

文本相似度计算基础算法BM25

原创 纳纳 网络安全探秘 2019-09-28

接着文本相似度计算基础算法（一），继续介绍BM25算法。

BM25是基于TF-IDF算法的改进方法，也是一种统计学的方法，算法简单易懂。相似度计算方法如下：

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

其中， d 表示相关文档， Q 是query查询中的所有的词， q_i 是查询语句中的一个关键词词， W_i 是这个词的权重， $R(q_i, d)$ 是这个词和文档的相关度值的查询权重， n 是 Q 的词总数。

W_i 默认是IDF的值，上一节已经介绍过，计算公式如下：

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

其中， N 表示所有文档数， $n(q_i)$ 表示包含查询词 q_i 的文档数，0.5是避免 $n(q_i)$ 为0的情况。大致的意思是关键词出现频率越小说明越重要。

$R(q_i, d)$ 的计算公式如下：

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

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

其中 k_1, k_2, b 都是调节因子，一般 $k_1=2, k_2=1, b=0.75$ ，其中 f_i 表示关键字出现的次数， dl 是文档的长度， $avgdl$ 是文档的平均长度（为了防止 K 过大，归一化处理）， qfi 是关键字在 $query$ 查询 Q 中出现的次数。

$R(q_i, d)$ 乘积左边表示关键词在文档中的关系，右边表示关键词在查询语句中的关系。一般情况下，关键词在查询语句只会出现一次，所以 qfi 可以看成1，公式右边经过简化后也等于1，忽略不计。

相似度计算总公式经过分析后，最终计算式如下：

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

从上面公式可以看出，影响BM25公式的因数有：

1. idf，idf得分越高，分值越高。
2. tf，tf得分越高，分值越高。
3. dl，文档长度越长，分值越低。

喜欢此内容的人还喜欢

医疗保险入门知识

网络安全探秘

今年，你一定是朋友圈最好运爆棚的人。

新世相

2021，一个女人最好的生活状态（25-55岁必看）

佳人