

资讯

指数

活动

创投罗盘

综合 | 专家 | 专访 | 学堂 | 报告 | 案例 | **大数据+** | 能源 | 环保 | 营销 | 金融 | 征信 | 医疗 | 零售 | 交通 | 通信 | 互娱 | 农业 | 政府

文本挖掘小探索：避孕药都有哪些内容主题

冯大福 | 2017-08-31 10:29

挖掘 避孕 内容 主题

【数据猿导读】舆情监测一直是众多品牌关注的地方，尤其品牌想知道在品牌推广，品牌策略，品牌广告中出现的问题，从而能进行策略上的改进，但是现在很多人都是读帖子，笔者在4年前做舆情分析时候就是读帖子



舆情监测一直是众多品牌关注的地方，尤其品牌想知道在品牌推广，品牌策略，品牌广告中出现的问题，从而能进行策略上的改进，但是现在很多人都是读帖子，笔者在4年前做舆情分析时候就是读帖子，至今没有太多改善，关注舆情监测中的主题挖掘部分，主题挖掘可以使数据分析师，减轻工作量，去掉读帖子等一系列等的复杂工作，大致了解主题规律。

本文是笔者早前发在某网站上的，由于笔者最近太忙，将本文修改下呈现给大家：

本文分析逻辑：

数据处理

数据探查

主题挖掘

数据处理

1.数据源：

从各大网站论坛，微博等爬虫关于某避孕药的内容

关键字段名称包含

content Author: 发帖作者(第D列)

Content Forward: 转发的内容(第F列)

Content_Main: 发帖内容(第G列)

Title：发帖内容(第H列)

精品栏目

[2017/07/27] 大数据24小时 [More>](#)

[2017/07/24-28] 大数据周周看 [More>](#)

[2017/07/24-28] 大数据投融资 [More>](#)

[2017/07/24-28] 大咖周语录 [More>](#)

[2017/07/24-28] 大数据周聘汇 [More>](#)

[2017/07/24-28] 每周一本书 [More>](#)

[2016/08-10] 大数据活动公告 [More>](#)

专家推荐

[More >](#)



涂子沛



傅志华



马亮



苏萌



崔晓波



韩涵



车品觉



刘雷鸣



董飞



郭伟



张涵诚



陈运文

人物专访

[More >](#)



曾被167家VC拒绝，如今公司估值百

活动推荐

[More >](#)

ACS 2017中国汽车CIO峰会10 2017-10-25

2017金融科技价值—数据驱动 2017-10-19

其他字段和本文不想关，不阐述

2.加载数据包(r语言)和需要在中文分词中插入的中文词语：

Rwordseg：(4年前用的分词包，不知道现在更新与否)，分词包就是让R语言识别中文，按照单词来视为一个值

插入单词：因为Rwordseg中文词性包含不了其他奇怪词汇，例如： 妈富隆、优思明、短期避孕药、治疗多囊等。插入单词作为模型的变量值

```
insertWords("毓婷")
insertWords("短期避孕药")
insertWords("多囊")
insertWords("治疗多囊")
insertWords("长期避孕药")
insertWords("调经")
insertWords("紧急避孕药")
insertWords("口服避孕药")
insertWords("Yasmin")
insertWords("Yas")
insertWords("Yaz")
insertWords("短效避孕药")
```

3.读入文本分析处理

去掉数字、特殊字符、标准符号

数据探索：大概了解下数据现状

1.根据变量值(单词)统计各个单词出现的次数

2017第二届中国国际大数据产业博览会	2017-08-17
GIEC2017全球互联网经济大会	2017-08-08
2017年第二届上海大数据与分享大会	2017-08-01

不容错过的资讯

- 1 大数据24小时：Facebook“神童”跳槽谷歌
- 2 金融科技&大数据产品推荐：神策分析—
- 3 分享：解析6个公司的大数据岗位的面试经
- 4 小白做数据分析的一点感悟
- 5 机器人即将抢走你的工作？数据表明你可
- 6 大数据贵在应用 或将成为中国弯道超车的
- 7 如何让大数据分析更有效？这里有5种技
- 8 大数据是什么？一文秒读懂大数据
- 9 如何将数据可视化技术应用于广告投放？
- 10 大数据将把人类带进怎样的新世界？

大数据学堂

More >



【每周一本书】之《深度学习算法

大数据企业推荐

More >



九次方 | 贡献中国数据智慧



星图数据 | Data turn biz



晶赞科技 | 数据推动产业智能化



TalkingData | 移动·数据·价



百分点 | 大数据践行者

热门职位

More >

避孕药	避孕	口服避孕药	服用	月经	女性	可以
5320	3571	3154	3066	2565	2542	2064
优思明	使用	作用	没有	激素	孕激素	紧急避孕药
2055	1782	1435	1434	1424	1423	1373
副作用	如果	可能	药物	周期	影响	雌激素
1359	1339	1260	1201	1198	1151	1136
卵巢	子宫	安全	一个	排卵	怀孕	增加
1121	1046	1033	1024	985	943	935
妇女	效果	开始	这个	方法	问题	目前
919	907	898	892	849	839	810
身体	发胖	出血	PostsEnd	医生	复方	健康
804	799	787	778	778	761	758
但是	自己	治疗	体重	导致	时间	因为
751	744	743	738	736	732	731
降低	发生	引起	出现	剂量	很多	我们
727	715	715	699	685	676	671
长期	所以	服药	妈富隆	什么	生育	那么
644	641	627	612	608	595	592
主要	风险	减少	建议	一种	一些	妊娠
591	584	579	559	555	542	541

2.根据单词量画词云图



3.重新转化用于聚类的数据格式

根据以上数据探索的词频，词作为colname，词频表示数值，每一行是帖子内容作为id标示

例如：

即每个帖子出现了某词的词频的次数，帖子1中出现避孕药2次，优思明4次，囊中1次

R语言tm包来作处理

即：分词之后生成一个列表变量，用列表变量构建语料库。

由于tm包中的停用词()都是英文(可以输入stopwords()查看)，所以大家可以去网上查找中文的停用词，用removeWords函数去除语料库中的停用词：

生成语料库之后，生成词项-文档矩阵(Term Document Matrix，TDM)，顾名思义，TDM是一个矩阵，矩阵的列对应语料库中所有的文档，矩阵的行对应所有文档中抽取的词项，该矩阵中，一个[i,j]位置的元素代表词项i在文档j中出现的次数。

4.注意：

北京 | 数据堂 大数据架构师&大数据分析

北京 | 聚合数据 业务拓展经理&JAVA工程

北京 | 慧米数据 三个职位

湖南 | 银杏数据科技有限公司 数据工程师

北京 | 中献电子技术开发中心 大数据分析

大家都在搜

阿里巴巴	医疗	互联网	
机器人	漏洞	科技部	融资
小米	数据挖掘	金融	
物联网	大数据应用	追随	
机器学习	陕西	春节	营销
百度	python	投资	大数据
互联网+	电商	大数据	
数据分析	中国	创业	
云计算	人工智能	北京	

默认的加权方式是TF，即词频，这里采用Tf-Idf，该方法用于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：

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次数。这个数字通常会被归一化，以防止它偏向长的文件。

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，可以由总文件数目除以包含该词语之文件的数目，再将得到的商取对数得到。

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倾向于保留文档中较为特别的词语，过滤常用词。

同时，需要用removeSparseTerms()函数进行降维

数据挖掘

1.查看频率&基本统计

其实就是在数据挖掘查看数据基本统计，目的看下数据逻辑符合不符合社会认知

2.LDA

LDA(Latent Dirichlet Allocation)是一种文档主题生成模型，也称为一个三层贝叶斯概率模型，包含词、主题和文档三层结构。所谓生成模型，就是说，我们认为一篇文章的每个词都是通过“以一定概率选择了某个主题，并从这个主题中以一定概率选择某个词语”

具体的算法核心在这里略，因为写太多可能读者看不懂。

VEM = LDA(sample.dtm2, k=10, control = list(seed= SEED)),

根据上面所形成的raw data形成主题词：

解读：

主题分类10个主题，在这个文本中，但是还需要优化

第一个主题是女性服用避孕药后作用会不会增加流产风险

第二个主题是女性避孕和激素的关系(需要优化)

第三个主题医生推荐优思明会不会影响月经(需要优化)

第四个主题口服避孕药会不会影响月经

第五个主题治疗痤疮，激素，多囊

第六个主题和第四个主题一样(需要优化)

第七个主题同上

第八个主题，杂文帖子

第九个主题，会不会是吃完避孕药后发胖

第十个主题，优思明女性服用避孕

确切来说，这10个主题还需要优化，文本经过人工看完应该提炼的的是优思明使用目的，大部分集中在避孕，安全，发胖，治疗痤疮等，少部分会集中副作用等。

结束语

由于4年前做脚本，因此好多需要优化，之后会将优化的和大家分享

优化内容包含

需要在文本库中添加月经不调，治疗痤疮，青春痘等词语

文本还需要继续处理改进去掉postend

主题数目需要加大

以及主题内容维度需要增加可以让他成为一句话

来源：36大数据

收藏

分享

声明：数据猿尊重媒体行业规范，相关内容都会注明来源与作者；转载我们原创内容时，也请务必注明“来源：数据猿”与作者名称，否则将会受到数据猿追责。

相关文章

刷新

七夕后大数据：CBNData发布线上避 如何更好地阅读很多数学相关内容的【大咖周语录】深耕数据挖掘，实现孕套十大品牌 机器学习论文？ 互联互通，在银行风控管理中 ...

我要评论

我想要评论.....

提交评论

ok

热点导航

大数据人物专访	大数据活动推荐
大数据学堂	商业智能
互联网广告	央行征信
检察系统	人工智能
内容为王	宽带资本
硅谷大数据	通联数据
京东金融	二次元大数据
信息安全	大数据风控
大数据研究基地	原生数据
大数据地形图	位置大数据
互联网	人工智能
大数据技术	快递
大数据入门	网站地图

关于数据猿
成为专栏专家
好文投递&寻求报道
广告推广与活动合作
数据支持&合作



数据合作伙伴：

