

Top N

Driver code:

```
import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class TopN {
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = (new GenericOptionsParser(conf, args)).getRemainingArgs();
        if (otherArgs.length != 2) {
            System.err.println("Usage: TopN <in> <out>");
            System.exit(2);
        }
        Job job = Job.getInstance(conf);
        job.setJobName("Top N");
        job.setJarByClass(TopN.class);
        job.setMapperClass(TopNMapper.class);
        job.setReducerClass(TopNReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(otherArgs[0]));
        FileOutputFormat.setOutputPath(job, new Path(otherArgs[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }

    public static class TopNMapper extends Mapper<Object, Text, Text, IntWritable> {
        private static final IntWritable one = new IntWritable(1);

        private Text word = new Text();

        private String tokens = "[_!$#<>\\^=\\[\\]\\|\\*\\/\\\\,;\\.\\-:()?!\\\"'"]";

        public void map(Object key, Text value, Mapper<Object, Text, Text, IntWritable>.Context
context) throws IOException, InterruptedException {
            String cleanLine = value.toString().toLowerCase().replaceAll(this.tokens, " ");
            StringTokenizer itr = new StringTokenizer(cleanLine);
            while (itr.hasMoreTokens()) {
```

```

        this.word.set(itr.nextToken().trim());
        context.write(this.word, one);
    }
}
}
}

```

TopN Combiner

```

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

public class TopNCombiner extends Reducer<Text, IntWritable, Text, IntWritable> {
    public void reduce(Text key, Iterable<IntWritable> values, Reducer<Text, IntWritable, Text,
IntWritable>.Context context) throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values)
            sum += val.get();
        context.write(key, new IntWritable(sum));
    }
}

```

TopN Mapper

```

import java.io.IOException;
import java.util.StringTokenizer;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class TopNMapper extends Mapper<Object, Text, Text, IntWritable> {
    private static final IntWritable one = new IntWritable(1);

    private Text word = new Text();

    private String tokens = "[_!$#<>\\^=\\[\\]\\*\\/\\\\\\,;\\.\\-:()?!\"']";

    public void map(Object key, Text value, Mapper<Object, Text, Text, IntWritable>.Context context)
throws IOException, InterruptedException {
        String cleanLine = value.toString().toLowerCase().replaceAll(this.tokens, " ");
        StringTokenizer itr = new StringTokenizer(cleanLine);
        while (itr.hasMoreTokens()) {
            this.word.set(itr.nextToken().trim());
            context.write(this.word, one);
        }
    }
}

```

```

    }
}
}

```

TopN Reducer

```

import java.io.IOException;
import java.util.HashMap;
import java.util.Map;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
import utils.MiscUtils;

public class TopNReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
    private Map<Text, IntWritable> countMap = new HashMap<>();

    public void reduce(Text key, Iterable<IntWritable> values, Reducer<Text, IntWritable, Text,
IntWritable>.Context context) throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values)
            sum += val.get();
        this.countMap.put(new Text(key), new IntWritable(sum));
    }

    protected void cleanup(Reducer<Text, IntWritable, Text, IntWritable>.Context context) throws
IOException, InterruptedException {
        Map<Text, IntWritable> sortedMap = MiscUtils.sortByValues(this.countMap);
        int counter = 0;
        for (Text key : sortedMap.keySet()) {
            if (counter++ == 20)
                break;
            context.write(key, sortedMap.get(key));
        }
    }
}

```

Utils/MiscUtils

```

package utils;
import java.util.*;
public class MiscUtils {
    /**
     * sorts the map by values. Taken from:
     * http://javarevisited.blogspot.it/2012/12/how-to-sort-hashmap-java-by-key-and-value.html

```

```

*/
public static <K extends Comparable, V extends Comparable> Map<K, V>
sortByValues(Map<K, V> map) {
    List<Map.Entry<K, V>> entries = new LinkedList<Map.Entry<K, V>>(map.entrySet());
    Collections.sort(entries, new Comparator<Map.Entry<K, V>>() {
        @Override
        public int compare(Map.Entry<K, V> o1, Map.Entry<K, V> o2) {
            return o2.getValue().compareTo(o1.getValue());
        }
    });
    //LinkedHashMap will keep the keys in the order they are inserted

    //which is currently sorted on natural ordering

    Map<K, V> sortedMap = new LinkedHashMap<K, V>();
    for (Map.Entry<K, V> entry : entries) {
        sortedMap.put(entry.getKey(), entry.getValue());
    }
    return sortedMap;
}
}

```

output:

```

hduser@bmsce-Precision-T1700:~$ start-all.sh
This script is Deprecated. Instead use start-dfs.sh and start-yarn.sh
Starting namenodes on [localhost]
hduser@localhost's password:
localhost: starting namenode, logging to /usr/local/hadoop/logs/hadoop-hduser-namenode-bmsce-Precision-T1700.ou
hduser@localhost's password:
localhost: starting datanode, logging to /usr/local/hadoop/logs/hadoop-hduser-datanode-bmsce-Precision-T1700.ou
Starting secondary namenodes [0.0.0.0]
hduser@0.0.0.0's password:
0.0.0.0: starting secondarynamenode, logging to /usr/local/hadoop/logs/hadoop-hduser-secondarynamenode-bmsce-Pr
starting yarn daemons
starting resourcemanager, logging to /usr/local/hadoop/logs/yarn-hduser-resourcemanager-bmsce-Precision-T1700.d
hduser@localhost's password:
localhost: starting nodemanager, logging to /usr/local/hadoop/logs/yarn-hduser-nodemanager-bmsce-Precision-T170
hduser@bmsce-Precision-T1700:~$ jps
8464 Jps
7817 SecondaryNameNode
7419 NameNode
7596 DataNode
7983 ResourceManager
8319 NodeManager
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$ hadoop fs -ls /
Found 27 items
drwxr-xr-x - hduser supergroup 0 2022-06-06 12:35 /CSE
drwxr-xr-x - hduser supergroup 0 2022-06-06 12:23 /FFF
drwxr-xr-x - hduser supergroup 0 2022-06-06 12:36 /LLL
drwxr-xr-x - hduser supergroup 0 2022-06-20 12:06 /amit_bda
drwxr-xr-x - hduser supergroup 0 2022-06-03 14:52 /bharath
drwxr-xr-x - hduser supergroup 0 2022-06-03 14:43 /bharath035
drwxr-xr-x - hduser supergroup 0 2022-06-24 14:54 /chi
drwxr-xr-x - hduser supergroup 0 2022-05-31 10:21 /example
drwxr-xr-x - hduser supergroup 0 2022-06-01 15:13 /foldernew
drwxr-xr-x - hduser supergroup 0 2022-06-06 15:04 /hemang061
drwxr-xr-x - hduser supergroup 0 2022-06-20 15:16 /input_khushil
drwxr-xr-x - hduser supergroup 0 2022-06-03 12:27 /irfan
drwxr-xr-x - hduser supergroup 0 2022-06-22 10:44 /lwde
drwxr-xr-x - hduser supergroup 0 2022-06-22 15:32 /muskan
drwxr-xr-x - hduser supergroup 0 2022-06-22 15:06 /muskan_op
drwxr-xr-x - hduser supergroup 0 2022-06-22 15:35 /muskan_output
drwxr-xr-x - hduser supergroup 0 2022-06-06 15:04 /new_folder
drwxr-xr-x - hduser supergroup 0 2022-05-31 10:26 /one
drwxr-xr-x - hduser supergroup 0 2022-06-24 15:30 /out55
drwxr-xr-x - hduser supergroup 0 2022-06-20 12:17 /output
drwxr-xr-x - hduser supergroup 0 2022-06-24 12:42 /r1
drwxr-xr-x - hduser supergroup 0 2022-06-24 12:24 /rgs
drwxr-xr-x - hduser supergroup 0 2022-06-03 12:08 /saurab
drwxrwxr-x - hduser supergroup 0 2019-08-01 16:19 /tmp
drwxr-xr-x - hduser supergroup 0 2019-08-01 16:03 /user
drwxr-xr-x - hduser supergroup 0 2022-06-01 09:46 /user1
-rw-r--r-- 1 hduser supergroup 2436 2022-06-24 12:17 /wc.jar
hduser@bmsce-Precision-T1700:~$ hadoop fs -mkdir /amit_lab
hduser@bmsce-Precision-T1700:~$ hadoop fs -ls /
Found 28 items

```

```

drwxr-xr-x - hduser supergroup 0 2022-06-22 15:35 /muskan_output
drwxr-xr-x - hduser supergroup 0 2022-06-06 15:04 /new_folder
drwxr-xr-x - hduser supergroup 0 2022-05-31 10:26 /one
drwxr-xr-x - hduser supergroup 0 2022-06-24 15:30 /out55
drwxr-xr-x - hduser supergroup 0 2022-06-20 12:17 /output
drwxr-xr-x - hduser supergroup 0 2022-06-24 12:42 /r1
drwxr-xr-x - hduser supergroup 0 2022-06-24 12:24 /rgs
drwxr-xr-x - hduser supergroup 0 2022-06-03 12:08 /saurab
drwxrwxr-x - hduser supergroup 0 2019-08-01 16:19 /tmp
drwxr-xr-x - hduser supergroup 0 2019-08-01 16:03 /user
drwxr-xr-x - hduser supergroup 0 2022-06-01 09:46 /user1
-rw-r--r-- 1 hduser supergroup 2436 2022-06-24 12:17 /wc.jar
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$ hadoop fs -copyFromLocal /home/hduser/Desktop/sample.txt /amit_lab/file.txt
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$ hadoop fs -ls /amit_lab
Found 1 items
-rw-r--r-- 1 hduser supergroup 51 2022-06-27 11:42 /amit_lab/file.txt
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$ hdfs fs -rmdir /bharath
Error: Could not find or load main class fs
hduser@bmsce-Precision-T1700:~$ hdfs fs -rmdir bharath
Error: Could not find or load main class fs
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$
hduser@bmsce-Precision-T1700:~$ hadoop jar /home/hduser/Desktop/TopN.jar TopN /amit_lab/file.txt /output_Topn
22/06/27 12:14:41 INFO Configuration.deprecation: session.id is deprecated. Instead, use dfs.metrics.session-id
22/06/27 12:14:41 INFO jvm.JvmMetrics: Initializing JVM Metrics with processName=JobTracker, sessionId=
22/06/27 12:14:41 INFO input.FileInputFormat: Total input paths to process : 1
22/06/27 12:14:41 INFO mapreduce.JobSubmitter: number of splits:1

```

```

hduser@bmsce-Precision-T1700:~$ hadoop fs -ls /output_Topn
Found 2 items
-rw-r--r-- 1 hduser supergroup 0 2022-06-27 12:14 /output_Topn/_SUCCESS
-rw-r--r-- 1 hduser supergroup 43 2022-06-27 12:14 /output_Topn/part-r-00000
hduser@bmsce-Precision-T1700:~$ hadoop fs -cat /output_Topn/part-r-00000
bms 2
college 2
computer 1
law 1
science 1

```