# Big Data and Large Scale Computing

Lab Report -02

July 21, 2025

**Name : Rohan Baghel**
**Student ID: 202116011**

# Question 1

To modification the "WordCount.java" program and run the same input file as in 1st lab using hadoop libraries and observe the output.

**Answer :**

We have done some modification in "WordCount.java" program to get the output format as.

- Tokenizing strings based on \t\n\r\f , . : ; ? ! [ ] '

- All the word are in small letters

- Counting and storing only those words which occurs more than 4 times in the document.

Modification:

1. In place of StringTokenizer itr = new StringTokenizer(line); use
   StringTokenizer itr = new StringTokenizer(line,"\t\n\r\f , . : ; ? ! [ ] ' ");

2. In place of word.set(itr.nextToken()); use
   word.set(itr.nextToken().toLowerCase());

3. While storing data add a condition
   if (sum > 4) output.collect(key, new IntWritable(sum)); it will sore the words which occurs more than 4 times.

## File Name: WordCount.java

This java WordCount.java program is modified version on the the lab1 WordCount;java program.
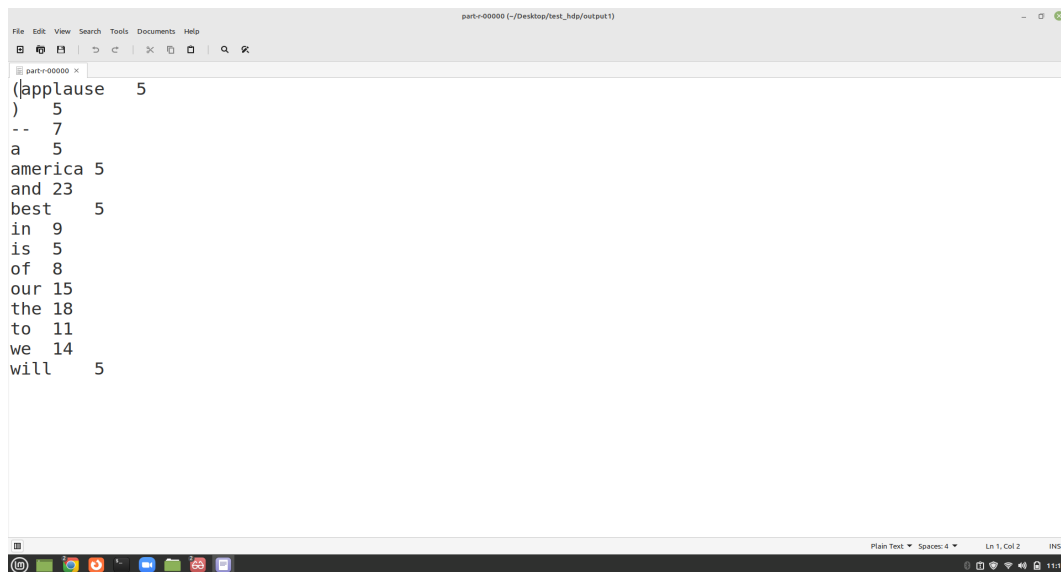
```java
import java.io.IOException;
import java.util.StringTokenizer;

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.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

  public static class TokenizerMapper
       extends Mapper<Object, Text, Text, IntWritable>{

    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(Object key, Text value, Context context
                    ) throws IOException, InterruptedException {
      StringTokenizer itr = new StringTokenizer(value.toString()," \t\n\r\f...
    ,.:;?![]'");
      while (itr.hasMoreTokens()) {
        word.set(itr.nextToken().toLowerCase());
        context.write(word, one);
      }
    }
  }

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

    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context
                       ) throws IOException, InterruptedException {
      int sum = 0;
      for (IntWritable val : values) {
        sum += val.get();
      }
      result.set(sum);
      if(sum>4) context.write(key, result);
    }
  }

  public static void main(String[] args) throws Exception {
    Configuration conf = new Configuration();
    Job job = Job.getInstance(conf, "word count");
    job.setJarByClass(WordCount.class);
```

```
52    job.setMapperClass(TokenizerMapper.class);
53    job.setCombinerClass(IntSumReducer.class);
54    job.setReducerClass(IntSumReducer.class);
55    job.setOutputKeyClass(Text.class);
56    job.setOutputValueClass(IntWritable.class);
57    FileInputFormat.addInputPath(job, new Path(args[0]));
58    FileOutputFormat.setOutputPath(job, new Path(args[1]));
59    System.exit(job.waitForCompletion(true) ? 0 : 1);
60  }
61 }
```
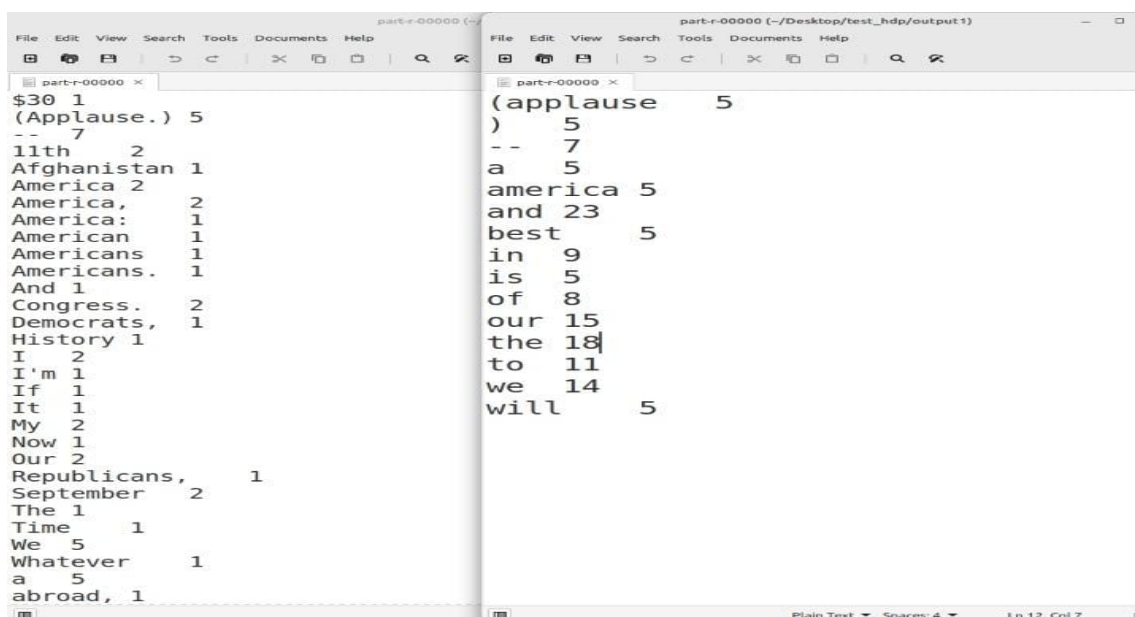
## Output



Figure 1: Output file



Figure 2: Difference in both output files

3

# Question 2

Working on patent data set and observe the output
Dataset Download link :-
https://www.nber.org/research/data/us-patents

## Answer :

### Finding Citing Patents
The data set is of size 264.1 MB
Name : cite75_99.txt

These data comprise detail information on almost 3 million U.S. patents granted between January 1963 and December 1999, all citations made to these patents between 1975 and 1999 (over 16 million), and a reasonably broad match of patents to Compustat (the data set of all firms traded in the U.S. stock market).

## File Name: FindCitingPatents.java

```java
1
2 import java.io.IOException;
3
4 import org.apache.hadoop.conf.Configuration;
5 import org.apache.hadoop.conf.Configured;
6 import org.apache.hadoop.fs.Path;
7 import org.apache.hadoop.io.LongWritable;
8 import org.apache.hadoop.io.Text;
9 import org.apache.hadoop.mapreduce.Job;
10 import org.apache.hadoop.mapreduce.Mapper;
11 import org.apache.hadoop.mapreduce.Reducer;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14 import org.apache.hadoop.util.Tool;
15 import org.apache.hadoop.util.ToolRunner;
16
17 public class FindCitingPatents extends Configured implements Tool {
18
19     public static enum Counters {
20         TOTAL_CITATIONS,
21         TOTAL_PATENTS
22     }
23
24     // Map inputs: (citing patent, cited patent)
25     // Map outputs: (cited patent, citing patent)
26     public static class MapClass extends Mapper<LongWritable, Text, Text, Text> ...
    {
27
28         private Text citing = new Text();
29         private Text cited = new Text();
30
31         @Override
```
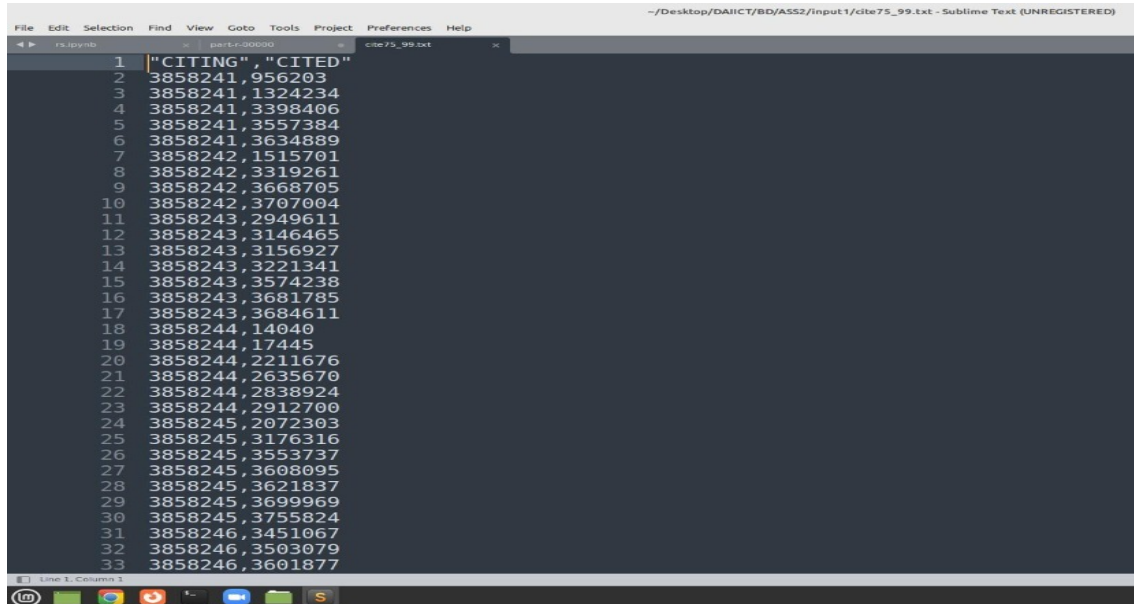
```java
        protected void map(LongWritable key, Text value, Context context)
                throws IOException, InterruptedException {

            String[] split = value.toString().split(",");
            citing.set(split[0]);
            cited.set(split[1]);

            context.write(cited, citing);
            context.getCounter(Counters.TOTAL_CITATIONS).increment(1L);
        }
    }

    // Reduce inputs: (cited patent, list(citing patent))
    // Reduce outputs: (cited patent, CSV of citing patents)
    public static class Reduce extends Reducer<Text, Text, Text, Text> {

        private Text citing = new Text();

        @Override
        protected void reduce(Text key, Iterable<Text> values, Context context)
                throws IOException, InterruptedException {

            StringBuilder builder = new StringBuilder();
            for (Text value : values) {
                if (builder.length() > 0) {
                    builder.append(",");
                }
                builder.append(value.toString());
            }

            citing.set(builder.toString());
            context.write(key, citing);
            context.getCounter(Counters.TOTAL_PATENTS).increment(1L);
        }
    }

    public int run(String[] args) throws Exception {
        Configuration conf = getConf();

        Job job = new Job(conf, FindCitingPatents.class.getSimpleName());
        job.setJarByClass(FindCitingPatents.class);
        job.setMapperClass(MapClass.class);
        job.setReducerClass(Reduce.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(Text.class);

        FileInputFormat.setInputPaths(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        return job.waitForCompletion(true) ? 0 : 1;
    }

    public static void main(String[] args) throws Exception {
        int result = ToolRunner.run(new Configuration(), new FindCitingPatents()...
, args);
        System.exit(result);
    }

}
```

This program will take the patent citation data and invert it. For each patent, we we will find and group the patents that cite it.
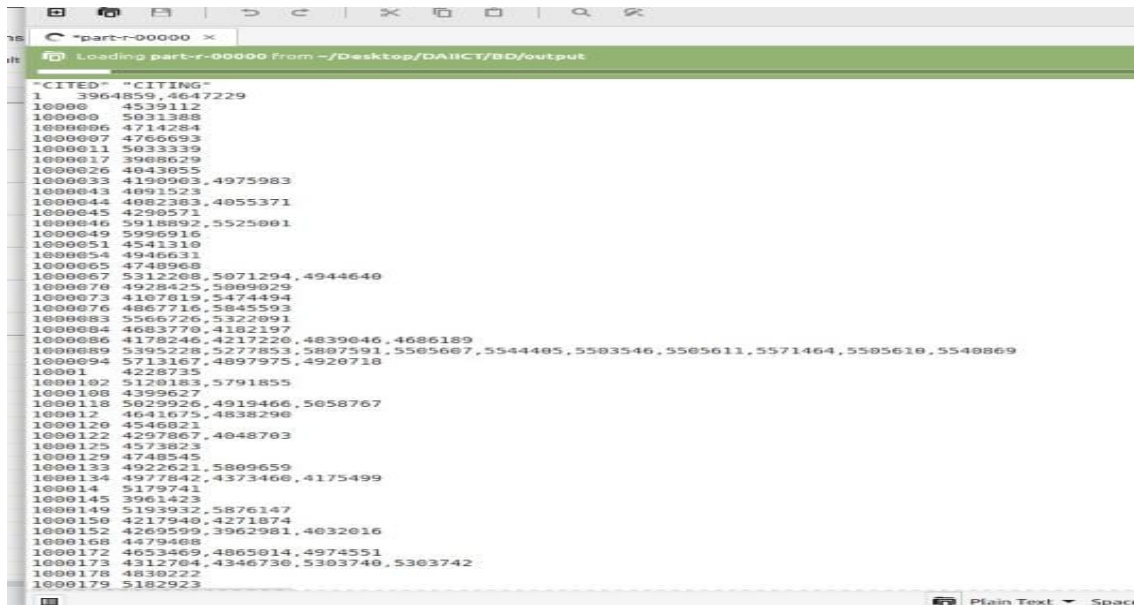
**Input**



Figure 3: Input

**Output**



Figure 4: citation data

The output is same the given in the Book

# Counting Citations of patent

Take input as used in previous input file

As the patent citation data set only covers patents issued between 1975 and 1999.
To build a new program just need some modification in previous program name "Find-CitingPatents.java".
in this java program modifiction required is in Reducer class.

## File Name: FindCitingPatentsCount.java

```java
1  import java.io.IOException;
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.LongWritable;
8  import org.apache.hadoop.io.Text;
9  import org.apache.hadoop.mapreduce.Job;
10 import org.apache.hadoop.mapreduce.Mapper;
11 import org.apache.hadoop.mapreduce.Reducer;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
14 import org.apache.hadoop.util.Tool;
15 import org.apache.hadoop.util.ToolRunner;
16
17 public class FindCitingPatentsCount extends Configured implements Tool {
18
19     public static enum Counters {
20         TOTAL_CITATIONS,
21         TOTAL_PATENTS
22     }
23
24     // Map inputs: (citing patent, cited patent)
25     // Map outputs: (cited patent, citing patent)
26     public static class MapClass extends Mapper<LongWritable, Text, Text, Text> ...
    {
27
28         private Text citing = new Text();
29         private Text cited = new Text();
30
31         @Override
32         protected void map(LongWritable key, Text value, Context context)
33                 throws IOException, InterruptedException {
34
35             String[] split = value.toString().split(",");
36             citing.set(split[0]);
37             cited.set(split[1]);
38
39             context.write(cited, citing);
40             context.getCounter(Counters.TOTAL_CITATIONS).increment(1L);
41         }
42     }
```

```
43
44      // Reduce inputs: (cited patent, list(citing patent))
45      // Reduce outputs: (cited patent, count of citing patents)
46      public static class Reduce extends Reducer<Text, Text, Text, IntWritable> {
47
48          private IntWritable citingCount = new IntWritable();
49
50          @SuppressWarnings({ "UnusedDeclaration" })
51          @Override
52          protected void reduce(Text key, Iterable<Text> values, Context context)
53                  throws IOException, InterruptedException {
54
55              int count = 0;
56              for (Text value : values) {
57                  count++;
58              }
59
60              citingCount.set(count);
61              context.write(key, citingCount);
62              context.getCounter(Counters.TOTAL_PATENTS).increment(1L);
63          }
64      }
65
66      public int run(String[] args) throws Exception {
67          Configuration conf = getConf();
68
69          Job job = new Job(conf, FindCitingPatentsCount.class.getSimpleName());
70          job.setJarByClass(FindCitingPatentsCount.class);
71          job.setMapperClass(MapClass.class);
72          job.setReducerClass(Reduce.class);
73          job.setMapOutputKeyClass(Text.class);
74          job.setMapOutputValueClass(Text.class);
75          job.setOutputKeyClass(Text.class);
76          job.setOutputValueClass(IntWritable.class);
77
78          FileInputFormat.setInputPaths(job, new Path(args[0]));
79          FileOutputFormat.setOutputPath(job, new Path(args[1]));
80
81          return job.waitForCompletion(true) ? 0 : 1;
82      }
83
84      public static void main(String[] args) throws Exception {
85          int result = ToolRunner.run(new Configuration(), new ...
    FindCitingPatentsCount(), args);
86          System.exit(result);
87      }
88
89 }
```
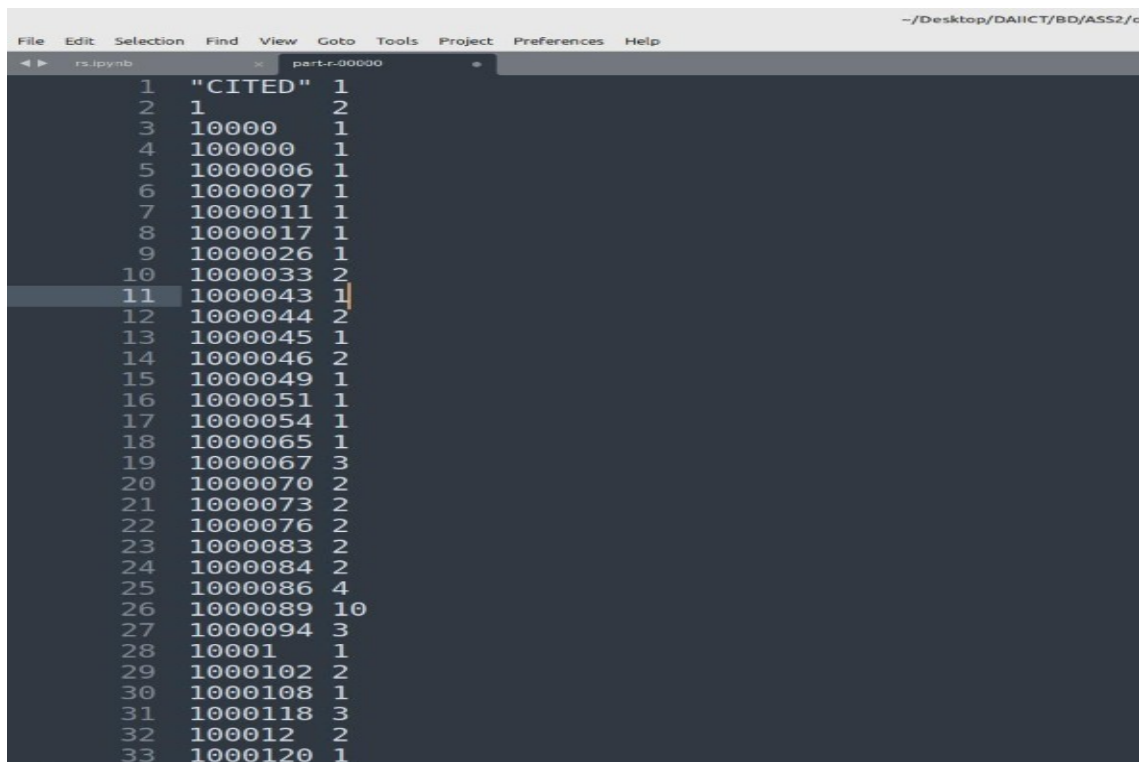
## Output

In each record, a patent number is associated with the number of citations it has received.

## Output File



Figure 5: count patent cited

## Citation Histogram

In this program we are taking input as output of the previous(patent citation count data) program.

Use the program given below to get the axis to plot histogram.

**Fine Name: CitationHistogram.java**

```
1  import java.io.IOException;
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.LongWritable;
8  import org.apache.hadoop.io.Text;
9  import org.apache.hadoop.mapreduce.Job;
10 import org.apache.hadoop.mapreduce.Mapper;
11 import org.apache.hadoop.mapreduce.Reducer;
12 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
13 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
```

```java
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class CitationHistogram extends Configured implements Tool {

    public static enum Counters {
        TOTAL_CITATIONS,
        TOTAL_PATENTS
    }

    public static class MapClass extends Mapper<LongWritable, Text, IntWritable,...
    IntWritable> {

        private final IntWritable one = new IntWritable(1);
        private IntWritable citationCount = new IntWritable();

        @Override
        protected void map(LongWritable key, Text value, Context context)
                throws IOException, InterruptedException {

            String[] split = value.toString().split("\t");
            citationCount.set(Integer.parseInt(split[1]));
            context.write(citationCount, one);
            context.getCounter(Counters.TOTAL_CITATIONS).increment(1L);
        }
    }

     public static class Reduce extends Reducer<IntWritable, IntWritable, ...
    IntWritable, IntWritable> {

        private IntWritable frequency = new IntWritable();

        @Override
        protected void reduce(IntWritable key, Iterable<IntWritable> values, ...
    Context context)
                throws IOException, InterruptedException {

            int count = 0;
            for (IntWritable value : values) {
                count += value.get();
            }
            frequency.set(count);
            context.write(key, frequency);
            context.getCounter(Counters.TOTAL_PATENTS).increment(1L);
        }
    }

    public int run(String[] args) throws Exception {
        Configuration conf = getConf();

        Job job = new Job(conf, CitationHistogram.class.getSimpleName());
        job.setJarByClass(CitationHistogram.class);
        job.setMapperClass(MapClass.class);
        job.setCombinerClass(Reduce.class);
        job.setReducerClass(Reduce.class);
        job.setOutputKeyClass(IntWritable.class);
        job.setOutputValueClass(IntWritable.class);

        FileInputFormat.setInputPaths(job, new Path(args[0]));
```

```
70        FileOutputFormat.setOutputPath(job, new Path(args[1]));
71
72        return job.waitForCompletion(true) ? 0 : 1;
73    }
74
75    public static void main(String[] args) throws Exception {
76        int result = ToolRunner.run(new Configuration(), new CitationHistogram()...
   , args);
77        System.exit(result);
78    }
79
80 }
```

Running the MapReduce job on the citation count data will show the output shown below.
As we suspect,

A large number (900K+) of patents have only one citation, whereas some have hundreds of citations. The most popular patent has 779 citations.
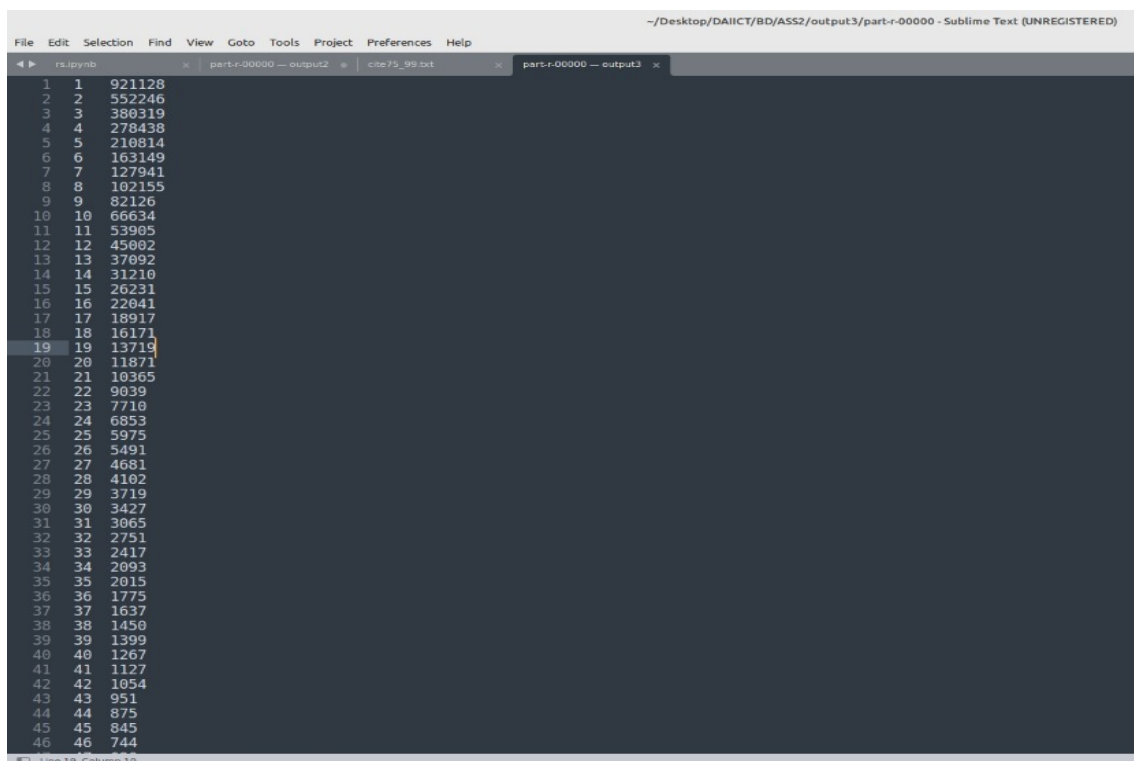
**Output File**



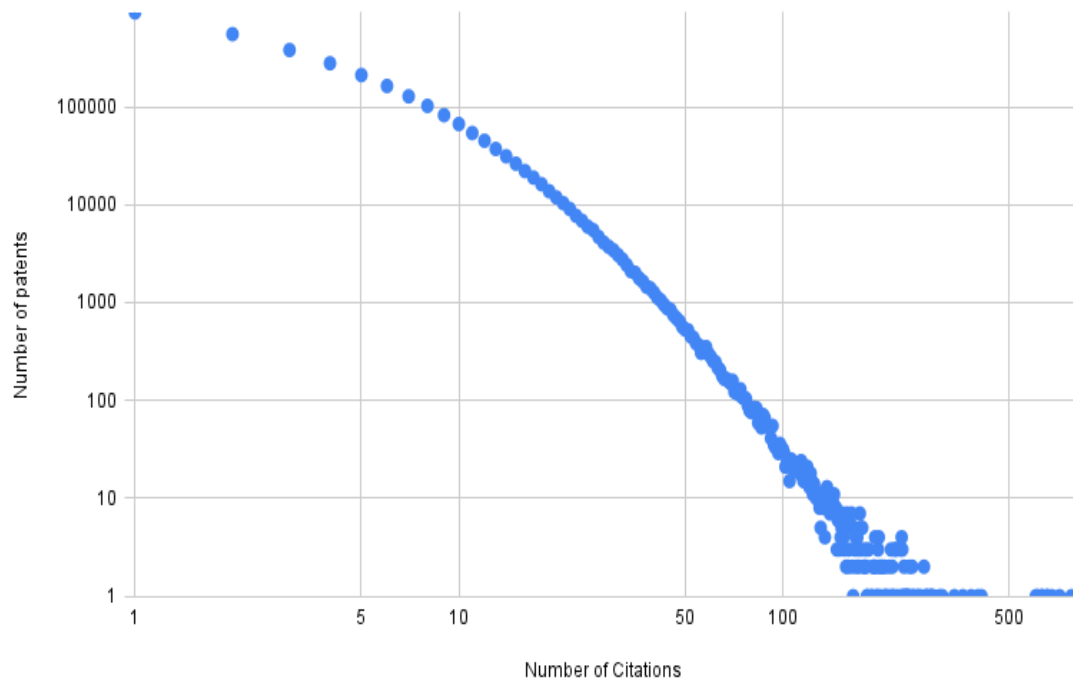Figure 6: Histogram Output

**Histogram**



Figure 7: Histogram

Figure shows the number of patents at various citation frequencies.

The plot is on a log-log scale. When a distribution shows as a line in a log-log plot, it's considered to be a power law distribution .
The citation count histogram seems to fit the description, although its approximately parabolic curvature also suggests a log-normal distribution .

## Observation

With the Help of Hadoop, large data Execute very fast.