

MapReduce 基础编程实验报告

151180045 侯汶昕

1. 实验目的

利用 MapReduce 编程实现对给定格式的文本文件进行基本的分词，词频统计，倒排索引和二次排序的操作

2. 实验重点与难点

1. 文本分词和去除停用词
2. 倒排索引
3. 降序排序及二次排序

3. 实验环境

实验中使用了两台虚拟机，虚拟机具有相同的环境和配置：

系统：ubuntu 16.04 LTS 64-bit

内存：1.9GB

4. 实验要求

4.1 针对股票新闻数据集中的新闻标题，去除其中的停用词（Stop-Word）后，对出现次数为 k 以上的单词进行词频统计，结果按词频从高到底排序输出。

4.2 针对股票新闻数据集，以新闻标题中的词组为 **key**，编写带 URL 属性和词频的文档倒排索引程序，并按照词频从大到小排序，将结果输出到指定文件。

规定输出格式为：

高速公路, 10, 股票代码, [url0, url1,...,url9]

5. 实验设计思路

该实验要求实现两个需求，考虑分别使用两个 MapReduce 程序进行实现。

5.1 中文词频统计

第一个需求的基本框架依然是 WordCount，有三个重点：

第一个重点也是一个难点，那就是对中文进行词频统计。词频统计的前提是分词，相对可以简单的依托空格来分词的英文来说，中文的分词明显要难得多。在尝试自己写程序进行分词失败后，我还是决定使用开源的 HanLP 中文自然语言处理包进行分词（感谢 HanLP）。

第二个重点在于统计除停用词以外的出现 k 次以上的单词。去除停用词的功能依然借助了 HanLP 进行实现（再次感谢 HanLP）。只统计 k 次以上的单词是比较简单的，只需要对统计完

的停用词出现次数进行一个条件判断，将出现次数大于 k 次的键值对写入 context 就行。

第三个重点在于要求将最后结果按照词频从高到低进行排序输出。考虑将词频统计完的结果输出到一个临时文件后，再新建一个 Job 对其进行倒序排序。

结合上述分析，最终实现的程序的大体框架是在 Map 阶段对于输入利用 HanLP 进行分词、去除停用词等基本操作，并以 $\langle \text{word}, 1 \rangle$ 键值对 (Key-Value) 的方式传输给 Combiner 和 Reducer。在 Combiner 和 Reducer 中以单词进行划分，对 Value 进行求和得到总出现次数 Sum，并将和大于 k 的键值对 $\langle \text{word}, \text{Sum} \rangle$ 输出得到一个临时文件。再利用系统给出的 InverseMapper 类将临时文件中的 $\langle \text{word}, \text{Sum} \rangle$ 进行键值对交换以便对 Sum 进行排序，再重写一个降序排序的 IntWritableDecreasingComparator 类并利用 setSortComparatorClass 进行倒序排序，最后输出到文件即可实现第一个需求。

5.2 新闻倒排索引

第二个需求相比第一个需求要难一些。由于输出格式受限，我首先就是考虑 Key-Value 对如何设计。

在 Map 阶段，Value 输入格式为：

股票代码 Code\t 股票名称\t 日期\t 时间\t 新闻标题 Title\t 链接 URL

在 Map 阶段对 Title 进行分词、去除停用词等基本操作后得到了 TitleWords，考虑到我下一步要进行词频统计，并且是按照不同的 Code 和不同的 TitleWord 进行统计，所以我决定将 Code 和 TitleWord 进行组合成 Key，URL 和出现次数 1 进行组合作为 Value，剩下的没用信息直接丢弃。于是 Map 的输出键值对为：

$\langle \text{TitleWord\#Code}, \text{URL\#1} \rangle$

这样可以保证同一个 Code 下的同一个 TitleWord 会被分到一起进行 Reduce。

那么在 Reducer 中，我就可以对每一个股票下的每一个单词进行词频统计得到 Sum，在实施过程中我还发现其实这里可以顺便把 URL 也拼起来，于是 Reducer 的输出键值对为：

$\langle \text{TitleWord\#Code\#[URL0, URL1, ...]}, \text{Sum} \rangle$

这样就得到了最后输出所需要的所有信息了，剩下要做的就是倒序排序和调整输出格式。于是我决定使用和上面第一个需求类似的方法，将刚才 Reducer 的输出放在一个临时文件中，再新建一个 Job (SecondSort) 对其进行格式重组和排序。

在新建 Job 的 Map 阶段，输入键值对格式为 $\langle \text{Text}, \text{IntWritable} \rangle$ ，内容和上面的 Reducer 输出一致，考虑反正要格式调整，不如在这里就把格式调一下，于是没有别的操作，仅仅更换了输出时的键值对格式为 $\langle \text{Text}, \text{NullWritable} \rangle$ ，键值对内容为：

$\langle \text{TitleWord\#Sum\#Code\#[URL0, URL1, ...]}, \text{NullWritable} \rangle$

因为最后的结果要将 TitleWord 一致的放在一起，所以还需要重写一个 SecondSortPartitioner，依据 TitleWord 而不是全部的 Key 进行划分。然后和第一个需求类似，重写一个基于二次排序并且是降序排序的 SecondSortComparator 类并利用 setSortComparatorClass 进行倒序排序，首先依据字符串形式的 TitleWord 进行排序，在 TitleWord 一致的情况下依据整型的 Sum 进行倒序排序。

在 Reduce 阶段，将 Mapper 输出键值对中的所有“#”替换成“，”，再输出到文件即可完成第二个需求。

6. 重要类的设计说明

6.1 中文词频统计

6.1.1 MainReducer 类

MainReducer 类是第一个 Job 的关键任务，主要实现的功能为词频统计和 URL 的拼接，由于便于后续处理的分割，采用“#”作为连接符。

```
public static class MainReducer extends
    Reducer<Text, Text, Text, IntWritable> {
    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<Text> values, Context context)
        throws IOException, InterruptedException {
        int sum = 0;
        Text word = new Text();
        String titleWord = key.toString().split("#")[0];
        String code = key.toString().split("#")[1];
        String URL;
        StringBuilder strBuilder = new StringBuilder("");
        int Occurrence;
        for (Text val : values) {
            URL = val.toString().split("#")[0];
            strBuilder.append(URL+", ");
            Occurrence = Integer.parseInt(val.toString().split("#")[1]);
            sum += Occurrence;
        }
        strBuilder.replace(strBuilder.length()-2, strBuilder.length(), "");
        word.set(titleWord + "#" + code + "#" + strBuilder.toString());
        result.set(sum);
        context.write(word, result);
    }
}
```

6.1.1 InverseMapper 类

由于 MapReduce 默认按照 Key 进行排序，为了使最后结果按词频进行排序，需要调用 InverseMapper 类。该类是系统自带的类，用于调换 Key 和 Value，输出<Value, Key>

该类源码如下：

```
public class InverseMapper<K, V> extends
```

```

        MapReduceBase implements Mapper<K, V, V, K> {
        /** The inverse function. Input keys and values are swapped.*/
        public void map(K key, V value,
                        OutputCollector<V, K> output, Reporter reporter)
            throws IOException {
            output.collect(value, key);
        }
    }
}

```

6.1.2 IntWritableDecreasingComparator 类

由于 MapReduce 默认对 Key 进行升序排序，因此为了实现降序排序，需要重写一个 IntWritableDecreasingComparator 类。该类通过继承 IntWritable.Comparator 类，将父类 compare 函数返回的结果分别取反，从而实现降序排序。

```

private static class IntWritableDecreasingComparator extends IntWritable.Comparator{
    public int compare(WritableComparable a, WritableComparable b){
        return -super.compare(a, b);
    }
    public int compare(byte[] b1, int s1, int l1, byte[] b2, int s2, int l2){
        return -super.compare(b1, s1, l1, b2, s2, l2);
    }
}

```

6.2 新闻倒排索引

6.2.1 SecondSortPartitioner 类

由于 MapReduce 默认是按照相同的完整的 Key 来分区，而在 Mapper 阶段之后输出的格式为 <TitleWord#Sum#Code#[URL0, URL1, ...], NullWritable>。为了保证最后同一个 TitleWord 能放在一起并且正常排序，这里需要重写一个 Partitioner，实现仅针对 TitleWord 进行分区。

```

public static class SecondSortPartitioner extends HashPartitioner<Text, NullWritable>
{
    //Partitioned by titleWord
    public int getPartition(Text key, NullWritable value, int numReduceTasks){
        String term = new String();
        term = key.toString().split("#")[0];
        return super.getPartition(new Text(term), value, numReduceTasks);
    }
}

```

6.2.2 SecondSortComparator 类

由于要实现二次排序和依据词频降序排序两个功能，所以这里必须重写一个 `Comparator` 类。该类首先比较 `TitleWord`，如果不一致则比较 `TitleWord`（`TitleWord` 并不需要反向排序）；如果一致则比较单词的词频，并且返回与正常的 `Compare` 函数相反的结果从而实现降序排序。

```
private static class SecondSortComparator extends WritableComparator{
    protected SecondSortComparator(){
        super(Text.class, true);
    }
    @Override
    public int compare(WritableComparable a, WritableComparable b) {
        Text o1 = (Text) a;
        Text o2 = (Text) b;
        int occurrence1 = Integer.parseInt(o1.toString().split("#")[1]);
        String titleWord1 = o1.toString().split("#")[0];
        int occurrence2 = Integer.parseInt(o2.toString().split("#")[1]);
        String titleWord2 = o2.toString().split("#")[0];
        if(! titleWord1.equals(titleWord2))
            return titleWord1.compareTo(titleWord2);
        else
            return occurrence2 - occurrence1;
    }
}
```

7. 实验结果及截图

7.1 中文分词结果

```
houwenxin@ubuntu:~/hadoop_installs/hadoop-2.9.1$ bin/hadoop jar /home/houwenxin/Desktop/HMX.WCN.jar wcn input output 30
18/11/13 23:27:30 INFO client.RMProxy: Connecting to ResourceManager at host-0/127.0.0.1:8032
18/11/13 23:27:32 INFO input.FileInputFormat: Total input files to process : 1
18/11/13 23:27:33 WARN hdfs.DataStreamer: Caught exception
java.lang.InterruptedException
    at java.lang.Object.wait(Native Method)
    at java.lang.Thread.join(Thread.java:1252)
    at java.lang.Thread.join(Thread.java:1326)
    at org.apache.hadoop.hdfs.DataStreamer.closeResponder(DataStreamer.java:980)
    at org.apache.hadoop.hdfs.DataStreamer.endBlock(DataStreamer.java:630)
    at org.apache.hadoop.hdfs.DataStreamer.run(DataStreamer.java:807)
18/11/13 23:27:33 INFO mapreduce.JobSubmitter: number of splits:1
18/11/13 23:27:33 INFO Configuration.deprecation: yarn.resourcemanager.system-metrics-publisher.enabled is deprecated. Instead, use yarn.system-metrics-publisher.enabled
18/11/13 23:27:33 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1542180344850_0001
18/11/13 23:27:34 INFO impl.YarnClientImpl: Submitted application application_1542180344850_0001
18/11/13 23:27:34 INFO mapreduce.Job: The url to track the job: http://ubuntu:8088/proxy/application_1542180344850_0001/
18/11/13 23:27:34 INFO mapreduce.Job: Running job: job_1542180344850_0001
18/11/13 23:27:49 INFO mapreduce.Job: Job job_1542180344850_0001 running in uber mode : false
18/11/13 23:27:49 INFO mapreduce.Job: map 0% reduce 0%
18/11/13 23:28:10 INFO mapreduce.Job: map 49% reduce 0%
18/11/13 23:28:15 INFO mapreduce.Job: map 100% reduce 0%
18/11/13 23:28:38 INFO mapreduce.Job: map 100% reduce 100%
18/11/13 23:28:39 INFO mapreduce.Job: Job job_1542180344850_0001 completed successfully
18/11/13 23:28:39 INFO mapreduce.Job: Counters: 49
  File System Counters
    FILE: Number of bytes read=44272
    FILE: Number of bytes written=484311
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0
    HDFS: Number of bytes read=20980068
    HDFS: Number of bytes written=66003
    HDFS: Number of read operations=6
    HDFS: Number of large read operations=0
    HDFS: Number of write operations=2
  Job Counters
    Launched map tasks=1
    Launched reduce tasks=1
    Data-local map tasks=1
    Total time spent by all maps in occupied slots (ms)=23666
    Total time spent by all reduces in occupied slots (ms)=18770
    Total time spent by all map tasks (ms)=23666
    Total time spent by all reduce tasks (ms)=18770
    Total vcore-milliseconds taken by all map tasks=23666
    Total vcore-milliseconds taken by all reduce tasks=18770
```

```
houwenxin@ubuntu:~/hadoop_installs/hadoop-2.9.1$ bin/hdfs dfs -cat output/*
28274 公告
18416 股份有限公司
15533 日
12665 公司
9107 上市公司
8686 %
8449 提示
8066 股份
7570 沪
7410 深
7349 股市
7318 业绩
6963 导航
6907 股海
6777 板块
6189 股东
5943 涨停
5796 科技
5036 晚间
4659 机构
4639 速
4638 集团
4505 成
4443 A
4437 快讯
4431 线
4120 超
3956 概念股
3893 年度
3539 中
3503 拟
3446 资
3429 新
3211 投
3168 控
2940 报
2914 行
2856 涨
2840 看
2814 近
```

```
31 心
31 飞
31 英
31 茂
31 固
31 喜
31 爱
31 道
31 具
31 索
31 出
31 自
31 红
31 售
31 蜂
31 虫
31 硕
31 闻
31 韵
31 横
31 店
31 盘
31 美
31 森
31 启
31 氧
31 果
31 汁
31 服
31 装
31 公
31 司
31 业
31 中
31 名
31 流
31 感
31 凯
31 大
31 升
31 嘉
31 世
31 趣
31 星
31 期
31 六
31 漏
31 洞
31 缓
31 解
31 旭
31 黄
31 浦
31 灵
31 Q
31 风
31 景
31 远
31 望
31 空
31 落
31 新
31 药
31 金
31 法
31 司
31 整
31 特
31 支
31 出
31 临
31 床
31 探
31 索
houwenxin@ubuntu:~/hadoop_installs/hadoop-2.9.1$
```

程序基本实现了词频统计和倒序排序功能，不过由于一些分词的原因，有一些符号也被统计了进来。测试时 k 取值为 30。

7.2 新闻倒排索引结果

```

housenxin@ubuntu:~/hadoop_installs/hadoop-2.9.15$ bin/hadoop jar /home/housenxin/Desktop/HMW_InvertedIndex.jar invertedIndex input output
18/11/15 00:32:45 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
18/11/15 00:32:45 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
18/11/15 00:32:46 INFO mapreduce.JobInputFormat: Total input files to process : 1
18/11/15 00:32:46 INFO mapreduce.JobSubmitter: number of splits=1
18/11/15 00:32:47 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local347624016_0001
18/11/15 00:32:47 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
18/11/15 00:32:47 INFO mapreduce.Job: Running job: job_local347624016_0001
18/11/15 00:32:47 INFO mapred.LocalJobRunner: OutputCommitter set in config null
18/11/15 00:32:47 INFO mapred.LocalJobRunner: FileOutputCommitter: FileOutputCommitter Algorithm version is 1
18/11/15 00:32:47 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
18/11/15 00:32:47 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
18/11/15 00:32:47 INFO mapred.LocalJobRunner: Waiting for map tasks
18/11/15 00:32:47 INFO mapred.LocalJobRunner: Start task: attempt_local347624016_0001_m_000000_0
18/11/15 00:32:47 INFO output.FileOutputCommitter: FileOutputCommitter Algorithm version is 1
18/11/15 00:32:47 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
18/11/15 00:32:48 INFO mapred.MapTask: Using ResourceCalculatorProcessTree : [ ]
18/11/15 00:32:48 INFO mapred.MapTask: Processing split: hdfs://host-0:9000/user/housenxin/input/fulldata.txt:0:420979951
18/11/15 00:32:48 INFO mapred.MapTask: (EQUATOR) 0 kvl 26214396(104857584)
18/11/15 00:32:48 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
18/11/15 00:32:48 INFO mapred.MapTask: sort.limit at 8386000
18/11/15 00:32:48 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
18/11/15 00:32:48 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
18/11/15 00:32:48 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
18/11/15 00:32:48 INFO mapreduce.Job: Job job_local347624016_0001 running in uber mode : false
18/11/15 00:32:48 INFO mapreduce.Job: map 0% reduce 0%
18/11/15 00:33:00 INFO mapred.LocalJobRunner: map > map
18/11/15 00:33:00 INFO mapred.MapTask: Spilling map output
18/11/15 00:33:00 INFO mapred.MapTask: bufstart = 0; bufend = 71002656; bufvoid = 104857600
18/11/15 00:33:00 INFO mapred.MapTask: kvstart = 26214396(104857584); kvoid = 22993548(91974192); length = 3220849/6553600
18/11/15 00:33:00 INFO mapred.MapTask: (EQUATOR) 74229040 kvl 18557256(74229024)
18/11/15 00:33:00 INFO mapreduce.Job: map 50% reduce 0%
18/11/15 00:33:05 INFO mapred.MapTask: Finished spill 0
18/11/15 00:33:05 INFO mapred.MapTask: (RESET) equator 74229040 kv 18557256(74229024) kvl 17751708(71066800)
18/11/15 00:33:05 INFO mapred.LocalJobRunner: map > map
18/11/15 00:33:05 INFO mapred.MapTask: Starting flush of map output
18/11/15 00:33:05 INFO mapred.MapTask: Spilling map output
18/11/15 00:33:05 INFO mapred.MapTask: bufstart = 74229040; bufend = 93351536; bufvoid = 104857600
18/11/15 00:33:05 INFO mapred.MapTask: kvstart = 18557256(74229024); kvoid = 17690008(70766032); length = 867249/6553600
18/11/15 00:33:06 INFO mapred.LocalJobRunner: map > sort
18/11/15 00:33:06 INFO mapred.MapTask: Finished spill 1
18/11/15 00:33:06 INFO mapred.Merger: Merging 2 sorted segments
18/11/15 00:33:06 INFO mapred.Merger: Down to the last merge-pass, with 2 segments left of total size: 92169187 bytes
18/11/15 00:33:06 INFO mapreduce.Job: map 67% reduce 0%

```

[illegible]

程序基本实现了新闻的倒排索引，并且能够把同一个词放在一起，按照降序排序进行输出。图示为测试时尝试搜索了“高速公路”这个词，得到的结果。

8. 实验中可以改进之处

本实验中，我完成了实验所要求的两个任务：面向中文的词频统计和针对金融新闻的倒排索引，并且基本完成了所有需求。但是该实验还有几个可以改进优化的地方：

1. 这两个需求在的实现和功能都是十分类似的，因而可以考虑把他们做成一个程序，并且还可以通过定制参数对功能进行选择（比如只进行词频统计），简化用户操作。
2. 该实验由于使用了开源的分词包 HanLP，并且使用的词典没有对金融专业的词语，符号进行特别设计，所以分词过程中还是有一些专业的词语或者符号被拆开来了，后续可以考虑通过应用自己的词典，或者将这些词语添加进默认词典。

3. 我在实验中尝试过使用自己的 Stop-Word 表和分词词典，但是由于 MapReduce 程序直接从 hdfs 上读取文件的效率太低，只能放弃，后续或许可以通过使用 DistributedCache 的方法加快文件读取的运行效率，从而实现使用自己的词表。