

Lesson 1: Introduction to the Spark Environment

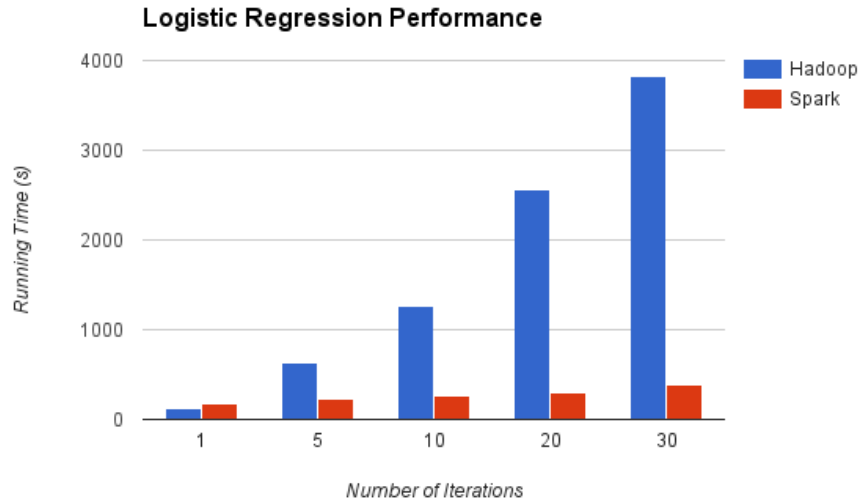
1.4 Why Spark?

Why Spark?

- Handles petabytes of data (and more!)
- Significantly faster than Hadoop Map-Reduce (for most jobs)
- Simple and intuitive APIs
- General Framework
 - Runs Anywhere
 - Handles (most) any I/O
 - Interoperable libraries for specific use-cases



Performance



- Very fast at iterative algorithms
- DAG scheduler supports cyclic flows (and graph computation)
- Intermediate results kept in memory when possible
- Bring computation to the data (data locality)



Rich API



<code>map()</code>	<code>reduce()</code>
<code>filter()</code>	<code>sortBy()</code>
<code>join()</code>	<code>groupByKey()</code>
<code>first()</code>	<code>count()</code>

`map()`
`reduce()`

... and more ...



Rich API



```
text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split())
            .map(lambda word: (word, 1))
            .reduceByKey(lambda a, b: a+b)
```

```
package org.myorg;

import java.io.IOException;
import java.util.*;

import org.apache.hadoop.fs.Path;
import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class WordCount {

    public static class Map extends Mapper<LongWritable, 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, InterruptedException {
            String line = value.toString();
            StringTokenizer tokenizer = new StringTokenizer(line);
            while (tokenizer.hasMoreTokens()) {
                word.set(tokenizer.nextToken());
                context.write(word, one);
            }
        }
    }

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

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

    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();

        Job job = new Job(conf, "wordcount");

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

        job.setMapperClass(Map.class);
        job.setReducerClass(Reduce.class);

        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

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

        job.waitForCompletion(true);
    }
}
```



Unified Platform

Spark SQL

Spark
Streaming

PySpark
SparkR

MLlib
spark.ml

GraphX

Spark Core

Standalone Scheduler

YARN

Mesos



Infrastructure



Data Sources



TACHYON



Review

- Framework for distributed processing
- In-memory, fault tolerant data structures
- Flexible APIs in Scala, Java, Python, SQL... and now R!
- Open Source



Next Up: Installation

