Week 9)

SPARK :-

1201512025

① using RDD and flatmap count new many times each word appears in file.

```
from pyspark import SparkContent
sc = SparkContent (appname = "word count")

filepath = "path. text"
text-file = sc.text-file (file_path)
words = text-File. flatmap (lamda: line
                                 split())


word-pairs : words. map(lambda)
result = filter-words. collect ()
for word, count in result:
    print (word count)


sc. stop()
```

output!
```
text_rdd = sc.parallelize ([
    "hello world",
    "hello spark",
    "hello hello world spark spark spark"
])
```

[("hello', 4), ('spark', u)]

**②** Java Hadoop Mapreduce Program: top 10 frequent words in Alphabetical order.

## Mapper class

```java
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;


public class WordCountMapper extends Mapper<LongWrite
    -able, Text, text, IntWritable> {
    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();


    public void map(LongWritable key, Text value, Context
        context) throws IOException, Interrupted Exception {
        String[] words = value.toString().toLowerCase().
            split("\\w+");
        for (String w : words) {
            if (!w.isEmpty()) {
                word.set(w);
                context.write(word, one);
            }
        }
    }
}
```

## Reducer class

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

```
Public class WordCountReducer extends Reducer
   <Text, IntWritable, Text, IntWritable> {
   private map<String, Integer> CountMap = new
                                     HashMap<> ();

   Public void reduce(Text key, Iterable < IntWritable>
       values, Context context) {
       int Sum = 0;
       for (IntWritable. val : values) {
           Sum += val.get();
       }
       CountMap. put (key.toString (), Sum);
   } }

   Protected void cleanup(Context context) throws
   IOException, InterruptedException {
   countMap. entryset ().Stream()
       .sorted(Map. Entry. <String. Integer> Comparing
           By Value (Comparator. revers order ())
           . then comparing (Map.Entry. Comparing By Key)

       .limit (10)
       . forEach (entry -> {
       try {
           context.write (new Text(entry.get key ()),
           new IntWritable (entry. get value ());
       } catch (IOException | Interrupted Except e){
           e. printStack Trace ();
       }
   });
 }
}
```

```
Public class Toplowords {
    Public static void main (String[] args),
        throws Exception {

    Configuration conf = new configuration();
    Job job = Job.getInstance(conf, "Top10 m");

    job.setMapOutPutKeyclass (xed.class);
    job.setMapOutPut ValueClass (IntWritable class);
    job.setOutPut Key Class (xed.t class);

    System.exit(job.waitForCompletion(true)? 0: 1);
    }
}
```

Output :-

apple banana apple mango banana orange apple apple banana grape mango apple banana grape.

apple     5
banana    4
grape     2
mango     2
orange    1

③ Cleaning the text!

```
from pyspark import Spark Context
import re.
```

```
def clean_text (rdd):
    clean_rdd = rdd.map (lambda : re.sub)
    return clean_rdd

def main():
    SC = SparkContext (appName = "Text clean")
    SSC = StreamingContext (SC, 1)
    Port = 999
    lines = SSC.socket Text Stream (port)
    cleaned_lines = clean_text (lines)
    cleaned_lines.pprint()
    SSC.start()
    SSC.awaitTermination()

if name == "main":
    main()
```

Output :-

This is an example of streaming text
spark is very useful for Big Data Processing

Example Streaming text
spark useful big data Processing.