

金融大数据 实验一 实验报告

151220006 陈安宇

一、实验需求

需求 1: 针对股票新闻数据集中的新闻标题, 编写 WordCount 程序, 统计所有除 Stop-word (如 “的”, “得”, “在” 等) 出现次数 k 次以上的单词计数, 最后的结果按照词频从高到低排序输出。

需求 2: 针对股票新闻数据集, 以新闻标题中的词组为 key, 编写带 URL 属性的文档倒排索引程序, 将结果输出到指定文件。

二、实验设计说明

实验 1: 对中文新闻标题进行分词并计数, 实现出现频率大于某一值的所有词的从高到低的输出。

设计思路: 程序框架: 该实验的主要部分是对词频计数, 因此采用 WordCount2 算法的框架。

中文进行分词程序: 编写一个简单的分词程序。首先把预先准备好的词典中的词读入一个字符串类型的链表中。并记录最长的字符串的长度。维护一个最终分词结果的链表。对读入的段落按照取指定的最大长度的文本去词典里面匹配, 如果匹配不到, 则长度减一继续匹配, 若匹配成功, 把该词加入分词结果列表, 并从待分词文本中去除已经分词的文本。直到待分词文本为空, 分词结束。

停词: 读取停词的文件, 将停词放入一个链表中。从取出的中文标题中逐一地去除停词。

词频从高到低输出: 调用本身有的比较器 Comparator, 在前面加上负号实现从词频从高到低的排序。

程序流程: 从文本中获取中文标题 → 去除停词 → 分词 → 去除长度为 1 的词 → 计数 → 与 K 比较 → 结果写入临时文件 → 排序 → 写入结果文件

设计说明: 分词类 WordSeg: List<String> DIC 存储词典
 List<String> result 存储分词结果
词频计数 WordCount3: Set<String> patternsToSkip 存储停词
 String[] itr = line.split("\t");用于截取中文文本
 List<String> lis 存储切分好的中文单词

程序运行: bin/hadoop jar ./share/hadoop/mapreduce/WordCount3.jar WordCount3
 -Dwordcount.case.sensitive=true inputdata output1 100
 -skip /user/hadoop/stopwords/patterns.txt

运行结果截图:

28275	公告	120	绿色
8449	提示	119	处理
7183	板块	119	变动
7107	股东	118	账户
6609	业绩	118	陕西
6292	涨停	117	民营
5036	晚间	116	宣布
4178	投资	115	黑马
3533	资金	114	募投
3499	控股	114	支持
3038	行业	113	期货
3022	减持	113	高速公路
2890	今日	112	幸福
2886	增持	111	投入
2877	年报	110	创意
2806	科技股份	110	特钢
2726	中国	110	违法
2668	年度	109	阅兵
2625	交易	109	暂时
2622	增长	109	贷款
2550	净利	109	精华
2317	集团股份	108	合肥
2299	资产	108	说明书
2260	股权	108	人事
2133	中报	108	格局
2093	股价	107	长期
2072	利好	107	电源
2031	股票	106	新低
1831	报告	106	交通
1779	拉升	106	存货
1729	次新股	106	行政处罚
1705	市场	105	机场
1701	突破	104	拓展
1647	披露	103	移动
1619	银行	102	营业部
1549	龙头	102	亮点
1508	发展	101	债务
1498	项目	101	具备

词频较高部分：

词频较低部分：

- 实验不足之处：** 1.用 StringTokenizer 比 String.split()效率高。整个程序在处理大数据集 fulldat 时效率不高，运行较慢。
- 2.在实现通过词频排序时，没有考虑当两个词词频相同时的排序情况。

实验 2：实现带 URL 属性的文档倒排索引程序。

设计思路： 该实验的主要目的是文档的倒排索引，借鉴带词频属性的文档倒排索引程序的框架。

中文分词程序与停词：同实验 1。

Map 类： 计算每一行的词与该行每个词出现的次数。每个词作为 KEY，该行的 URL+@+该词在该行的词频作为 VALUE。

Combiner 类： 通过@符号抽取出 URL 和词频，合并 URL 和词频，合并本地 mapper

函数的输出

Reducer:合并 URL 和词频，合并 mapper 函数的输出,将结果写入输出文件。

程序流程: 从文本中获取中文标题→去除停词→分词→去除长度为 1 的词→计数并记录该词的 URL→合并 URL 和词频→写入结果文件

设计说明: 分词类和停词类同实验 1

在 map 中维护一个 `hashmap()`，统计一行出现的词，以及这一行出现该词的次数。

`Text fileName_frequency` 记录 URL 与词频

`String UriList` 记录合并的 URL。URL 与 URL 之间用 `*****` 隔开

程序运行: `bin/hadoop jar ./share/hadoop/mapreduce/InvertedIndexer.jar InvertedIndexer -Dwordcount.case.sensitive=true inputdata output2 -skip /user/hadoop/stopwords/patterns.txt`

运行结果截图:

```
余颖 *****http://finance.sina.com.cn/roll/2016-12-15/doc-ifywtkcf7712268.shtml*****http://finance.sina.com.cn/stock/t/2016-12-15/doc-ifywtqgn8611247.shtml**
使命 *****http://finance.sina.com.cn/stock/s/2017-08-01/doc-ifyinryq7378957.shtml*****http://finance.sina.com.cn/stock/s/2017-08-01/doc-ifyinryq7378957.shtm
供应国 *****http://finance.sina.com.cn/stock/t/2017-03-30/doc-ifycwunr8176980.shtml*****http://finance.sina.com.cn/stock/t/2017-03-30/doc-ifycwunr8176980.shtm
供电 *****http://finance.sina.com.cn/roll/2017-08-04/doc-ifyitayr8956265.shtml*****http://finance.sina.com.cn/stock/t/2016-12-28/doc-ifyxvzr7798331.shtml**
供货 *****http://finance.sina.com.cn/stock/t/2017-04-19/doc-ifyeimzx7024842.shtml*****http://finance.sina.com.cn/stock/s/2016-11-08/doc-ifyxnmety7695783.shtml
促进 *****http://finance.sina.com.cn/stock/gujlayidong/2016-02-25/doc-ifykpyxsk1641347.shtml*****http://finance.sina.com.cn/stock/companyresearch/2017-06-05/c
俄罗斯 *****http://finance.sina.com.cn/stock/t/2016-10-27/doc-ifyxwuff6963591.shtml*****http://cj.sina.com.cn/article/detail/12097297631/323086*****http://fini
报价 *****http://finance.sina.com.cn/stock/hkstock/hkstocknews/2017-07-10/doc-ifyhvyie0851753.shtml*****http://finance.sina.com.cn/stock/s/2017-06-23/doc-ify
*****http://finance.sina.com.cn/stock/s/2017-07-27/doc-ifyinwmp0166312.shtml*****http://finance.sina.com.cn/stock/s/2017-08-11/doc-ifyixhyw7203106.shtml
保健品 *****http://finance.sina.com.cn/roll/2017-07-08/doc-ifyhwehx5353240.shtml*****http://cj.sina.com.cn/article/detail/5894240270/301235*****http://finan
保利 *****http://cj.sina.com.cn/article/detail/1644983660/146633*****http://finance.sina.com.cn/stock/s/2017-04-18/doc-ifyeifqx6210095.shtml*****http://fini
保持 *****http://finance.sina.com.cn/stock/hyyj/2017-04-17/doc-ifyeimzx6590462.shtml*****http://finance.sina.com.cn/stock/hyyj/2017-04-17/doc-ifyeimzx6590462.
保持 *****http://finance.sina.com.cn/stock/t/2017-01-14/doc-ifyxgnva3528520.shtml*****http://finance.sina.com.cn/stock/t/2017-01-13/doc-ifyxgnva3403809.shtml
保障 *****http://finance.sina.com.cn/stock/s/2017-04-25/doc-ifyepnea4981324.shtml*****http://finance.sina.com.cn/chanying/geneve/2017-07-05/doc-ifyhvyie02266.
保驾护航 *****http://finance.sina.com.cn/roll/2017-06-19/doc-ifyhthrt4715645.shtml*****http://finance.sina.com.cn/roll/2017-06-19/doc-ifyhthrt4715645.shtml*
信心 *****http://finance.sina.com.cn/stock/t/2017-04-25/doc-ifyepnea4981324.shtml*****http://finance.sina.com.cn/stock/t/2017-07-19/doc-ifyiamif3602878.shtml
信息 *****http://finance.sina.com.cn/stock/s/2017-08-03/doc-ifyiswpt5085441.shtml*****http://finance.sina.com.cn/stock/t/2017-04-21/doc-ifyepnea4411181.shtml
信息化 *****http://finance.sina.com.cn/stock/hyyj/2016-12-27/doc-ifyxusa5567410.shtml*****http://finance.sina.com.cn/stock/hyyj/2017-03-20/doc-ifycnpiu916515.
信息系统 *****http://finance.sina.com.cn/stock/t/2017-04-01/doc-ifycwjxr9026310.shtml*****http://finance.sina.com.cn/roll/2016-12-23/doc-ifyxqk6386165.shtu
信托 *****http://finance.sina.com.cn/stock/t/2017-01-18/doc-ifyxrunxf1362139.shtml*****http://finance.sina.com.cn/roll/2017-08-08/doc-ifyiswpt5949012.shtml**
信托 *****http://cj.sina.com.cn/article/detail/1131398582/317369*****http://finance.sina.com.cn/stock/t/2017-04-27/doc-ifyetec6778900.shtml*****http://fini
倒挂 *****http://finance.sina.com.cn/stock/t/2017-02-20/doc-ifyarxrc50212570.shtml*****http://finance.sina.com.cn/stock/s/2017-05-25/doc-ifyfqgh8233540.shtml
*****http://finance.sina.com.cn/stock/t/2017-05-23/doc-ifyfkme0121209.shtml*****http://finance.sina.com.cn/stock/t/2017-04-14/doc-ifyeifqx6210096.shtml
*****http://finance.sina.com.cn/stock/t/20151214/020424007476.shtml*****http://finance.sina.com.cn/stock/s/2017-08-11/doc-ifyixipt1004590.shtml*****ht
*****http://finance.sina.com.cn/stock/t/2017-03-22/doc-ifycnphv521637.shtml*****http://finance.sina.com.cn/stock/t/2017-04-08/doc-ifyeayz7179614.
倾斜 *****http://finance.sina.com.cn/roll/2017-03-31/doc-ifycwunr8369682.shtml*****http://finance.sina.com.cn/stock/hyyj/2016-08-23/doc-ifyxvzr7608735.shtml
原期 *****http://cj.sina.com.cn/article/detail/5995014987/207219*****http://finance.sina.com.cn/stock/hyyj/2017-06-02/doc-ifyfuzmy1125900.shtml*****http://f
*****http://finance.sina.com.cn/stock/s/2017-02-22/doc-ifyarxrc5439940.shtml*****http://finance.sina.com.cn/stock/t/2017-02-22/doc-ifyarxrc5439940.shtml
*****http://finance.sina.com.cn/stock/s/2017-05-03/doc-ifyetwml829276.shtml*****http://tech.sina.com.cn/lt/2017-05-03/doc-ifyetwml829276.shtml82
*****http://finance.sina.com.cn/roll/2017-08-05/doc-ifyitawv5347458.shtml*****http://finance.sina.com.cn/stock/s/2017-02-24/doc-ifyavvcv679354.shtml**
*****http://finance.sina.com.cn/stock/s/2017-08-07/doc-ifyitayr944825.shtml*****http://finance.sina.com.cn/roll/2017-06-19/doc-ifyhfpac5266870.shtml**
*****http://finance.sina.com.cn/roll/2017-02-28/doc-ifyavrxs539060.shtml*****http://finance.sina.com.cn/roll/2016-08-29/doc-ifyvixeq0632605.shtml*****
*****http://finance.sina.com.cn/stock/s/2017-04-18/doc-ifyeimqx4760452.shtml*****http://finance.sina.com.cn/roll/2017-04-01/doc-ifycwyns4071524.shtml**
*****http://finance.sina.com.cn/stock/t/2017-01-10/doc-ifyxkfun6479226.shtml*****http://finance.sina.com.cn/stock/s/2017-05-31/doc-ifyfqgh9058043.shtml
*****http://finance.sina.com.cn/stock/s/2017-08-07/doc-ifyitapp2307855.shtml*****http://finance.sina.com.cn/stock/s/2017-06-28/doc-ifyhmpew3488525.shtml**
*****http://cj.sina.com.cn/article/detail/231077472/240292*****http://cj.sina.com.cn/article/detail/231077472/240292*****http://cj.sina.com.cn/artic
*****http://finance.sina.com.cn/stock/s/2017-07-03/doc-ifyhrxsk1570990.shtml*****http://finance.sina.com.cn/stock/s/2017-07-03/doc-ifyhrxsk1570990.shtm
```

实验不足之处: 1.用 `StringTokenizer` 比 `String.split()`效率高。整个程序在处理大数据集 `fulldat` 时效率不高，运行较慢。
2.对于关键字的筛选做的不够细致。