
目录

文档关键词的提取和文本相似度的计算.....	1
一、数据集的获取.....	1
数据集介绍	1
数据集格式	1
二、数据的清洗	2
三、分词	2
四、相似度 simhash 算法与相似度的计算.....	3
相似度 simhash 算法	3
相似度的计算.....	4
五、计算结果.....	4
相似度和索引 csv.....	5
新闻文本 txt.....	5
词云	6

文档关键词的提取和文本相似度的计算

一、数据集的获取

数据集介绍

数据集使用搜狐新闻数据(SogouCS)Ver.2012, 来自搜狐新闻 2012 年 6 月—7 月期间国内, 国际, 体育, 社会, 娱乐等 18 个频道的新闻数据, 提供 URL 和正文信息。

数据集格式

```
<doc>
<url>页面 URL</url>
<docno>页面 ID</docno>
<contenttitle>页面标题</contenttitle>
<content>页面内容</content>
</doc>
```

二、数据的清洗

因为语料数据没有采用 utf-8 编码，这里我们使用 gb18030 编码方式，gb18030 拥有更广的中文编码，能够更好地读取语料数据。

我们数据集中的新闻数据进行遍历，将<content>页面内容</content>筛选出来，同时去除两边的<content>标识符，共筛选得到 1,367,845 个新闻语料。

接着我们对内容长度过少的语料数据进行清洗，规则是文本长度≤100。

接着再对新闻文本中的格式字符和 Unicode 字符编码进行处理，如\n、\u3000和\u201c等特殊字符。

数据集语料数目	清洗后数据集语料数目
1,367,845	1,093,388

三、分词

分词我们选择 jieba 分词，停用词库我们选择了 stars 人数较多的中文常用停用词表(<https://github.com/goto456/stopwords>)，再根据我们的新闻与了数据对停用词表进行适当的增删，以达到良好的分词效果。

节选的一片分词效果如下。

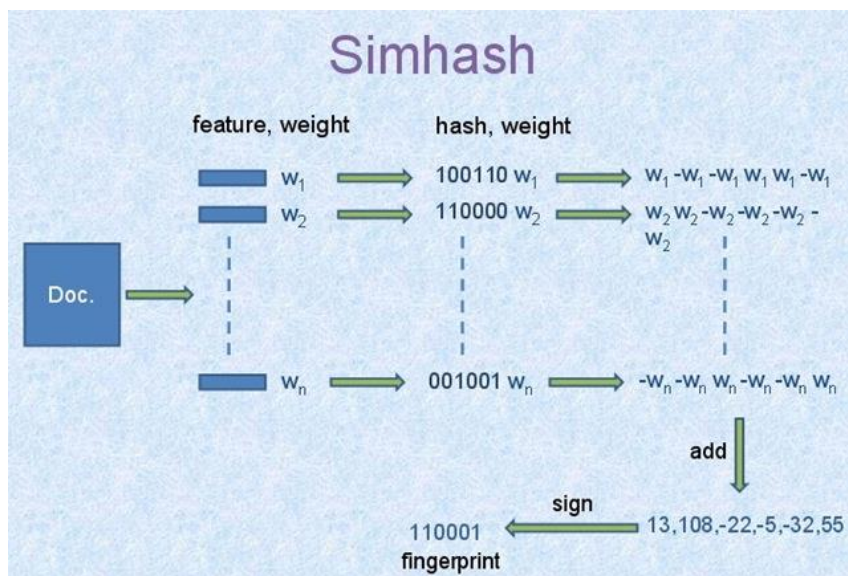
深圳市 法制办 网站 公布 深圳经济特区 社会 救助 条例 以下 简称 条例 再次 公开 征求 社会 意见 条例 增加 流浪 乞讨 人员 生活 救助 灾民 生活 救助 两章 内容 规定 救助站 应 乞讨 人员 提供 必需 食物 住处 医疗 救治 灾民 每人 每天 领取 元 基本 生活费 投靠 亲友 解决 住宿 每人每天 领取 元 补助 职业 行乞 进行 救助 教育 劝返 异地 分流 安置 条例 规定 流浪 乞讨 人员 实行 分类 救助 管理 无力解决 食宿 亲友 投靠 享受 最低 生活 保障 正在 流浪 乞讨 度 日 生活 着落 自愿 求助 流浪 乞讨 人员 列为 求助 流浪 乞讨 人员 救助 管理 部门 给予 救助 护送 购票 返乡 二 十八 岁 以下 未成年人 流浪 乞讨 行为 实行 保护性 救助 列为 未成年 流浪 乞讨 人员 此类 人员 及时 予以 进站 保护 三 乞讨 职业 流浪 乞讨 人员 列为 职业 流浪 乞讨 人员 救助 管理 部门 进行 救助 教育 劝返 异地 分流 安置 四 生命危 险 流浪 乞讨 危重病 精神病人 列为 患病 流浪 乞讨 人员 先 救治 救助 原则 实行 人道主义 救治 民政 城管 公安 送往 定点医院 救治 五 违法犯罪 流浪 乞讨 人员 列为 有害 流浪 乞讨 人员 公安部门 依法 严厉 处置 打击 一时 处理 不了 无法 返乡 人员 进行 身份 核实 严加 管理 逐一 甄别 滞留 站 送 异地 安置 家属 抚慰金 每名 遇难者 低于 元 条例 规 定 灾害 救助 采取 提供 救灾 生活 物资 基本 生活费 补助 临时 安置 住所 住宿 补助 医疗 救助 危机 心理咨询 服务 就业 服务 生产资料 救助 方式 自然灾害 救助 保障 灾民 基本 生活 需要 原则 困难 程度 每人每天 维持 基本 生活 费 用 二十元 天 标准 二 住宿 补助 邻里 借住 投靠亲友 方式 解决 住宿 五十元 天 标准 补助 三 衣被 救助 保证 灾民 有 衣 穿 天 冷 重点 救助 御寒衣 四 伤病 救助 伤病 家庭 困难 程度 适当 救助 五 遇难 人员 家属 抚慰金 每名 遇难者 低于 五千元 灾民 生活 特 困难 可予 临时 救助金 条例 提及 过渡性 安置 期 灾民 家庭 生活 特别 困难 人民政府 困难 程度 给予 一次性 临时 救助 补助金 本市 户籍 居民 持有 本市 居住证 深圳 连续 工作 生活 满 一年 申请 前 连续 本 市 缴纳 养老保险费 超过 一年 本市 非 户籍 居民 遭受 疾病 意外事故 诉讼 失踪 死亡 突发 情况 致使 基本 生活 暂时 出现 较大 困难 申请 临时 救助 符合 临时 救助 条件 居民 遭遇 上述 突发 情况 申请 能力 无人 代理 申请 社会 救助 机构 了解 情况 后应 协助 申请 临时 救助 死亡 无遗 属 遗产 居民 应由 殡仪馆 办理 丧葬 事宜 财政 支付 丧葬费

四、相似度 simhash 算法与相似度的计算

相似度 simhash 算法

simhash 是由 Charikar 在 2002 年提出来的，论文名为《Similarity estimation techniques from rounding algorithms》。Google 基于 simhash 在海量网页中进行相似度计算并去重。通常对比两个文档是否相同时，会计算对应的 hash 值，常见的算法包括 md5 和 sha256。实际使用中，对于检测文档是否被篡改时，使用 hash 值具有不错的表现。但是当文档内容因为修改少许文字，插入广告甚至只是修改了标点符合和错别字，都会导致 hash 值改变，可是文档的核心内容并未发生改变。如何使用数学的方法表征这种文档相似性呢？simhash 的设计初衷就是使用一种所谓局部 hash 的方法，可以既可以敏感的识别文档的少许修改又可以识别出文档的大多数内容相同。

simhash 的一种典型实现就是将一个文档最后转换成一个 64 位的字节的特征字或者说 simhash 值，然后判断重复只需要判断他们的特征字的距离是不是小于 3，就可以判断两个文档是否相似。这个距离使用海明距离，即两个 simhash 值取异或后二进制中 1 的个数。大家可以结合自身业务特点修改 simhash 值的位数以及判断文档相似性的海明距离的值。



如图所示，计算 6 位 simhash 值典型的实现算法为：

将 Doc 分词和计算权重，抽取出 n 个(关键词，权重)对，即图中的(feature, weight)。

计算关键词的 hash，生成图中的(hash,weight)，并将 hash 和 weight 相乘，这一过程是对 hash 值加权。

将 hash 和 weight 相乘的值相加，比如图中的[13, 108, -22, -5, -32, 55]，并最终转换成 simhash 值 110001，转换的规则为正数为 1 负数为 0。

相似度的计算

针对 1,093,388 条新闻语料，大范围的计算相似度需要耗费大量的时间，这显然是不可取的。我们采取的方法是，利用 random 库从 1,093,388 条新闻语料随机抽取 200 条新闻语料，这 200 条语料分别与 1,093,388 条新闻语料做相似度计算，记录下两两之间的新闻索引和计算得出的 score。

```
simhashListSorted: [(18, 0), (8, 21), (9, 27), (17, 27), (6, -
> special variables
> function variables
> 00: (18, 0)
> 01: (8, 21)
> 02: (5, 27)
> 03: (17, 27)
> 04: (6, 28)
> 05: (11, 28)
> 06: (4, 29)
> 07: (9, 29)
> 08: (18, 29)
> 09: (14, 29)
> 10: (3, 30)
> 11: (16, 30)
> 12: (7, 31)
> 13: (8, 32)
> 14: (13, 32)
> 15: (19, 32)
> 16: (2, 34)
> 17: (12, 34)
> 18: (1, 36)
> 19: (15, 38)
```

score 越低则代表着两条新闻越相似。我们将相似度最高的前 20000 条新闻文本以 txt 的形式保存下来，把他们的索引和相似度有 csv 格式存储。

五、计算结果

相似度计算的结果由 3 部分组成，相似度 score 和新闻索引组成的 csv、新闻文本 txt、新闻和新闻的词云。

相似度和索引 csv

	index	sim
0	37025	0
1	33773	7
2	85977	7
3	560586	8
4	1012555	8
5	611862	10
6	262697	11
7	680565	11
8	143239	12
9	745701	12
10	1175	13
11	155568	13
12	281028	13
13	372854	13
14	419894	13
15	804060	13
16	983921	13
17	6183	14

新闻文本 txt

新闻 37025 新闻 33773 新闻 85977 新闻 560586 新闻 1012555 新闻 611862 新闻 262697 新闻 680565 新闻 143239 新闻 745701 新闻 1175 新闻 155568 新闻 281028 新闻 372854 新闻 419894 新闻 804060 新闻 983921 新闻 6183

新闻 37025

新闻 33773

新闻 85977

新闻 560586

词云

