# TOPIC:

# <u>Project on NCDC (National Climate Data Center)</u>
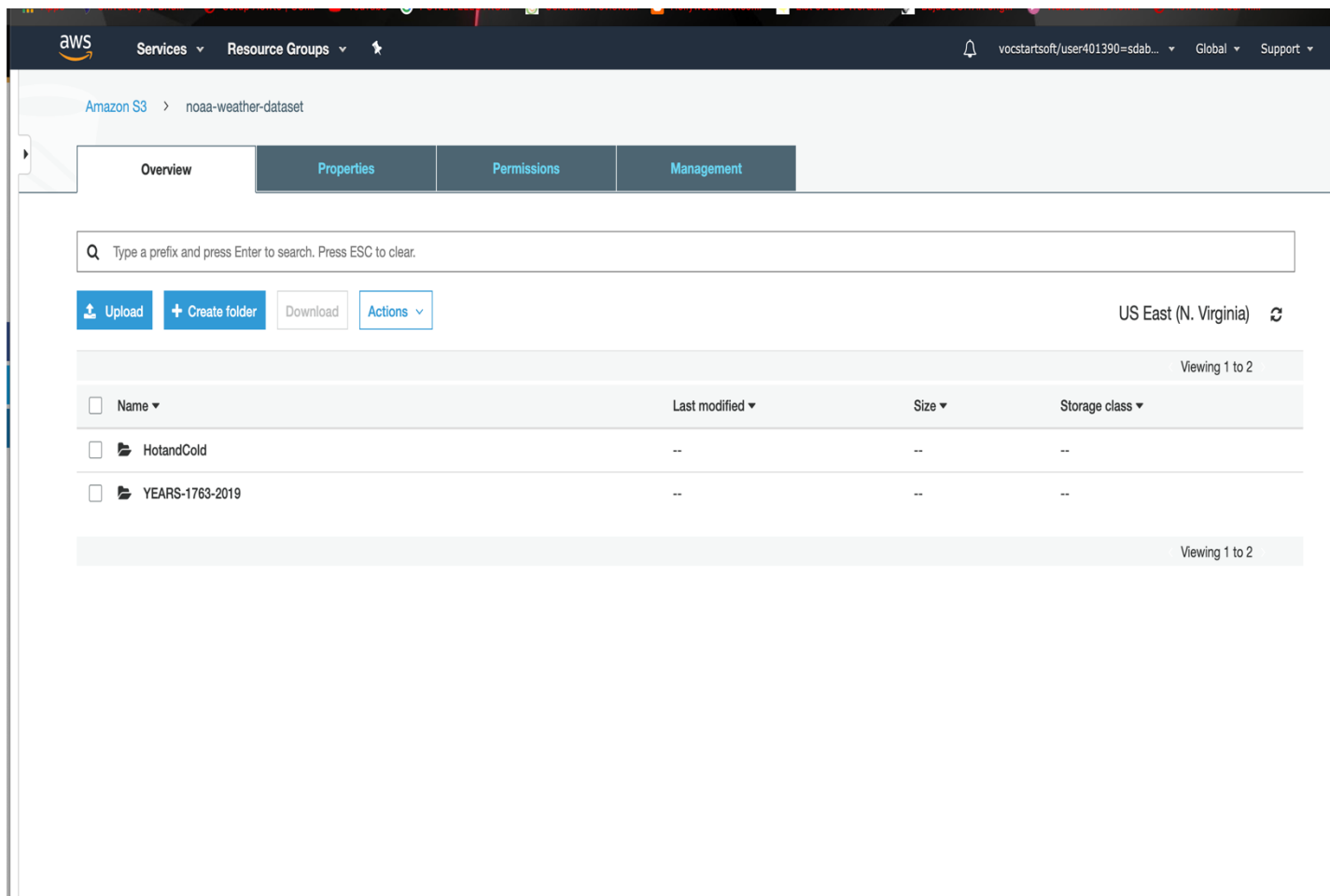
**Submitted by: Srikanth Dabbiru (UB ID 1046112) & Aditya Lavu (UB ID 1066404 )**

- **Testing the Project on Cloud with AWS:**

1. After testing the sample dataset on Cloudera VM, the project solution is set using AWS.

2. All the dataset files uploaded in .csv format to S3 after creation of a bucket by the name noaa-weather-dataset.

3. Data from years 1763-2019 was uploaded along with the source code jar file.

4. Snapshot of the project on AWS :

5. Dataset used in almost 17GB in size with more than 20M instances.

6. Clusters were created on EC2 and initiated after making the necessary changes.

- **<u>Mapper-Reducer Code:</u>**

For the mapper, the max temperature class is static, this method takes the input as text data type. Leaving the first five tokens, the 6$^{th}$ token is taken as the temp_max and the 7$^{th}$ as temp_min.

Now temp_max value is set to be >35.0 and the temp_min is set to be <10.0 and are now passed to the reducer step.

If the temp values for the day are >35.0 output as Hot Day and if <10.0 output as a Cold Day.

For the Reducer method, it takes the input as key and the pairs would be the list of values from the Mapper.

Now **Aggregation** is applied, and it produce the next result.

For the main method, it is used for setting up all the configuration properties. This will be acting as the driver for our Map Reduce code.

- **Below is the Complete Source Code used:**

```java
import java.io.IOException;
import java.util.Iterator;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;

import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.conf.Configuration;

public class MyMaxMin {


    //Mapper



    public static class MaxTemperatureMapper extends
            Mapper<LongWritable, Text, Text, Text> {


        @Override
        public void map(LongWritable arg0, Text Value, Context context)
                throws IOException, InterruptedException {

        //To Convert the record (single line) to String and storing it

            String line = Value.toString();

        //To Check if the line is not empty

            if (!(line.length() == 0)) {

                //date

                String date = line.substring(6, 14);

                //maximum temperature

                float temp_Max = Float
                        .parseFloat(line.substring(39, 45).trim());

                //minimum temperature

                float temp_Min = Float
                        .parseFloat(line.substring(47, 53).trim());

                //if maximum temperature is greater than 35.0 , its a h

                if (temp_Max > 35.0) {
                    // Hot day
                    context.write(new Text("Hot Day " + date),
                            new Text(String.valueOf(temp_Max)));
                }
```

```java
60                 }
61
62                 //if minimum temperature is less than 10.0 , its a cold day
63
64             if (temp_Min < 10) {
65                 // Cold day
66                 context.write(new Text("Cold Day " + date),
67                         new Text(String.valueOf(temp_Min)));
68             }
69         }
70     }
71
72     }
73
74 //Reducer
75
76
77     public static class MaxTemperatureReducer extends
78             Reducer<Text, Text, Text, Text> {
79
80
81         public void reduce(Text Key, Iterator<Text> Values, Context context)
82                 throws IOException, InterruptedException {
83
84
85             //To put all the values in temperature variable of type String
86
87             String temperature = Values.next().toString();
88             context.write(Key, new Text(temperature));
89         }
90
91     }
92
93
94
95     public static void main(String[] args) throws Exception {
96
97         //reads the default configuration of cluster from the configuration xml files
98         Configuration conf = new Configuration();
99
100         //Initializing the job with the default configuration of the cluster
101         Job job = new Job(conf, "weather example");
102
103         //Assigning the driver class name
104         job.setJarByClass(MyMaxMin.class);
105
106         //Key type coming out of mapper
107         job.setMapOutputKeyClass(Text.class);
108
109         //value type coming out of mapper
110         job.setMapOutputValueClass(Text.class);
111
112         //Defining the mapper class name
113         job.setMapperClass(MaxTemperatureMapper.class);
114
115         //Defining the reducer class name
116         job.setReducerClass(MaxTemperatureReducer.class);
117
118         //Defining input Format class which is responsible to parse the dataset into a key value pair
119         job.setInputFormatClass(TextInputFormat.class);
```
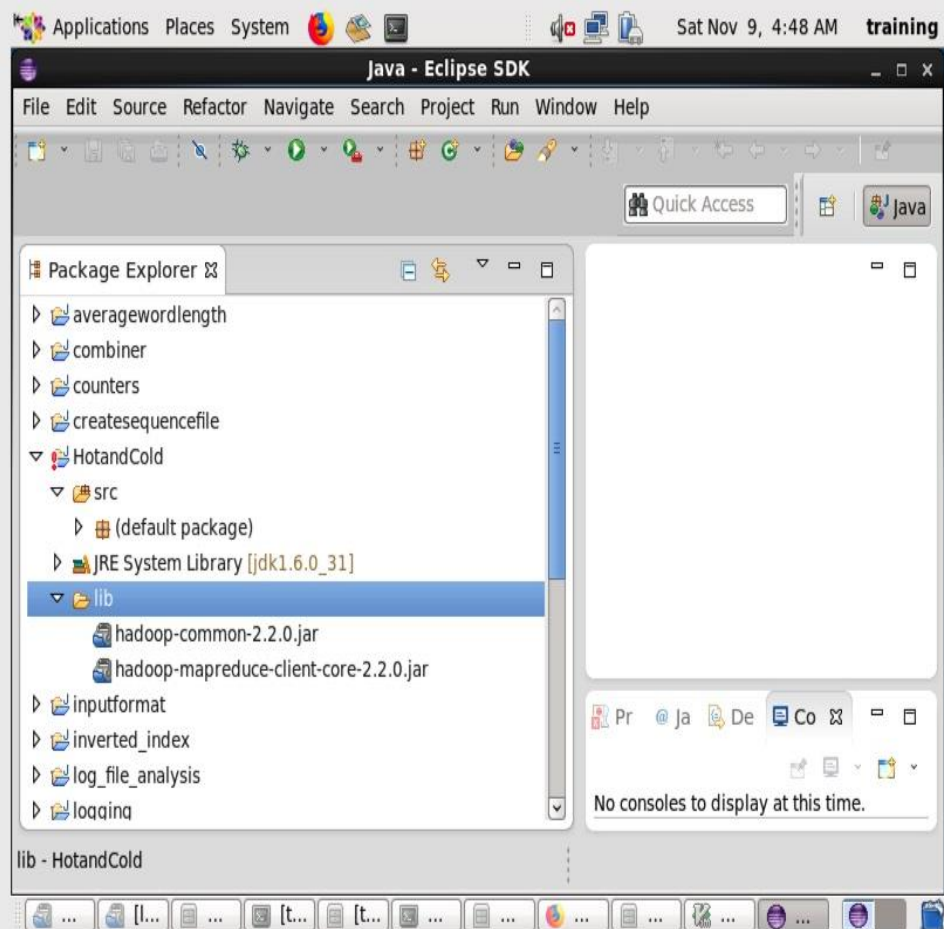
```java
113        job.setMapperClass(MaxTemperatureMapper.class);
114
115        //Defining the reducer class name
116        job.setReducerClass(MaxTemperatureReducer.class);
117
118        //Defining input Format class which is responsible to parse the dataset into a key value pair
119        job.setInputFormatClass(TextInputFormat.class);
120
121        //Defining output Format class which is responsible to parse the dataset into a key value pair
122        job.setOutputFormatClass(TextOutputFormat.class);
123
124        //setting the second argument as a path in a path variable
125        Path OutputPath = new Path(args[1]);
126
127        //Configuring the input path from the filesystem into the job
128        FileInputFormat.addInputPath(job, new Path(args[0]));
129
130        //Configuring the output path from the filesystem into the job
131        FileOutputFormat.setOutputPath(job, new Path(args[1]));
132
133        //deleting the context path automatically from hdfs so that we don't have delete it
134        OutputPath.getFileSystem(conf).delete(OutputPath);
135
136        //exiting the job only if the flag value becomes false
137        System.exit(job.waitForCompletion(true) ? 0 : 1);
138
139    }
140 }
141
142
```

- **<u>Eclipse IDE:</u>**

Now the project is created in the Eclipse IDE to analyze the sample dataset.

The jar file is then exported after having no issues with the project files.

The next step was to send the sample dataset onto HDFS.

**Hdfs dfs -put Downloads/Noaa_Weather_data.txt /**

Run the Jar file for output.

**Hadoop jar temperature.jar /Noaa_weather_data.txt /output_hotandcold**

Check the Output directory in the HDFS.

## Results analysis:

Depending on the 2015 sample dataset only two days above 35.0 recorded.