

Part I

深度学习

Chapter 1

自然语言处理 (Natural language processing, NLP)

1.1 tf-idf

在一份给定的文件里，词频（term frequency, tf）指的是某一个给定的词语在该文件中出现的频率。这个数字是对词数（term count）的标准化，以防止它偏向长的文件。（同一个词语在长文件里可能会比短文件有更高的次数，而不管该词语重要与否）对于在某一个特定文件里的词语 t_i 来说，它的重要性可表示：

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}} \quad (1.1)$$

式中假设文件 d_j 中共有 k 个词语， $n_{k,j}$ 是 t_k 在文件 d_j 中出现的次数。分子 $n_{i,j}$ 是该词在文件 d_j 中出现次数，而分母则是在文件 d_j 中所有字词出现的次数之和。

逆向文件频率（inverse document frequency, idf）是一个词语普遍重要性的度量。某一特定词语的 idf，可以有总文件数目除以包含该词语之文件的数目，再将得到的上取以10为底的对数得到：

$$idf_i = \log_{10} \frac{|D|}{|\{j : t_i \in d_j\}|} \quad (1.2)$$

式中， $|D|$ 表示语料库中的文件总数， $|\{j : t_i \in d_j\}|$ 表示包含词语 t_i 的文件数目，如果词语不再资料中，就导致分母为零，因此一般情况下使用 $1 + |\{j :$

$t_i \in d_j\}$ ，然后

$$tfidf_{i,j} = tf_{i,j} \times idf_i \quad (1.3)$$

某一特定文件内的高词语频率，以及该词语在整个文件集合中的低文件频率，可以产生出高权重的 $tf-idf$ ，因此， $tf-idf$ 倾向于过滤掉常见的词语，保留重要的词语。

例子：假如一篇文件的总词数是100个，而词语母牛出现了3次，那么母牛一次在该文件中的词频就 $3/100 = 0.03$ ，而计算文件频率的方法就是以文件集的文件总数除以出现母牛一词的文件数，所以，如果母牛一次在1000份文件出现过，而文件总数是10000000份的话，其你想文件频率就是 $\log_{10}(10000000/1000) = 4$ ，最后的 $tf-idf$ 的分数为 $0.03 * 4 = .012$ ，更多访问 [tf-idf](#)。

text