

Questions

1 Setting up a stand-alone Spark instance

- Download and install *Spark 2.2.1* on your machine: <https://www.apache.org/dyn/closer.lua/spark/spark-2.2.1/spark-2.2.1-bin-hadoop2.7.tgz>
- Unpack the compressed TAR ball.

Spark requires JDK 8, which is installed by default on the rice.stanford.edu clusters. **Do not install JDK 9**; Spark is currently incompatible with JDK 9. If you need to download the JDK, please visit Oracle's download site: <http://www.oracle.com/technetwork/java/javase/downloads/index.html>. If you plan to use Scala, you will also need Scala 2.11, and if you plan to use python, you will need python 2.7 or higher (which is also preinstalled on rice) or 3.4 or higher.

To make Spark, Scala, Maven, etc. work on rice, we had to add the following to our .profile file:

```
JAVA_HOME="/usr/bin/java"
SCALA_HOME="$HOME/scala-2.12.4"
SPARK_HOME="$HOME/spark-2.2.1-bin-hadoop2.7"
SPARK_LOCAL_IP="127.0.0.1"
PATH="$HOME/bin:$HOME/.local/bin:$SCALA_HOME/bin:$SPARK_HOME/bin:
$HOME/apache-maven-3.5.2/bin:$PATH"
```

These commands are just guidelines, and what you have to add to your file may vary depending on your specific computer setup.

2 Running the Spark shell

Spark gives you two different ways to run your applications. The easiest is using the Spark shell, a REPL that let's you interactively compose your application. The Spark shell supports two languages: Scala/Java and python. In this tutorial we will only discuss using python and Scala. We highly recommend python as both the language itself and the python Spark API are straightforward.

2.1 Spark Shell for Python

To start the Spark shell for python, do the following:

1. Open a terminal window on Mac or Linux or a command window on Windows.
2. Change into the directory where you unpacked the Spark binary.
3. Run: `bin/pyspark` on Mac or Linux or `bin\pyspark` on Windows.

As the Spark shell starts, you may see large amounts of logging information displayed on the screen, possibly including several warnings. You can ignore that output for now. Regardless, the startup is complete when you see something like:

Welcome to

```

      _--_
     /  _/_-_-_-_-_/  /  _-
    _\  \/_- \/_-' /  _/'  _/
   /__ / .__/\_,_/_/ /_\_\_ version 2.2.1
       /_/_

```

```
Using Python version 2.7.10 (default, Feb  7 2017 00:08:15)
SparkSession available as 'spark'.
>>>
```

The Spark shell is a full python interpreter and can be used to write and execute regular python programs. For example:

```
>>> print "Hello!"
Hello!
```

The Spark shell can also be used to write Spark applications in python. (Surprise!) To learn about writing Spark applications, please read through the Spark programming guide: <https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html>

2.2 Spark Shell for Scala

To start the Spark shell for Scala, do the following:

1. Open a terminal window on Mac or Linux or a command window on Windows.
2. Change into the directory where you unpacked the Spark binary.
3. Run: `bin/spark-shell` on Mac or Linux or `bin\spark-shell` on Windows.

As the Spark shell starts, you may see large amounts of logging information displayed on the screen, possibly including several warnings. You can ignore that output for now. Regardless, the startup is complete when you see something like:

```

      _ _ _ _ _      _ _
     /  _  /  _ _ _ _ _ /  /  _ _
    _ \  \ /  _ \ /  _ ' /  _ /  ' _ /
 / _ _ /  . _ _ \  _ , _ / _ /  / _ \  _
      /  /

```

version 2.2.1

```
scala>
```

```
scala> print("Hello!")
Hello!
```

3 Submitting Spark applications

The Spark shell is great for exploring a data set or experimenting with the API, but it's often best to write your Spark applications outside of the Spark interpreter using an IDE or other smart editor. One of the advantages of this approach for this class is that you will have created a submittable file that contains your application, rather than having to piece it together from the Spark shell's command history. Spark accepts applications written in four languages: Scala, Java, python, and R. In this tutorial we will discuss using Java, Scala, and python. If you choose to use R, you will find resources online, but the TA staff may not be able to help you if you run into issues. We highly recommend python as both the language itself and the python Spark API are straightforward.

3.1 Submitting Python Applications

Python is a convenient choice of language as your python application doesn't need to be compiled or linked. Assume you have the following program in a text file called `myapp.py`:

```
import sys
```

```
from pyspark import SparkConf, SparkContext

conf = SparkConf()
sc = SparkContext(conf=conf)
print "%d lines" % sc.textFile(sys.argv[1]).count()
```

This short application opens the file path given as the first argument from the local working directory and prints the number of lines in it. To run this application, do the following:

1. Open a terminal window on Mac or Linux or a command window on Windows.
2. Change into the directory where you unpacked the Spark binary.
3. Run: `bin/spark-submit path/to/myapp.py path/to/pg100.txt` on Mac or Linux or `bin\spark-submit path\to\myapp.py path\to\pg100.txt` on Windows.

(See section 4 for where to find the `pg100.txt` file.)

As Spark starts, you may see large amounts of logging information displayed on the screen, possibly including several warnings. You can ignore that output for now. Regardless, near the bottom of the output you will see the output from the application:

```
17/12/18 11:55:40 INFO TaskSetManager: Finished task 1.0 in stage 0.0 (TID 1)
in 633 ms on localhost (executor driver) (2/2)
17/12/18 11:55:40 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks
have all completed, from pool
17/12/18 11:55:40 INFO DAGScheduler: ResultStage 0 (count at myapp.py:6)
finished in 0.690 s
17/12/18 11:55:40 INFO DAGScheduler: Job 0 finished: count at myapp.py:6,
took 0.841068 s
124787 lines
17/12/18 11:55:40 INFO SparkContext: Invoking stop() from shutdown hook
17/12/18 11:55:40 INFO SparkUI: Stopped Spark web UI at
http://192.168.86.218:4040
17/12/18 11:55:40 INFO MapOutputTrackerMasterEndpoint: MapOutputTrackerMasterEndpoint
stopped!
17/12/18 11:55:41 INFO MemoryStore: MemoryStore cleared
```

Executing the application this way causes it to be run single-threaded. To run the application with 4 threads, launch it as `bin/spark-submit --master 'local[4]' path/to/myapp.py path/to/pg100.txt`. You can replace the “4” with any number. To use as many threads as are available on your system, launch the application as `bin/spark-submit --master 'local[*]' path/to/myapp.py path/to/pg100.txt`.

To learn about writing Spark applications, please read through the Spark programming guide: <https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html>

3.2 Submitting Java Applications

Building a Spark application in Java will require you to have Maven installed: <https://maven.apache.org/install.html>. Follow these steps:

1. Create (or clone) a Maven project. If you're using an IDE, this step can typically be accomplished by choosing to create a new project from the IDE's menus and selecting "Maven project" as the type. If you're not using an IDE, you can follow these steps: <https://maven.apache.org/guides/getting-started/>. Using an IDE is highly recommended. IntelliJ is generally a good choice, though NetBeans and Eclipse will also work.
2. Modify the `pom.xml` file to include the Spark artifact. In an IDE this can typically be accomplished through a menu-driven process or by editing the `pom.xml` file directly. You can find instructions for IntelliJ and Eclipse here: <https://sparktutorials.github.io/2015/04/02/setting-up-a-spark-project-with-maven.html>. If you're not using an IDE, you'll need to add the following snippet to the `pom.xml` file using an editor:

```
<dependencies>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_2.11</artifactId>
    <version>2.2.1</version>
  </dependency>
</dependencies>
```

Your resulting `pom.xml` file should look something like:

```
<?xml version="1.0" encoding="UTF-8"?>
  <project xmlns="http://maven.apache.org/POM/4.0.0"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
      http://maven.apache.org/xsd/maven-4.0.0.xsd">
    <modelVersion>4.0.0</modelVersion>
    <groupId>edu.stanford.cs246</groupId>
    <artifactId>MyApp</artifactId>
    <version>1.0-SNAPSHOT</version>
    <packaging>jar</packaging>
    <properties>
      <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
      <maven.compiler.source>1.8</maven.compiler.source>
      <maven.compiler.target>1.8</maven.compiler.target>
    </properties>
```

```
<dependencies>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_2.11</artifactId>
    <version>2.2.1</version>
  </dependency>
</dependencies>
</project>
```

Note that the exact file contents may vary depending on how you created your project.

3. Create your application. Assume that you have the following Java file in your Maven project:

```
package edu.stanford.cs246;

import org.apache.spark.api.java.JavaSparkContext;
import org.apache.spark.api.java.JavaRDD;
import org.apache.spark.SparkConf;

public class MyApp {
    public static void main(String[] args) throws Exception {
        SparkConf conf = new SparkConf();
        JavaSparkContext sc = new JavaSparkContext(conf);
        System.out.printf("%d lines\n", sc.textFile(args[0]).count());
    }
}
```

This short application opens the file path given as the first argument from the local working directory and prints the number of lines in it.

4. Build your project. Most IDEs provide a way to build a Maven project from the menu or toolbar. You can also build your project from the command line in a terminal or command window by changing to the directory that contains your `pom.xml` file and running: `mvn clean package`. This command will create a JAR file from your project in the `target` directory. Note that the first time you build your Maven project, Maven will try to download all needed dependencies. For Spark that list can be quite long. Be sure you're on a network with reasonable bandwidth and allow yourself enough time.
5. To run your application, do the following:
 - (a) Open a terminal window on Mac or Linux or a command window on Windows.
 - (b) Change into the directory where you unpacked the Spark binary.
 - (c) Run:

```
bin/spark-submit --class edu.stanford.cs246.MyApp
    path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt
```

on Mac or Linux or

```
bin\spark-submit --class edu.stanford.cs246.MyApp
    path\to\MyApp-1.0-SNAPSHOT.jar path\to\pg100.txt
```

on Windows.

(See section 4 for where to find the pg100.txt file.)

As Spark starts, you may see large amounts of logging information displayed on the screen, possibly including several warnings. You can ignore that output for now. Regardless, near the bottom of the output you will see the output from the application:

```
17/12/18 13:18:16 INFO TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0)
in 406 ms on localhost (executor driver) (2/2)
17/12/18 13:18:16 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks
have all completed, from pool
17/12/18 13:18:16 INFO DAGScheduler: ResultStage 0 (count at MyApp.java:10)
finished in 0.441 s
17/12/18 13:18:16 INFO DAGScheduler: Job 0 finished: count at MyApp.java:10,
took 0.694674 s
124787 lines
17/12/18 13:18:16 INFO SparkContext: Invoking stop() from shutdown hook
17/12/18 13:18:16 INFO SparkUI: Stopped Spark web UI at
http://192.168.86.218:4040
17/12/18 13:18:16 INFO MapOutputTrackerMasterEndpoint:
MapOutputTrackerMasterEndpoint stopped!
17/12/18 13:18:16 INFO MemoryStore: MemoryStore cleared
```

Executing the application this way causes it to be run single-threaded. To run the application with 4 threads, launch it as `bin/spark-submit --master 'local[4]' --class edu.stanford.cs246.MyApp path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt`. You can replace the “4” with any number. To use as many threads as are available on your system, launch the application as `bin/spark-submit --master 'local[*]' --class edu.stanford.cs246.MyApp path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt`.

To learn about writing Spark applications, please read through the Spark programming guide: <https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html>

3.3 Submitting Scala Applications

Building a Spark application in Scala will require you to have Maven installed: <https://maven.apache.org/install.html>. Follow these steps:

1. Create (or clone) a Maven project. If you're using an IDE, this step can typically be accomplished by choosing to create a new project from the IDE's menus and selecting "Maven project" as the type. If you're not using an IDE, you can follow these steps: <https://maven.apache.org/guides/getting-started/>. Using an IDE is highly recommended. IntelliJ is generally a good choice, though Eclipse will also work. If you want to use Netbeans, you will need to install the Scala plugin: <https://github.com/dcaoyuan/nbscala>.
2. Modify the `pom.xml` file to include the Spark and Scala artifacts and enable the Scala build process. In an IDE this can typically be accomplished through a menu-driven process or by editing the `pom.xml` file directly. You can find instructions for IntelliJ here: <http://knowdimension.com/en/data/create-a-spark-application-with-Scala-using-maven>. If you're not using an IDE, you'll need to add the following snippet to the `pom.xml` file using an editor:

```
<dependencies>
  <dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_2.11</artifactId>
    <version>2.2.1</version>
  </dependency>
  <dependency>
    <groupId>org.scala-tools</groupId>
    <artifactId>maven-scala-plugin</artifactId>
    <version>2.11</version>
  </dependency>
</dependencies>
<build>
  <plugins>
    <plugin>
      <groupId>org.scala-tools</groupId>
      <artifactId>maven-scala-plugin</artifactId>
      <executions>
        <execution>
          <goals>
            <goal>compile</goal>
            <goal>testCompile</goal>
          </goals>
        </execution>
      </executions>
    </plugin>
  </plugins>
</build>
```



```

        </executions>
    </plugin>
</plugins>
</build>

```

Your resulting pom.xml file should look something like:

```

<?xml version="1.0" encoding="UTF-8"?>
  <project xmlns="http://maven.apache.org/POM/4.0.0"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://maven.apache.org/POM/4.0.0
      http://maven.apache.org/xsd/maven-4.0.0.xsd">
    <modelVersion>4.0.0</modelVersion>
    <groupId>edu.stanford.cs246</groupId>
    <artifactId>MyApp</artifactId>
    <version>1.0-SNAPSHOT</version>
    <packaging>jar</packaging>
    <properties>
      <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
      <maven.compiler.source>1.8</maven.compiler.source>
      <maven.compiler.target>1.8</maven.compiler.target>
    </properties>
    <dependencies>
      <dependency>
        <groupId>org.apache.spark</groupId>
        <artifactId>spark-core_2.11</artifactId>
        <version>2.2.1</version>
      </dependency>
      <dependency>
        <groupId>org.scala-tools</groupId>
        <artifactId>maven-scala-plugin</artifactId>
        <version>2.11</version>
      </dependency>
    </dependencies>
    <build>
      <plugins>
        <plugin>
          <groupId>org.scala-tools</groupId>
          <artifactId>maven-scala-plugin</artifactId>
          <executions>
            <execution>
              <goals>
                <goal>compile</goal>
                <goal>testCompile</goal>
              </goals>
            </execution>
          </executions>
        </plugin>
      </plugins>
    </build>
  </project>

```

```
        </execution>
      </executions>
    </plugin>
  </plugins>
</build>
</project>
```

Note that the exact file contents may vary depending on how you created your project.

3. Create your application. Assume that you have the following Scala file in your Maven project:

```
package edu.stanford.cs246

import org.apache.spark.{SparkConf, SparkContext}
import org.apache.spark.SparkContext._

object MyApp {
  def main(args: Array[String]) {
    val conf = new SparkConf();
    val sc = new SparkContext(conf)
    printf("%d lines\n", sc.textFile(args(0)).count);
  }
}
```

This short application opens the file path given as the first argument from the local working directory and prints the number of lines in it.

4. Build your project. Most IDEs provide a way to build a Maven project from the menu or toolbar. You can also build your project from the command line in a terminal or command window by changing to the directory that contains your `pom.xml` file and running: `mvn clean package`. This command will create a JAR file from your project in the `target` directory. Note that the first time you build your Maven project, Maven will try to download all needed dependencies. For Spark that list can be quite long. Be sure you're on a network with reasonable bandwidth and allow yourself enough time.
5. To run your application, do the following:
 - (a) Open a terminal window on Mac or Linux or a command window on Windows.
 - (b) Change into the directory where you unpacked the Spark binary.
 - (c) Run:

```
bin/spark-submit --class edu.stanford.cs246.MyApp
  path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt
```

on Mac or Linux or

```
bin\spark-submit --class edu.stanford.cs246.MyApp  
path\to\MyApp-1.0-SNAPSHOT.jar path\to\pg100.txt
```

on Windows.

(See section 4 for where to find the pg100.txt file.)

As Spark starts, you may see large amounts of logging information displayed on the screen, possibly including several warnings. You can ignore that output for now. Regardless, near the bottom of the output you will see the output from the application:

```
17/12/18 13:18:16 INFO TaskSetManager: Finished task 0.0 in stage 0.0 (TID 0)  
in 406 ms on localhost (executor driver) (2/2)  
17/12/18 13:18:16 INFO TaskSchedulerImpl: Removed TaskSet 0.0, whose tasks  
have all completed, from pool  
17/12/18 13:18:16 INFO DAGScheduler: ResultStage 0 (count at MyApp.java:10)  
finished in 0.441 s  
17/12/18 13:18:16 INFO DAGScheduler: Job 0 finished: count at MyApp.java:10,  
took 0.694674 s  
124787 lines  
17/12/18 13:18:16 INFO SparkContext: Invoking stop() from shutdown hook  
17/12/18 13:18:16 INFO SparkUI: Stopped Spark web UI at  
http://192.168.86.218:4040  
17/12/18 13:18:16 INFO MapOutputTrackerMasterEndpoint:  
MapOutputTrackerMasterEndpoint stopped!  
17/12/18 13:18:16 INFO MemoryStore: MemoryStore cleared
```

Executing the application this way causes it to be run single-threaded. To run the application with 4 threads, launch it as `bin/spark-submit --master 'local[4]' --class edu.stanford.cs246.MyApp path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt`. You can replace the “4” with any number. To use as many threads as are available on your system, launch the application as `bin/spark-submit --master 'local[*]' --class edu.stanford.cs246.MyApp path/to/MyApp-1.0-SNAPSHOT.jar path/to/pg100.txt`.

To learn about writing Spark applications, please read through the Spark programming guide: <https://spark.apache.org/docs/2.2.0/rdd-programming-guide.html>

4 Word Count

The typical “Hello, world!” app for Spark applications is known as word count. The map/reduce model is particularly well suited to applications like counting words in a document. In