

과제1

목표

1. 알파벳 단위로 분산처리를 해보자

간단설명

hadoop을 이용하여 데이터를 처리한뒤, 알파벳 앞글자들의 개수를 분산처리하였다

결과

```
hadoop@ubuntu:~/Project$ hdfs dfs -cat wordcount_test_out/part-r-00000 | more
"      2
(      1
.      3
A     14
B      5
C      1
D      1
E      1
F      5
G      3
H      3
I      7
L      8
M      1
N      4
S      3
T     10
U      2
W      7
hadoop@ubuntu:~/Project$ hdfs dfs -cat wordcount_test_out/part-r-00001 | more
a     140
b      73
c      56
d      32
e      25
f      72
g      23
h      45
i      57
j       6
k       4
l      26
m      31
n      42
o     126
p      54
q       1
r      24
s      75
t     193
u      20
v       6
w      95
y      15
hadoop@ubuntu:~/Project$ s
```

code

```
package ssafy;

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;
import org.apache.hadoop.util.GenericOptionsParser;

public class wordcount1char {
    /*
        Object, Text : input key-value pair type (always same (to get a line of input
        file))
        Text, IntWritable : output key-value pair type
    */
    public static class TokenizerMapper
        extends Mapper<Object,Text,Text,IntWritable> {

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

        // map function (Context -> fixed parameter)
        public void map(Object key, Text value, Context context)
            throws IOException, InterruptedException {

            // value.toString() : get a line
            StringTokenizer itr = new StringTokenizer(value.toString());
            while ( itr.hasMoreTokens() ) {
                word.set(itr.nextToken().substring(0,1));

                // emit a key-value pair
                context.write(word,one);
            }
        }
    }

    /*
        Text, IntWritable : input key type and the value type of input value list
        Text, IntWritable : output key-value pair type
    */
    public static class IntSumReducer
        extends Reducer<Text,IntWritable,Text,IntWritable> {

        // variables
        private IntWritable result = new IntWritable();

        // key : a disticnt word
```

```

        // values : Iterable type (data list)
        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);
            context.write(key,result);
        }
    }

    /* Main function */
    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: <in> <out>");
            System.exit(2);
        }
        Job job = new Job(conf,"word count");
        job.setJarByClass(Wordcount1char.class);

        // let hadoop know my map and reduce classes
        job.setMapperClass(TokenizerMapper.class);
        job.setReducerClass(IntSumReducer.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        // set number of reduces
        job.setNumReduceTasks(2);

        // set input and output directories
        FileInputFormat.addInputPath(job,new Path(otherArgs[0]));
        FileOutputFormat.setOutputPath(job,new Path(otherArgs[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1 );
    }
}

```