

文本挖掘系统介绍

Introduction of Text Mining System

李宁

<http://www.lining0806.com/>

文本挖掘

- 词频统计
- 词性标注
- 文本分类
- 标签提取
- 实体发现和识别
- 情感趋势分析

实现目标

网易新闻 新闻排行榜

网易首页 > 新闻中心 > 新闻排行榜

请输入关键词 新闻 搜索

最新 | 排行 | 国内 | 国际 | 社会 | 评论 | 深度 | 军事 | 历史 | 探索 | 图片 | 博客 | 媒体 | 视频 | 公益 | 手机版

新闻日历

快速跳转：

新闻

娱乐

体育

财经

科技

汽车

女人

房产

读书

游戏

旅游

教育

公益

校园

传媒

视频

移动

全站

图集排行榜

全站

更多

点击榜

24小时点击排行

本周点击排行

本月点击排行

跟贴榜

今日跟贴排行

本周跟贴排行

本月跟贴排行

标题	点击数	标题	跟贴数
1 心疼！刘亦菲宣传新戏遭疯狂粉丝推倒在地	4361873	1 女医生为男子检查私处遭暴打	236204
2 网传云南一副教授与女生不雅照 尺度相当大(图)	2250116	2 北京拦截电信诈骗资金10亿元	138828
3 王思聪林更新被曝在售楼处看房子 16万一平	1404072	3 郑州警方深夜突查涉黄五星会所带走多名女子	135977
4 俄男子晨穿警服去扫黄 妓女被逼裸体游街(图)	1304362	4 老人不顾劝阻备用水源地游泳	125966

苹果要在印度建立4000人的研发中心，专门开发苹果地图

2016年05月19日15:23 新浪科技 微博 我有话说 收藏本文

新浪科技讯 北京时间5月19日下午消息，为了进一步吸引印度用户和开发者，苹果CEO蒂姆·库克(Tim Cook)周四宣布在印度南部的海德巴拉建设一个新的研发中心，专门开发该公司的地图产品。

苹果周三早些时候宣布，将于明年年初在印度班加罗尔建设一处新的设施，专门帮助开发者采用最佳的开发模式，并改进iOS应用的设计、质量和性能。

库克目前正在进行对印度的第一次访问，而在iPhone全球销量和整个公司的营收双双下滑的背景下，苹果iPhone今年第一季度在该国实现56%的销量增长。

苹果的这个新研发中心将专注于为iPhone、iPad、MAC和Apple Watch等产品开发地图功能。苹果表示，这笔投资将加快苹果地图的开发速度，并在当地创造4000个就业岗位。

苹果并未披露这笔投资的具体规模，但有报道称可能达到2500万元。

科技？财经？娱乐？匹配关键词靠谱不？

科技：{电脑，互联网，…}

财经：{股票，黄金，白银，…}

娱乐：{游戏，音乐，电影，…}

问题1：关键词重复？

问题2：新生词语的出现？

问题3：“苹果公司开始进军游戏行业，股价进一步上涨。”

实现目标

地区

☐ 中国☐ 美国☐ 英国☐ 日本☐ 澳大利亚☐ 德国☐ 法国☐ 韩国☐ 西班牙☐ 意大利☐ 其他

类型

☐ 新闻☐ 经济数据

重要性

☐ 低☐ 中☐ 高

搜索

查看详情

开

短讯

中文

筛选关键字

白名单

黑名单

添加

内容

搜索

发布时间

暂无数据

	发布时间	内容		来源	操作			
	内容	类型	操作	来源	内容	类型	操作	
	耶伦	白名单	修改 删除	ioaba	办文凭	黑名单	修改 删除	
	日本	白名单	修改 删除		踏纳税	黑名单	修改 删除	
	日经	白名单	修改 删除		华门开	黑名单	修改 删除	
	道琼斯	白名单	修改 删除		帝国之梦	黑名单	修改 删除	
	原油	白名单	修改 删除		代办学	黑名单	修改 删除	
	外汇	白名单	修改 删除		骑单车出	黑名单	修改 删除	
	欧盟	白名单	修改 删除		专业助	黑名单	修改 删除	
	美股	白名单	修改 删除		纯度黄	黑名单	修改 删除	
	布伦特	白名单	修改 删除		换妻	黑名单	修改 删除	
					麻果配	黑名单	修改 删除	
				三秒倒	黑名单	修改 删除		

关键字:

输入关键字,若有多个关键字请用 ","(英文状态下) 分隔

标普

美股

日本

日股

日经

道琼斯

纳斯达克

黄金

现货黄金

贵金属

布伦特

原油

美元

美联储

外汇

欧元

欧央行

欧盟

撤销同步

置顶

实现流程

前提!

文本分词

文本过滤

关键词提取

文本自动分类

文本推荐

长度?
去重?

过滤词黑
名单

邮件实
时通知
系统

停用词,
关键词白
名单

特征?
模型?
参数?

基于
统计的
大数据

停用词,
特征词

基于
规则的
大数据

时间?
数字?
长度?

主词典文
件

李小福是创新办主任也是云计算方面的专家

李小福 / 是 / 创新 / 办 / 主任 / 也 / 是 / 云 / 计算 / 方面 / 的 / 专家 /

李小福 / 是 / 创新办 / 主任 / 也 / 是 / 云计算 / 方面 / 的 / 专家 /

分词白名
单

小结:

1. “自学习”功能
2. 扩展到其他应用

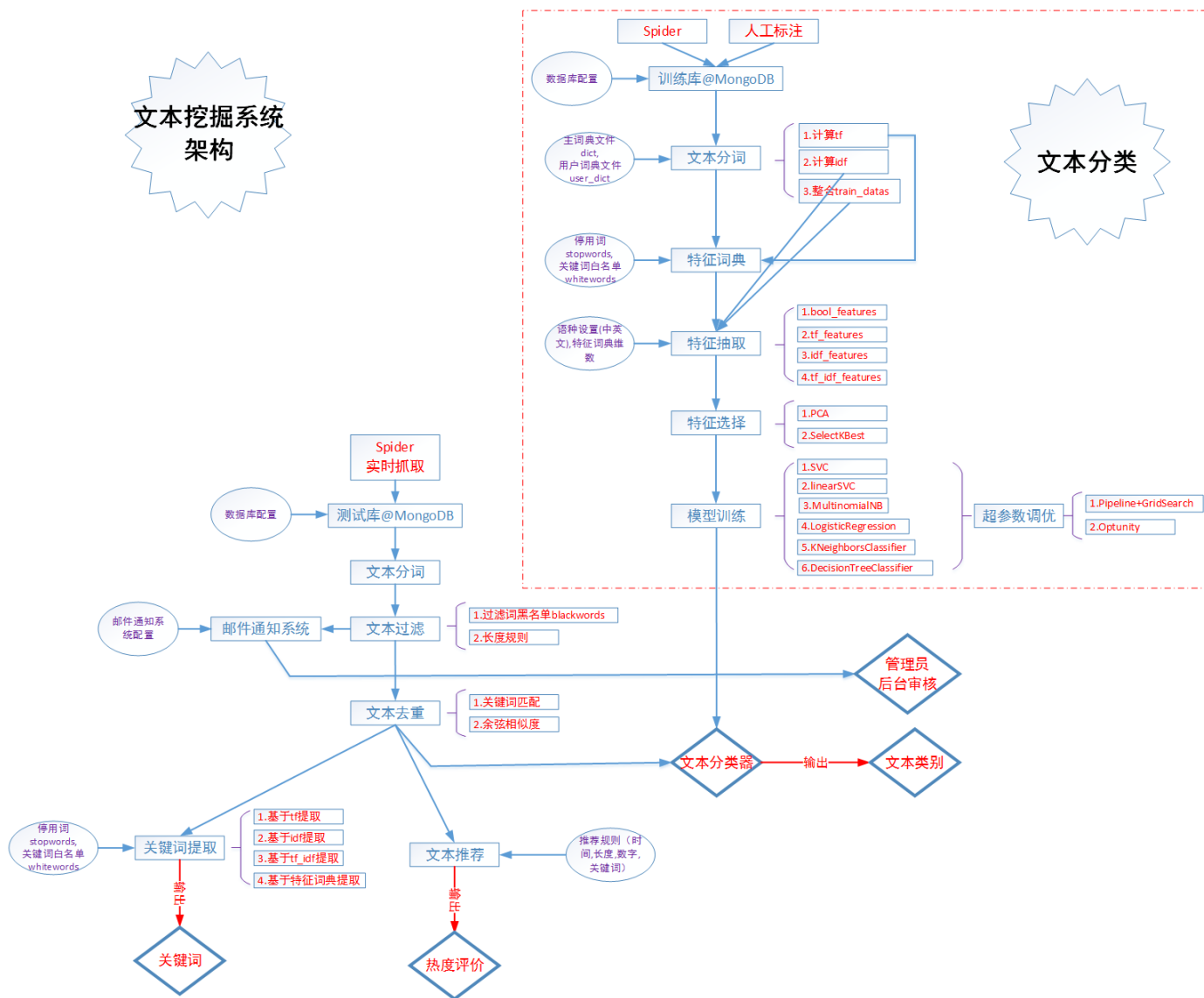
我们如何去做优化?

1. 分词白名单 (一般不搞)
2. 过滤词黑名单及过滤规则的总结
3. 停用词及关键词白名单
4. 分类统计的优化 (机器学习)
5. 推荐规则的总结

系统框架

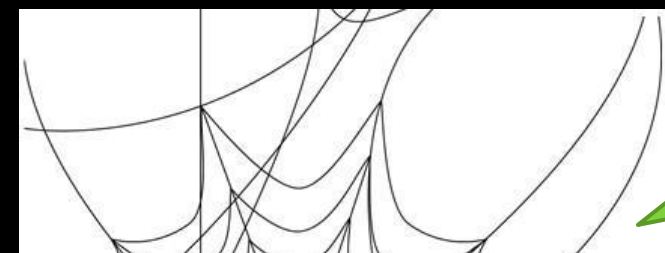
文本挖掘系统架构

文本分类



数据抓取——Spider

指的是通过程序实现访问某个URL地址，然后获得其所返回的内容（HTML源码，Json格式的字符串等）。然后通过解析规则，分析出我们需要的数据并取出来。



```
<html xmlns="http://www.w3.org/1999/xhtml">
<head>
  <meta http-equiv="Content-Type" content="text/html; charset=gb2312" />
  <meta http-equiv="Content-Language" content="zh-CN" />
  <title>新闻排行榜_网易新闻</title>
  <base target="_blank" />
  <meta name="keywords" content="" />
  <meta name="description" content="" />
  <meta name="author" content="网易" />
  <meta name="Copyright" content="网易版权所有" />
  <script>if(!/auto|house|home|bbs|blog/.test(location.host)&&!(document.documentElement&&document.documentElement.getAttribute("phone"))&&!/pc=1/.test(location.mobile|ipod|blackberry|bb|d+|phone/.i.test(navigator.userAgent)))document.write("<meta name='viewport' content='width=device-width, initial-scale=1, maximum-scale=1'"
  <div style="position: absolute; top: 50%; left: 0; width: 100%; height: 40px; margin-top: -40px; text-align: center; background: url(http://img1.cache.netease.com/utf8/endpoint/ima
  载中 ...</div></div><script src="http://img1.cache.netease.com/f2e/system/touchall/collect/foot~3Cwae6Po0Sne.js" + " defer"></script><plaintext style="display:none
  <script>var _ntes_const={stime: new Date()};</script>
  <link href="http://img1.cache.netease.com/cnews/css07/style.css" rel="stylesheet" type="text/css" />
  <link href="http://img1.cache.netease.com/cnews/img09/channel_nav.css" rel="stylesheet" type="text/css"/>
  <link href="http://img1.cache.netease.com/cnews/img/subscribe0304/rank.css" rel="stylesheet" type="text/css" />
  <style type="text/css">
    .gg735 {width:735px; overflow:hidden; float:left;}
    .gg210 {width:210px; overflow:hidden; float:right;}
    .channel h2 { width: auto; }
    .subNav a.photoset-icon { background: url(http://img1.cache.netease.com/cnews/img10/bbs0114/photoset.gif) no-repeat; width:128px; }
  </style>
  <script type="text/javascript" language="javascript" src="http://img1.cache.netease.com/cnews/js/ntes_islib_l.x.js" charset="gb2312"></script>

  <script language="javascript" type="text/javascript" src="http://202.102.100.100/35ff706fd57d11c141cdefcd58d6562b.js" charset="gb2312"></script><script type="text/j
  vvvvvvvvvv(' .btm-ad,.gg300,.content-ad,.gg,[class$='bottom_ad'],[class$='ad2'],.ggarea,.top-gg,#addiv,[class$='ad'],.top-gg-area,[class$='gg2'],[class$='ad ad-'],
  .youdao,iframe[src$='http://g.163.com/']&[src$='&affiliate='],iframe[src$='http://img1.126.net/'],iframe[src$='http://x.jd.com/'],img[src$='http://img1.126.net/'],if
  <body>
  <script>if(!/auto|house|home|bbs|blog/.test(location.host)&&!(document.documentElement&&document.documentElement.getAttribute("phone"))&&!/pc=1/.test(location.sear
  mobile|ipod|blackberry|bb|d+|phone/.i.test(navigator.userAgent)))document.write("<meta name='viewport' content='width=device-width, initial-scale=1, maximum-scale=1"
  <div style="position: absolute; top: 50%; left: 0; width: 100%; height: 40px; margin-top: -40px; text-align: center; background: url(http://img1.cache.netease.com/utf8/endpoint/ima
  载中 ...</div></div><script src="http://img1.cache.netease.com/f2e/system/touchall/collect/foot~3Cwae6Po0Sne.js" + " defer"></script><plaintext style="display:none
  <link type="text/css" rel="stylesheet" media="screen" href="http://img1.cache.netease.com/common/css/common_nav_v1.0.8.css" />
  <style>
```

数据抓取——Spider

取

存



使用Cookie登录

AJAX?
反盗链?
断线重连?
.....

终极武器:
Selenium!!

解析规则: 正则表达式, xpath, csspath,

数据抓取——Spider

微信公众号自动登录，抓取图文统计数据

——技术难点：

1.Selenium获取动态token

2.Api的时效性

指定关键词搜索微信文章并存入数据库

中文分词

分词就是将连续的字串或字符序列按照一定的规范重新组合成词序列的过程。

主要的分词方法

◆ 简单的模式匹配:

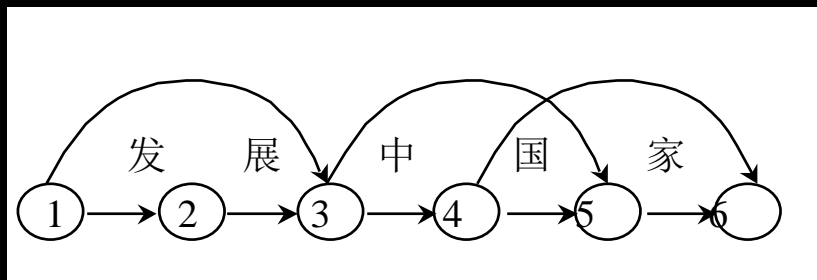
- 正向最大匹配
- 逆向最大匹配
- 双向匹配法

◆ 基于规则的方法:

- 最少分词算法

◆ 基于统计的方法:

- 统计语言模型分词
- 串频统计和词形匹配相结合
- 无词典分词



特征抽取

TF-IDF的主要思想

如果某个词或短语在一篇文章中出现的频率高，并且在其他文章中很少出现，则认为此词或者短语具有很好的类别区分能力，适合用来分类。

TF (Term Frequency): 词频

IDF (Inverse Document Frequency): 逆文档频率

$$tf_{i,j} = \frac{n_{i,j}}{\sum_k n_{k,j}}$$

$$idf_i = \log \frac{|D|}{|\{j : t_i \in d_j\}|}$$

$$tfidf_{i,j} = tf_{i,j} \times idf_i$$

特征选择

特征选择方法

◆ 去除无用特征：

- 不必要的特征对训练无用。

◆ 去除相关分量：

- 相关的多个分量可以变换成较少的不相关分量。

PCA (Principal Component Analysis)

- 主成分分析是设法将原来众多具有一定相关性（比如P个指标），重新组合成一组新的互相无关的综合指标来代替原来的指标。

假定有 n 个样本，每个样本共有 p 个变量，
构成一个 $n \times p$ 阶的数据矩阵

$$X = \begin{bmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \vdots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{bmatrix}$$

降维处理

$$\begin{cases} Z_1 = l_{11}X_1 + l_{12}X_2 + \cdots + l_{1p}X_p \\ Z_2 = l_{21}X_1 + l_{22}X_2 + \cdots + l_{2p}X_p \\ \vdots \\ Z_m = l_{m1}X_1 + l_{m2}X_2 + \cdots + l_{mp}X_p \end{cases}$$

分类器设计

朴素贝叶斯 (Naive Bayes)

Naive Bayes是一个生成模型，在计算 $P(y|x)$ 之前，先要从训练数据中计算 $P(x|y)$ 和 $P(y)$ 的概率，从而利用贝叶斯公式计算 $P(y|x)$ 。

Naive Bayes满足 $P(y=1|x) = P(y=1) * P(x|y=1) / p(x)$

采用使得后验概率 $P(y|x)$ 最大的输出，作为最佳的输出 y 。

$$\begin{aligned}P(X|Y) &= P(X_1, X_2|Y) \\&= P(X_1|X_2, Y)P(X_2|Y) \\&= P(X_1|Y)P(X_2|Y)\end{aligned}$$

$$P(X_1 \dots X_n|Y) = \prod_{i=1}^n P(X_i|Y)$$

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) P(X_1 \dots X_n | Y = y_k)}{\sum_j P(Y = y_j) P(X_1 \dots X_n | Y = y_j)}$$

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

$$Y \leftarrow \arg \max_{y_k} \frac{P(Y = y_k) \prod_i P(X_i | Y = y_k)}{\sum_j P(Y = y_j) \prod_i P(X_i | Y = y_j)}$$

$$Y \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i | Y = y_k)$$



文章相似度

余弦相似度

北京气象专家解释“泥雪”：长期无降水空气脏

金羊网 - 4小时前

两人合撑一把伞在雨中打车。昨天，京城迎来一场雨夹雪。记者陶冉摄。今天是春分节气，时中到大雪，而平原地区由于气温原因以雨夹雪为主。截至昨晚8点，城区 ...



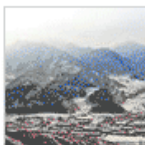
凤凰网



搜狐



每日甘肃



搜狐



腾讯网



北国网

北京暴雪清污染京城三月飘雪好预兆【组图】

www.591hx.com - 3小时前

飞雪迎春袭北京京城今晨或现“堵城”

大洋网 - 3小时前

北京普降瑞雪银装素裹树挂景观成春日美景

艾拉家居网 - 7小时前

延庆迎春雪城区下泥雪专家称系内蒙古沙尘被卷来

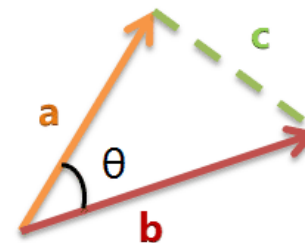
凤凰网 - 9小时前

昨夜北京普降大雪道路结冰早高峰注意出行安全

张家界在线 - 11小时前

北京春分降雪空气净化专家称三月下雪很正常

腾讯网 - 11小时前



$$\cos\theta = \frac{a^2 + b^2 - c^2}{2ab}$$

更多资源

项目链接: <https://github.com/lining0806/TextMining>

```
53     ## -----
54     ## 生成分类器模型
55     feature_selection_flag = False
56     my_selector = None
57     if test_speedup and os.path.exists(fea_dict_file) and os.path.exists(best_clf_file):
58         words_feature = []
59         with open(fea_dict_file, 'r') as fp:
60             for line in fp.readlines():
61                 word_feature = line.strip().decode("utf-8")
62                 words_feature.append(word_feature)
63         if feature_selection_flag:
64             with open(best_clf_file, "rb") as fp_pickle:
65                 my_selector, best_clf = pickle.load(fp_pickle)
66         else:
67             with open(best_clf_file, "rb") as fp_pickle:
68                 best_clf = pickle.load(fp_pickle)
69     else:
70         ## -----
71         words_feature = MakeFeatureWordsDict(all_words_tf_dict, stopwords_set, writewords_set, lag, fea_dict_size)
72         train_features = []
73         train_class = []
74         for train_data in train_datas:
75             TextFeatureClass = TextFeature(words_feature, train_data[0])
76             train_features.append(TextFeatureClass.TextBool()) ##### 可以调整特征抽取, 训练集与测试集保持一致
77             train_class.append(int(train_data[1])) # str转为int
78         train_features = np.array(train_features)
79         train_class = np.array(train_class)
80         if feature_selection_flag:
81             FeatureSelectorClass = FeatureSelector(train_features, train_class)
82             my_selector, train_features = FeatureSelectorClass.PCA_Selector() ##### 可以调整特征选择
83         start_time_train = datetime.datetime.now()
84         ClassifierTrainClass = ClassifierTrain(train_features, train_class)
85         best_clf = ClassifierTrainClass.LR() ##### 可以调整分类器训练
86         end_time_train = datetime.datetime.now()
87         print "best_clf training last time:", end_time_train-start_time_train
88         if not os.path.exists(Classifier_Dir):
89             os.makedirs(Classifier_Dir)
90         with open(fea_dict_file, 'w') as fp:
91             for word_feature in words_feature:
```

谢谢！