0.30
                            all
                                     0.0115
 iprec_at_recall_0.40
                            all
                                     0.0009
 iprec_at_recall_0.50
                           all
                                     0.0000
                                     0.0000
 iprec_at_recall_0.60
                            all
 iprec_at_recall_0.70
                            all
                                     0.0000
 iprec_at_recall_0.80
                            all
                                     0.0000
                                     0.0000
 iprec_at_recall_0.90
                            all
                                     0.0000
 iprec_at_recall_1.00
                            all
                            all
                                     0.0640
P_10
P_15
                                     0.0740
                            all
                           all
                                     0.0960
P_20
P_30
                            all
                                     0.1020
                            all
                                     0.1147
P_100
                           all
                                     0.1224
P_200
                                     0.0957
                           all
P_500
                           all
                                     0.0514
P_1000
                                     0.0257
                            all
```

三. 实验结果

10152130138_Unigram_0.5.res:

```
10152130138 Unigram 0.5.res
                                                                                                                          UNREGISTERED
     10152130138_Unigram_0.5.res ×
      151 0 WSJ870428-0145 0 1.3008190732032897e-125 10152130138_dingwanning_unigram
      151 0 AP890406-0100 1 5.896700272629206e-126 10152130138_dingwanning_unigram
      151 0 AP890523-0138 2 3.553181629073012e-126 10152130138_dingwanning_unigram 151 0 AP890125-0145 3 7.531226362386954e-127 10152130138_dingwanning_unigram
      151 0 AP890913-0046 4 1.9074778899955163e-127 10152130138_dingwanning_unigram 151 0 AP891030-0018 5 8.144039778966308e-128 10152130138_dingwanning_unigram
      151 0 AP890105-0262 6 1.4636779391863454e-128 10152130138_dingwanning_unigram
151 0 AP890408-0117 7 1.0173370520313563e-128 10152130138_dingwanning_unigram
151 0 AP890408-0070 8 5.776391485679636e-129 10152130138_dingwanning_unigram
      151 0 AP880815-0076 9 4.268883393873505e-129 10152130138_dingwanning_unigram
      151 0 AP881004-0197 10 2.349420486325124e-129 10152130138_dingwanning_unigram
      151 0 AP890116-0061 11 2.2997460463828658e-129 10152130138_dingwanning_unigram
      151 0 AP890105-0166 12 1.5516801742835816e-129 10152130138_dingwanning_unigram 151 0 AP890411-0092 13 1.43831556843307e-129 10152130138_dingwanning_unigram
14
      151 0 AP880608-0313 14 1.4145921081424862e-129 10152130138_dingwanning_unigram
      151 0 AP881005-0001 15 1.2120340590503986e-129 10152130138_dingwanning_unigram
      151 0 AP881108-0087 16 8.0572678265047e-130 10152130138_dingwanning_unigram
      151 0 AP890309-0047 17 4.245453886030364e-130 10152130138_dingwanning_unigram
      151 0 AP881002-0015 18 2.9974456071351656e-130 10152130138_dingwanning_unigram 151 0 AP890912-0107 19 1.6730402843477924e-130 10152130138_dingwanning_unigram
20
Line 1, Column 1
                                                                                                        Tab Size: 4
                                                                                                                             Plain Text
```

10152130138_Unigram_500.res:



10152130138_Unigram.res:

