

数据挖掘-TF-IDF

概述

一种统计方法，TF意思是词频，IDF意思是逆文本频率指数，用以评估一字词对于一个文件集或一个语料库中的其中一份文件的重要程度。字词的重要性随着它在文件中出现的次数成正比增加，但同时会随着它在语料库中出现的频率成反比下降。

TF-IDF值越高，说明该词越重要。

- 概述
- 输入/输出
- 参数设置
- 示例

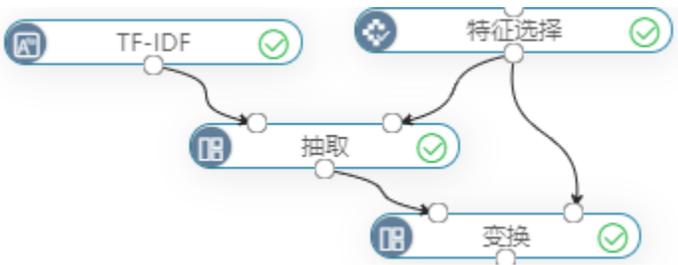
输入/输出

输入	没有输入端口
输出	一个输出端口，与抽取、变换节点组合使用。

参数设置

参数名称	说明	备注
特征项数	输入的数值n，代表算法最终会筛选出TF-IDF值最高的n个词	必填

示例



效果

使用“垃圾短信识别”数据，选择分词后的文本列_c2_seg_words_filtered，设置的特征项数为10，进行统计TF-IDF，输出结果如下图：

当前显示 100 条 / 总共有 295755 条数据 提示:点击单元格可查看超出的内容

#_c2_seg_words_filtered	Ab_c2_seg_words_filtered_tf	Ab_c2_seg_words_filtered_idf
减,拉...	(10,[0,1,3,5,6],[1.0,1.0,1.0,1.0,2.0])	(10,[0,1,3,5,6],[0.48142663540005...])
镇,市...	WrappedArray([乌兰察,布丰镇,市,...])	(10,[2,3,6,7,8,9],[1.0697945618004...])
保障...	WrappedArray([服务,保障,一路,建...])	(10,[0,2,5,7,9],[0.48142663540005...])
去,不...	WrappedArray([predictionio])	(10,[7],[0.5103782829293275])
楼,那...	WrappedArray([一层,楼,多个,ludo...])	(10,[1,2,3,4,6,7,9],[1.44203175015...])
我,我...	WrappedArray([南京,暴雨,四川,享...])	(10,[0,3,4],[0.9628532708001165...])
好,你...	WrappedArray([保卫,强奸,全家,权...])	(10,[0,1,5],[0.48142663540005826...])
由,球...	WrappedArray([美,职,篮,自由,球...])	(10,[0,1,7,8,9],[1.44427990620017...])
无,眼...	WrappedArray([晶体,全,净,无眼...])	(10,[2,4,5,6],[0.534897280900208...])
滴,是...	WrappedArray([距离,30,70公里,间])	(10,[1,4,6,9],[0.480677250052393...])
的,形...	WrappedArray([服务,形式,提供,wi...])	(10,[1,2,8],[0.961354500104787,0...])

注意:表头中#表示特征列,*表示标签列

表头真名 表头别名