**Program One**: Frequencies of Words and their Frequencies. (Only words greater than 5 frequency)

**Data Preparation**: Removed Special characters and converted all to lower case.

**Output**: First and Last 50 words.

```java
1.  // importing the Required Libraries -- Remain same for first and second programs
2.
3.  import java.io.IOException;
4.  import java.util.StringTokenizer;
5.
6.  import org.apache.hadoop.conf.Configuration;
7.  import org.apache.hadoop.fs.Path;
8.  import org.apache.hadoop.io.IntWritable;
9.  import org.apache.hadoop.io.Text;
10. import org.apache.hadoop.mapreduce.Job;
11. import org.apache.hadoop.mapreduce.Mapper;
12. import org.apache.hadoop.mapreduce.Reducer;
13. import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
14. import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
15.
16. public class first_program {
17.
18. // My Mapper Class
19.
20.   public static class TokenizerMapper
21.        extends Mapper<Object, Text, Text, IntWritable>{
22.
23.     private final static IntWritable one = new IntWritable(1);
24.     private Text word = new Text();
25.
26. // Passing the string of tokens to the mapper as (key, value) pair
27.
28.     public void map(Object key, Text value, Context context
29.                    ) throws IOException, InterruptedException {
30.
31.      // Filtering the data for alpha numerics and then converting them to lower case.
32.
33.       String line = value.toString().replaceAll("[^\\p{L}\\p{Z}]","").toLowerCase();
34.       StringTokenizer itr = new StringTokenizer(line);
35.       while (itr.hasMoreTokens()) {
36.         word.set(itr.nextToken());
37.         context.write(word, one);
38.       }
39.     }
40.   }
41.
42. //My Reducer Class
43.
44.   public static class IntSumReducer
45.        extends Reducer<Text,IntWritable,Text,IntWritable> {
46.     private IntWritable result = new IntWritable();
47.
48.     public void reduce(Text key, Iterable<IntWritable> values,
49.                    Context context
50.                    ) throws IOException, InterruptedException {
51.       int sum = 0;
52.       for (IntWritable val : values) {
53.         sum += val.get();
```

```java
54.        }
55.     // Filtering for words with frequency greater than 5
56.        if(sum>5) {
57.        result.set(sum);
58.        context.write(key, result);
59.      }
60.    }
61.   }
62.
63. //My Driver Section
64.
65.   public static void main(String[] args) throws Exception {
66.      Configuration conf = new Configuration();
67.      Job job = Job.getInstance(conf, "word count");
68.      job.setJarByClass(first_program.class);
69.      job.setMapperClass(TokenizerMapper.class);
70.      job.setCombinerClass(IntSumReducer.class);
71.      job.setReducerClass(IntSumReducer.class);
72.      job.setOutputKeyClass(Text.class);
73.      job.setOutputValueClass(IntWritable.class);
74.      FileInputFormat.addInputPath(job, new Path(args[0]));
75.      FileOutputFormat.setOutputPath(job, new Path(args[1]));
76.      System.exit(job.waitForCompletion(true) ? 0 : 1);
77.    }
78. }
```

# First 50 words and their Frequencies:

| | |
|---|---|
| a | 34572 |
| aaron | 30 |
| abandon | 8 |
| abandond | 6 |
| abandoned | 17 |
| abbey | 9 |
| abhorred | 12 |
| abhorrence | 12 |
| abilities | 23 |
| ability | 41 |
| able | 368 |
| aboard | 49 |
| abode | 25 |
| abodes | 27 |
| abolished | 7 |
| abominable | 6 |
| abounding | 7 |
| about | 2455 |
| above | 414 |
| abraham | 46 |
| abroad | 37 |
| absence | 60 |
| absent | 32 |
| absolute | 73 |
| absolutely | 79 |
| absolution | 7 |
| absolved | 6 |

| | |
|---|---|
| absorbed | 7 |
| absurd | 45 |
| absurdities | 9 |
| absurdity | 12 |
| abuse | 15 |
| acamas | 8 |
| accents | 12 |
| accept | 57 |
| acceptable | 6 |
| acceptation | 6 |
| accepted | 48 |
| accepteth | 7 |
| access | 224 |
| accident | 43 |
| accidents | 28 |
| acclaim | 6 |
| accompanied | 21 |
| accompany | 6 |
| accomplished | 28 |
| accord | 6 |
| according | 161 |
| accordingly | 31 |
| account | 258 |

**Last 50 words and their Frequencies:**

| | |
|---|---|
| here | 12 |
| his | 6 |
| how | 10 |
| i | 571 |
| if | 46 |
| ill | 23 |
| im | 13 |
| in | 16 |
| is | 15 |
| it | 120 |
| its | 32 |
| ive | 9 |
| jim | 42 |
| just | 9 |
| look | 8 |
| my | 13 |
| next | 6 |
| no | 22 |
| now | 34 |
| of | 9 |
| oh | 8 |
| one | 14 |
| pap | 6 |
| pretty | 15 |
| say | 6 |
| says | 9 |
| she | 61 |

| | |
|---|---|
| so | 146 |
| some | 10 |
| sometimes | 8 |
| that | 37 |
| thats | 15 |
| the | 141 |
| then | 103 |
| there | 64 |
| these | 6 |
| they | 75 |
| this | 13 |
| tom | 21 |
| we | 111 |
| well | 87 |
| what | 29 |
| whats | 11 |
| when | 43 |
| where | 10 |
| who | 6 |
| why | 63 |
| yes | 13 |
| you | 98 |
| your | 7 |
| | 14 |
| æschere | 6 |
| þær | 8 |
| ćneas | 58 |
| ćtolian | 6 |

**Program Two**: Frequencies of Bigrams and their Frequencies. (Only 2-words greater than 5 frequency). Driver class and Libraries are same as First Program.

**Data Preparation**: Removed Special characters.

**Output**: First and Last 50 2-words.

```java
// Mapper Class

1.
2.     public static class TokenizerMapper
3.            extends Mapper<Object, Text, Text, IntWritable>{
4.
5.         private final static IntWritable one = new IntWritable(1);
6.         private static final Text bigram = new Text();
7.
8.         public void map(Object key, Text value, Context context
9.                         ) throws IOException, InterruptedException {
10.
11. // Filtering the text for any special characters only
12.
13.             String line = value.toString().replaceAll("[^\\p{L}\\p{Z}]","");
14.             String prev = null;
15.
16.             StringTokenizer itr = new StringTokenizer(line);
17.
18.             while (itr.hasMoreTokens())
19.             {
20.                String cur = itr.nextToken();
21.
22.                // Emit only if we have an actual bigram.
23.
24.                if (prev != null)
25.                {
26.                   bigram.set(prev + " " + cur);
27.                   context.write(bigram, one);
28.                }
29.                prev = cur;
30.             }
31.          String prev1 = prev; // concatenating the first word of a line with the last
    word of previous line
32.     }
33.
34. }
35.
36. // Reducer Class
37.
38.    public static class IntSumReducer extends Reducer<Text,IntWritable,Text,IntWritable>
    {
39.
40.        private IntWritable result = new IntWritable();
41.
42.        public void reduce(Text key, Iterable<IntWritable> values,
43.                           Context context
44.                           ) throws IOException, InterruptedException {
45.          int sum = 0;
46.          for (IntWritable val : values) {
47.            sum += val.get();
48.          }
```

```java
49.    // Filtering the words with frequency greater than 5
50.
51.    if (sum >5){
52.
53.        result.set(sum);
54.        context.write(key, result);
55.      }
56.
57.      }
58.
59.    }
60.
61. // Driver Class
62.
63. public class second {
64.
65.    public static void main(String[] args) throws Exception {
66.
67.        Configuration conf = new Configuration();
68.        Job job = Job.getInstance(conf, "word count");
69.        job.setJarByClass(second.class);
70.        job.setMapperClass(TokenizerMapper.class);
71.        job.setCombinerClass(IntSumReducer.class);
72.        job.setReducerClass(IntSumReducer.class);
73.        job.setOutputKeyClass(Text.class);
74.        job.setOutputValueClass(IntWritable.class);
75.        FileInputFormat.addInputPath(job, new Path(args[0]));
76.        FileOutputFormat.setOutputPath(job, new Path(args[1]));
77.        System.exit(job.waitForCompletion(true) ? 0 : 1);
78.
79.    }
80. }
```

## First 50 Bigrams and their Frequencies :

| Bigram | Frequency |
|---|---|
| A Commonwealth | 9 |
| A Law | 6 |
| A Man | 7 |
| A chief | 6 |
| A man | 6 |
| A moment | 6 |
| A mortal | 6 |
| A sudden | 9 |
| ADVENTURE OF | 6 |
| AND OF | 6 |
| AND THE | 18 |
| Above the | 9 |
| According to | 6 |
| Achilles and | 14 |
| Achilles arms | 6 |
| Achilles hand | 7 |
| Achilles in | 6 |
| Achilles thus | 8 |
| Achilles to | 9 |
| Achilles with | 6 |
| Actions and | 6 |
| Acts of | 11 |
| Administration of | 8 |
| Adventure of | 6 |
| After a | 51 |
| After all | 6 |
| After this | 6 |

| | |
|---|---|
| Again I | 6 |
| Again he | 6 |
| Again the | 9 |
| Against The | 7 |
| Against the | 11 |
| Ah said | 16 |
| Ah yes | 6 |
| Ahab and | 19 |
| Ahab did | 7 |
| Ahab had | 12 |
| Ahab in | 9 |
| Ahab is | 11 |
| Ahab now | 8 |
| Ahab said | 6 |
| Ahab stood | 11 |
| Ahab the | 7 |
| Ahab to | 9 |
| Ahab was | 16 |
| Ahab who | 8 |
| Ahab with | 6 |
| Ahab would | 6 |
| Ajax and | 8 |
| Ajax the | 8 |

## Last 50 Bigrams and their Frequencies:

| Bigram | Frequency |
|---|---|
| I reckoned | 9 |
| I said | 13 |
| I says | 35 |
| I see | 11 |
| I set | 6 |
| I thought | 6 |
| I told | 6 |
| I took | 8 |
| I was | 41 |
| I went | 15 |
| I wish | 6 |
| If I | 7 |
| If he | 6 |
| If you | 9 |
| It aint | 6 |
| It warnt | 8 |
| It was | 50 |
| Its a | 7 |
| Its the | 7 |
| Jim said | 13 |
| Of course | 8 |
| Pretty soon | 12 |
| She said | 9 |
| She says | 6 |
| She was | 11 |
| So I | 40 |
| So he | 7 |
| So the | 9 |

| | |
|---|---|
| So then | 12 |
| So we | 14 |
| That was | 12 |
| The duke | 7 |
| The king | 13 |
| The old | 8 |
| Then I | 20 |
| Then he | 20 |
| Then the | 15 |
| Then we | 12 |
| There aint | 10 |
| There warnt | 11 |
| There was | 33 |
| They all | 7 |
| We got | 6 |
| We was | 6 |
| We went | 6 |
| Well I | 19 |
| Well then | 6 |
| When I | 11 |
| When we | 11 |
| You see | 6 |

> **Program Three**: 100 Most Frequent words. 2 phases of map and
> reduce used where the second mapper swaps the key-value pairs
> and passes to the reducer. Comparator to get the output in
> descending order.
>
> **Data Preparation**: Filtered for only words.

```
1.  // Libraries - Additional are the Job control libraries
2.
3.  import org.apache.hadoop.conf.Configuration;
4.  import org.apache.hadoop.conf.Configured;
5.  import org.apache.hadoop.fs.Path;
6.  import org.apache.hadoop.io.IntWritable;
7.  import org.apache.hadoop.io.Text;
8.  import org.apache.hadoop.mapreduce.Job;
9.  import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
10. import org.apache.hadoop.mapreduce.lib.input.KeyValueTextInputFormat;
11. import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
12. import org.apache.hadoop.util.Tool;
13. import org.apache.hadoop.util.ToolRunner;
14. import org.apache.hadoop.mapreduce.lib.jobcontrol.ControlledJob;
15. import org.apache.hadoop.mapreduce.lib.jobcontrol.JobControl;
16.
17. // First Mapper Class
18.
19. public class WordMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
20.
21.   IntWritable intWritable = new IntWritable(1);
22.   Text text = new Text();
23.
24.   @Override
25.   public void map(LongWritable key, Text value, Context context) throws IOException, In
    terruptedException {
26.
27. // Filtering the text without numbers or spcial characters
28.
29.     String line = value.toString().replaceAll("[^\\p{L}\\p{Z}]","");
30.
31.     for (String word : line.split("\\W+")) {
32.       if (word.length() > 0) {
33.         text.set(word);
34.         context.write(text, intWritable);
35.       }
36.     }
37.   }
38. }
39.
40. // First Reducer Class
41.
42.     public class SumReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
43.
44.         IntWritable intWritable = new IntWritable();
45.
46.         @Override
47.         public void reduce(Text key, Iterable<IntWritable> values, Context context) thr
    ows IOException, InterruptedException {
48.
49.             int wordCount = 0;
50.             for (IntWritable value : values) {
```

```
51.                    wordCount += value.get();
52.                }
53.              intWritable.set(wordCount);
54.              context.write(key, intWritable);
55.            }
56.          }
57.
58. // Second Mapper Class -- Swaps the output of first reducer
59.
60.      public class WordMapper2 extends Mapper< Text, Text, IntWritable, Text> {
61.
62.           IntWritable frequency = new IntWritable();
63.
64.            @Override
65.            public void map(Text key, Text value, Context context)
66.              throws IOException, InterruptedException {
67.
68.              int newVal = Integer.parseInt(value.toString());
69.              frequency.set(newVal);
70.              context.write(frequency, key);
71.            }
72.          }
73.
74. // Second reducer CLass
75.
76.      public class SumReducer2 extends Reducer<IntWritable, Text, IntWritable, Text> {
77.
78.           Text word = new Text();
79.
80.            @Override
81.            public void reduce(IntWritable key, Iterable<Text> values, Context context)
82.                 throws IOException, InterruptedException {
83.
84.              for (Text value : values) {
85.                  word.set(value);
86.                  context.write(key, word);
87.              }
88.            }
89.          }
90.
91.  // Comparator Class – For Descending Order.
92.
93. package stubs;
94. import java.nio.ByteBuffer;
95. import org.apache.hadoop.io.IntWritable;
96. import org.apache.hadoop.io.WritableComparator;
97.
98. public class IntComparator extends WritableComparator {
99.
100.          public IntComparator() {
101.            super(IntWritable.class);
102.          }
103.
104.          @Override
105.          public int compare(byte[] b1, int s1, int l1, byte[] b2,
106.               int s2, int l2) {
107.            Integer v1 = ByteBuffer.wrap(b1, s1, l1).getInt();
108.            Integer v2 = ByteBuffer.wrap(b2, s2, l2).getInt();
109.            return v1.compareTo(v2) * (-1);
110.          }
111.        }
```

```java
112.        // Driver Class
113.
114.        public class WordCombined extends Configured implements Tool {
115.
116.         public int run(String[] args) throws Exception {
117.
118.            JobControl jobControl = new JobControl("jobChain");
119.            Configuration conf1 = getConf();
120.
121.            Job job1 = Job.getInstance(conf1);
122.            job1.setJarByClass(WordCombined.class);
123.            job1.setJobName("Word Combined");
124.
125.            FileInputFormat.setInputPaths(job1, new Path(args[0]));
126.            FileOutputFormat.setOutputPath(job1, new Path(args[1] + "/temp"));
127.
128.        //Job 1
129.            job1.setMapperClass(WordMapper.class);
130.            job1.setReducerClass(SumReducer.class);
131.            job1.setCombinerClass(SumReducer.class);
132.
133.            job1.setOutputKeyClass(Text.class);
134.            job1.setOutputValueClass(IntWritable.class);
135.
136.            ControlledJob controlledJob1 = new ControlledJob(conf1);
137.            controlledJob1.setJob(job1);
138.
139.            jobControl.addJob(controlledJob1);
140.            Configuration conf2 = getConf();
141.
142.        //Job2
143.
144.            Job job2 = Job.getInstance(conf2);
145.            job2.setJarByClass(WordCombined.class);
146.            job2.setJobName("Word Invert");
147.
148.            FileInputFormat.setInputPaths(job2, new Path(args[1] + "/temp"));
149.            FileOutputFormat.setOutputPath(job2, new Path(args[1] + "/final"));
150.
151.            job2.setMapperClass(WordMapper2.class);
152.            job2.setReducerClass(SumReducer2.class);
153.            job2.setCombinerClass(SumReducer2.class);
154.
155.            job2.setOutputKeyClass(IntWritable.class);
156.            job2.setOutputValueClass(Text.class);
157.            job2.setInputFormatClass(KeyValueTextInputFormat.class);
158.
159.            job2.setSortComparatorClass(IntComparator.class);
160.            ControlledJob controlledJob2 = new ControlledJob(conf2);
161.            controlledJob2.setJob(job2);
162.
163.            // make job2 dependent on job1
164.            controlledJob2.addDependingJob(controlledJob1);
165.
166.            // add the job to the job control
167.            jobControl.addJob(controlledJob2);
168.            Thread jobControlThread = new Thread(jobControl);
169.            jobControlThread.start();
170.
171.        // Printing the Job States
172.
```

```java
173.        while (!jobControl.allFinished()) {
174.        System.out.println("Jobs in waiting state: " + jobControl.getWaitingJobList)
     .size());
175.            System.out.println("Jobs in ready state: " + jobControl.getReadyJobsList().s
     ize());
176.            System.out.println("Jobs in running state: " + jobControl.getRunningJobList(
     ).size());
177.            System.out.println("Jobs in success state: " + jobControl.getSuccessfulJobLi
     st().size());
178.            System.out.println("Jobs in failed state: " + jobControl.getFailedJobList().
     size());
179.        try {
180.            Thread.sleep(5000);
181.            } catch (Exception e) {
182.
183.            }
184.
185.        }
186.         System.exit(0);
187.         return (job1.waitForCompletion(true) ? 0 : 1);
188.        }
189.        public static void main(String[] args) throws Exception {
190.        int exitCode = ToolRunner.run(new WordCombined(), args);
191.        System.exit(exitCode);
192.        }
193.    }
194.
```

## 100 Most Frequent Words with their Frequencies:

| | |
|---|---|
| 99423 | the |
| 58906 | and |
| 52249 | of |
| 47666 | to |
| 33190 | a |
| 28356 | in |
| 27772 | I |
| 23083 | that |
| 18836 | it |
| 18387 | was |
| 18232 | his |
| 17844 | he |
| 14776 | is |
| 14173 | with |
| 12982 | as |
| 12369 | for |
| 12031 | you |
| 11447 | not |
| 10698 | be |
| 10106 | had |
| 9407 | s |
| 9397 | on |
| 9385 | by |
| 9168 | but |
| 9107 | The |
| 9042 | him |
| 9035 | all |

| | |
|---|---|
| 8924 | at |
| 8362 | her |
| 7943 | have |
| 7937 | my |
| 7568 | from |
| 7510 | or |
| 7416 | they |
| 7202 | this |
| 7117 | which |
| 7001 | me |
| 6873 | so |
| 5915 | she |
| 5798 | said |
| 5641 | are |
| 5489 | their |
| 5484 | And |
| 5461 | one |
| 5329 | were |
| 5238 | them |
| 5139 | no |
| 5102 | there |
| 4774 | we |
| 4755 | when |
| 4392 | out |
| 4361 | t |
| 4090 | an |
| 4080 | would |
| 4077 | if |
| 3986 | up |

| 3984 | d |
|------|------|
| 3821 | He |
| 3819 | But |
| 3777 | been |
| 3742 | what |
| 3693 | will |
| 3623 | any |
| 3566 | man |
| 3540 | more |
| 3517 | then |
| 3419 | could |
| 3306 | into |
| 3283 | some |
| 3280 | do |
| 3258 | who |
| 3166 | other |
| 3140 | It |
| 3119 | your |
| 2968 | now |
| 2952 | time |
| 2838 | can |
| 2781 | very |
| 2668 | such |
| 2643 | upon |
| 2625 | may |
| 2567 | down |
| 2559 | see |
| 2505 | like |
| 2500 | than |

2450    before

2419    shall

2401    our

2388    little

2375    about

2360    must

2306    has

2252    know

2238    did

2227    over

2215    Mr

2170    only

2081    should

2070    men

2033    again


**References**:

https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html

https://stackoverflow.com/questions/2499585/chaining-multiple-mapreduce-jobs-in-hadoop

https://coe4bd.github.io/HadoopHowTo/multipleJobsSingle/multipleJobsSingle.html

**COMMANDS USED :**

**Compiling Java file:**

$hadoop com.sun.tools.javac.Main <filename.java> -d <destination directory>

**Creating Jar file:**

$jar -cvf <filename.jar> -C <destination directory> <target directory>

**Putting file to hdfs:**

$hadoop fs -put <filename> <destination directory on hdfs>

**Running job on hadoop:**

hadoop jar <filename.jar> <filename without extension> <input data location on hdfs> <output data location on hdfs>

**Output of first 50 words :**

hadoop       fs -cat <output_file> | head -50 >  output_to local_file

**Output of last 50 words :**

hadoop       fs -cat <output_file> | head -50 >  output_to local_file