

文本相似度-bm25算法原理及实现

Python Web与Django大咖之路 2019-11-18

点击上方 “Python知音阁” 关注我!



自从Elasticsearch升级至6.X后，原本的相关性评分算法从TF-IDF（词频-逆文档频率）变为了BM25，那么BM25到底有什么过人之处，下面就来探索。

原理

BM25算法，通常用来作搜索相关性平分。一句话概况其主要思想：对Query进行语素解析，生成语素 q_i ；然后，对于每个搜索结果D，计算每个语素 q_i 与D的相关性得分，最后，将 q_i 相对于D的相关性得分进行加权求和，从而得到Query与D的相关性得分。

BM25算法的一般性公式如下：

$$Score(Q, d) = \sum_i^n W_i \cdot R(q_i, d)$$

其中，Q表示Query， q_i 表示Q解析之后的一个语素（对中文而言，我们可以把对Query的分词作为语素分析，每个词看成语素 q_i 。）；d表示一个搜索结果文档； W_i 表示语素 q_i 的权重； $R(q_i, d)$ 表示语素 q_i 与文档d的相关性得分。

下面我们来看如何定义Wi。判断一个词与一个文档的相关性的权重，方法有多种，较常用的是IDF。这里以IDF为例，公式如下：

$$IDF(q_i) = \log \frac{N - n(q_i) + 0.5}{n(q_i) + 0.5},$$

其中，N为索引中的全部文档数，n(qi)为包含了qi的文档数。

根据IDF的定义可以看出，对于给定的文档集合，包含了qi的文档数越多，qi的权重则越低。也就是说，当很多文档都包含了qi时，qi的区分度就不高，因此使用qi来判断相关性时的重要度就较低。

我们再来看语素qi与文档d的相关性得分R (qi, d) 。首先来看BM25中相关性得分的一般形式：

$$R(q_i, d) = \frac{f_i \cdot (k_1 + 1)}{f_i + K} \cdot \frac{qf_i \cdot (k_2 + 1)}{qf_i + k_2}$$

$$K = k_1 \cdot (1 - b + b \cdot \frac{dl}{avgdl})$$

其中，k1, k2, b为调节因子，通常根据经验设置，一般k1=2, b=0.75；fi为qi在d中的出现频率，qfi为qi在Query中的出现频率。dl为文档d的长度，avgdl为所有文档的平均长度。由于绝大部分情况下，qi在Query中只会出现一次，即qfi=1，因此公式可以简化为：

$$R(q_i, d) = \frac{f_i \cdot (k_1 + 1)}{f_i + K}$$

从K的定义中可以看到，参数b的作用是调整文档长度对相关性影响的大小。b越大，文档长度的对相关性得分的影响越大，反之越小。而文档的相对长度越长，K值将越大，则相关性得分会越小。这可以理解为，当文档较长时，包含qi的机会越大，因此，同等fi的情况下，长文档与qi的相关性应该比短文档与qi的相关性弱。

综上，BM25算法的相关性得分公式可总结为：

$$Score(Q, d) = \sum_i^n IDF(q_i) \cdot \frac{f_i \cdot (k_1 + 1)}{f_i + k_1 \cdot (1 - b + b \cdot \frac{dl}{avgdl})}$$

从BM25的公式可以看到，通过使用不同的语素分析方法、语素权重判定方法，以及语素与文档的相关性判定方法，我们可以衍生出不同的搜索相关性得分计算方法，这就为我们设计算法提供了较大的灵活性。

代码实现

```
1 import math
2 import jieba
3 from utils import utils
4
5 # 测试文本
6 text = '''
7 自然语言处理是计算机科学领域与人工智能领域中的一个重要方向。
8 它研究能实现人与计算机之间用自然语言进行有效通信的各种理论和方法。
9 自然语言处理是一门融语言学、计算机科学、数学于一体的科学。
10 因此，这一领域的研究将涉及自然语言，即人们日常使用的语言，
11 所以它与语言学的研究有着密切的联系，但又有重要的区别。
12 自然语言处理并不是一般地研究自然语言，
13 而在于研制能有效地实现自然语言通信的计算机系统，
14 特别是其中的软件系统。因而它是计算机科学的一部分。
15 '''
16
17 class BM25(object):
18
19     def __init__(self, docs):
20         self.D = len(docs)
21         self.avgdl = sum([len(doc)+0.0 for doc in docs]) / self.D
22         self.docs = docs
23         self.f = [] # 列表的每一个元素是一个dict, dict存储着一个文档中每个词的出现
24         self.df = {} # 存储每个词及出现了该词的文档数量
25         self.idf = {} # 存储每个词.idf值
26         self.k1 = 1.5
27         self.b = 0.75
28         self.init()
29
30     def init(self):
31         for doc in self.docs:
32             tmp = {}
33             for word in doc:
34                 tmp[word] = tmp.get(word, 0) + 1 # 存储每个文档中每个词的出现次
35             self.f.append(tmp)
36             for k in tmp.keys():
37                 self.df[k] = self.df.get(k, 0) + 1
38         for k, v in self.df.items():
```

```
39         self.idf[k] = math.log(self.D-v+0.5)-math.log(v+0.5)
40
41     def sim(self, doc, index):
42         score = 0
43         for word in doc:
44             if word not in self.f[index]:
45                 continue
46             d = len(self.docs[index])
47             score += (self.idf[word]*self.f[index][word]*(self.k1+1)
48                     / (self.f[index][word]+self.k1*(1-self.b+self.b*d
49                     / self.avgdl)))
50         return score
51
52     def simall(self, doc):
53         scores = []
54         for index in range(self.D):
55             score = self.sim(doc, index)
56             scores.append(score)
57         return scores
58
59 if __name__ == '__main__':
60     sents = utils.get_sentences(text)
61     doc = []
62     for sent in sents:
63         words = list(jieba.cut(sent))
64         words = utils.filter_stop(words)
65         doc.append(words)
66     print(doc)
67     s = BM25(doc)
68     print(s.f)
69     print(s.idf)
70     print(s.simall(['自然语言', '计算机科学', '领域', '人工智能', '领域']))
```

分段再分词结果

```
1  [['自然语言', '计算机科学', '领域', '人工智能', '领域', '中', '一个', '方向'],
2  ['研究', '人', '计算机', '之间', '自然语言', '通信', '理论', '方法'],
3  ['自然语言', '一门', '融', '语言学', '计算机科学', '数学', '一体', '科学'],
4  []]
```

```

5  ['这一', '领域', '研究', '涉及', '自然语言'],
6  ['日常', '语言'],
7  ['语言学', '研究'],
8  ['区别'],
9  ['自然语言', '研究', '自然语言'],
10 ['在于', '研制', '自然语言', '通信', '计算机系统'],
11 ['特别', '软件系统'],
12 ['计算机科学', '一部分']]

```

s.f

列表的每一个元素是一个dict，dict存储着一个文档中每个词的出现次数

```

1  [{'中': 1, '计算机科学': 1, '领域': 2, '一个': 1, '人工智能': 1, '方向': 1, '自然
2  {'之间': 1, '方法': 1, '理论': 1, '通信': 1, '计算机': 1, '人': 1, '研究': 1, '自
3  {'融': 1, '一门': 1, '一体': 1, '数学': 1, '科学': 1, '计算机科学': 1, '语言学':
4  {}},
5  {'领域': 1, '这一': 1, '涉及': 1, '研究': 1, '自然语言': 1},
6  {'日常': 1, '语言': 1},
7  {'语言学': 1, '研究': 1},
8  {'区别': 1},
9  {'研究': 1, '自然语言': 2},
10 {'通信': 1, '计算机系统': 1, '研制': 1, '在于': 1, '自然语言': 1},
11 {'软件系统': 1, '特别': 1},
12 {'一部分': 1, '计算机科学': 1}]

```

s.df

存储每个词及出现了该词的文档数量

```

1  {'在于': 1, '人工智能': 1, '语言': 1, '领域': 2, '融': 1, '日常': 1, '人': 1, '这

```

s.idf

存储每个词的idf值

```

1  {'在于': 2.0368819272610397, '一部分': 2.0368819272610397, '一个': 2.0368819272

```

s.simall(['自然语言', '计算机科学', '领域', '人工智能', '领域'])

```
1  ['自然语言', '计算机科学', '领域', '人工智能', '领域']与每一句的相似度[5.07699198143,
2
```

END



喜欢此内容的人还喜欢

Django 开发学习笔记

Python Web与Django大咖之路

北大医生：你们用焦虑养出来的娃，最后都送到我这里了

诗评万象

查理·芒格震撼演讲：如何低成本地从错误中爬出来？

金融精读