Jeffrey Aven

Sams **Teach Yourself**

# Apache Spark™

in **24**
**Hours**

**SAMS**

Jeffrey Aven

Sams **Teach Yourself**

# Apache Spark™

in **24**

# Hours

## Sams Teach Yourself Apache Spark™ in 24 Hours

### Trademarks

All terms mentioned in this book that are known to be trademarks or service marks have been appropriately capitalized. Sams Publishing cannot attest to the accuracy of this information. Use of a term in this book should not be regarded as affecting the validity of any trademark or service mark.

### Warning and Disclaimer

Every effort has been made to make this book as complete and as accurate as possible, but no warranty or fitness is implied. The information provided is on an "as is" basis. The author and the publisher shall have neither liability nor responsibility to any person or entity with respect to any loss or damages arising from the information contained in this book.

### Special Sales

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact

governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact

intlcs@pearsoned.com.

# Contents at a Glance

# Table of Contents

# Preface

This book assumes nothing, unlike many big data (Spark and Hadoop) books before it, which are often shrouded in complexity and assume years of prior experience. I don't assume that you are a seasoned software engineer with years of experience in Java, I don't assume that you are an experienced big data practitioner with extensive experience in Hadoop and other related open source software projects, and I don't assume that you are an experienced data scientist.

By the same token, you will not find this book patronizing or an insult to your intelligence either. The only prerequisite to this book is that you are "comfortable" with Python. Spark includes several application programming interfaces (APIs). The Python API was selected as the basis for this book as it is an intuitive, interpreted language that is widely known and easily learned by those who haven't used it.

This book could have easily been titled *Sams Teach Yourself Big Data Using Spark* because this is what I attempt to do, taking it from the beginning. I will introduce you to Hadoop, MapReduce, cloud computing, SQL, NoSQL, real-time stream processing, machine learning, and more, covering all topics in the context of how they pertain to Spark. I focus on core Spark concepts such as the Resilient Distributed Dataset (RDD), interacting with Spark using the shell, implementing common processing patterns, practical data engineering/analysis approaches using Spark, and much more.

I was first introduced to Spark in early 2013, which seems like a short time ago but is a lifetime ago in the context of the Hadoop ecosystem. Prior to this, I had been a Hadoop consultant and instructor for several years. Before writing this book, I had implemented and used Spark in several projects ranging in scale from small to medium business to enterprise implementations. Even having substantial exposure to Spark, researching and writing this book was a learning journey for myself, taking me further into areas of Spark that I had not yet appreciated. I would like to take you on this journey as well as you read this book.

Spark and Hadoop are subject areas I have dedicated myself to and that I am passionate about. The making of this book has been hard work but has truly been a labor of love. I hope this book launches your career as a big data practitioner and inspires you to do amazing things with Spark.

# Why Should I Learn Spark?

Spark is one of the most prominent big data processing platforms in use today and is one of the most popular big data open source projects ever. Spark has risen from its roots in academia to Silicon Valley start-ups to proliferation within traditional businesses such as banking, retail, and telecommunications. Whether you are a data analyst, data engineer, data scientist, or data steward, learning Spark will help you to advance your career or embark on a new career in the booming area of big data.

# How This Book Is Organized

This book starts by establishing some of the basic concepts behind Spark and Hadoop, which are covered in Part I, "Getting Started with Apache Spark." I also cover deployment of Spark both locally and in the cloud in Part I.

Part II, "Programming with Apache Spark," is focused on programming with Spark, which includes an introduction to functional programming with both Python and Scala as well as a detailed introduction to the Spark core API.

Part III, "Extensions to Spark," covers extensions to Spark, which include Spark SQL, Spark Streaming, machine learning, and graph processing with Spark. Other areas such as NoSQL systems (such as Cassandra and HBase) and messaging systems (such as Kafka) are covered here as well.

I wrap things up in Part IV, "Managing Spark," by discussing Spark management, administration, monitoring, and logging as well as securing Spark.

# Data Used in the Exercises

Data for the Try It Yourself exercises can be downloaded from the book's Amazon Web Services (AWS) S3 bucket (if you are not familiar with AWS, don't worry—I cover this topic in the book as well). When running the exercises, you can use the data directly from the S3 bucket or you can download the data locally first (examples of both methods are shown). If you choose to download the data first, you can do so from the book's download page at http://sty-spark.s3-website-us-east-1.amazonaws.com/.

# Conventions Used in This Book

Each hour begins with "What You'll Learn in This Hour," which provides a list of bullet points highlighting the topics covered in that hour. Each hour concludes with a "Summary" page summarizing the main points covered in the hour as well as "Q&A" and "Quiz" sections to help you consolidate your learning from that hour.

Key topics being introduced for the first time are typically *italicized* by convention. Most hours also include programming examples in numbered code listings. Where functions, commands, classes, or objects are referred to in text, they appear in `monospace` type.

Other asides in this book include the following:

### NOTE

Content not integral to the subject matter but worth noting or being aware of.

### TIP

**TIP Subtitle**

A hint or tip relating to the current topic that could be useful.

### CAUTION

**Caution Subtitle**

Something related to the current topic that could lead to issues if not addressed.

### ▼ TRY IT YOURSELF

**Exercise Title**

An exercise related to the current topic including a step-by-step guide and descriptions of expected outputs.

# About the Author

**Jeffrey Aven** is a big data consultant and instructor based in Melbourne, Australia. Jeff has an extensive background in data management and several years of experience consulting and teaching in the areas or Hadoop, HBase, Spark, and other big data ecosystem technologies. Jeff has won accolades as a big data instructor and is also an accomplished consultant who has been involved in several high-profile, enterprise-scale big data implementations across different industries in the region.

# Dedication

*This book is dedicated to my wife and three children. I have been burning the candle at both ends during the writing of this book and I appreciate your patience and understanding…*

# Acknowledgments

# We Want to Hear from You

As the reader of this book, *you* are our most important critic and commentator. We value your opinion and want to know what we're doing right, what we could do better, what areas you'd like to see us publish in, and any other words of wisdom you're willing to pass our way.

We welcome your comments. You can email or write to let us know what you did or didn't like about this book—as well as what we can do to make our books better.

*Please note that we cannot help you with technical problems related to the topic of this book.*

When you write, please be sure to include this book's title and author as well as your name and email address. We will carefully review your comments and share them with the author and editors who worked on the book.

E-mail:    feedback@samspublishing.com

Mail:      Sams Publishing
           ATTN: Reader Feedback
           800 East 96th Street
           Indianapolis, IN 46240 USA

# Reader Services

Visit our website and register this book at **informit.com/register** for convenient access to any updates, downloads, or errata that might be available for this book.

*This page intentionally left blank*

# Installing Spark

**What You'll Learn in This Hour:**

▶ What the different Spark deployment modes are

▶ How to install Spark in Standalone mode

▶ How to install and use Spark on YARN

Now that you've gotten through the heavy stuff in the last two hours, you can dive headfirst into Spark and get your hands dirty, so to speak.

This hour covers the basics about how Spark is deployed and how to install Spark. I will also cover how to deploy Spark on Hadoop using the Hadoop scheduler, YARN, discussed in Hour 2.

By the end of this hour, you'll be up and running with an installation of Spark that you will use in subsequent hours.

## Spark Deployment Modes

There are three primary deployment modes for Spark:

▶ Spark Standalone

▶ Spark on YARN (Hadoop)

▶ Spark on Mesos

*Spark Standalone* refers to the built-in or "standalone" scheduler. The term can be confusing because you can have a single machine or a multinode fully distributed cluster both running in Spark Standalone mode. The term "standalone" simply means it does not need an external scheduler.

With Spark Standalone, you can get up an running quickly with few dependencies or environmental considerations. Spark Standalone includes everything you need to get started.

Spark on YARN and Spark on Mesos are deployment modes that use the resource schedulers YARN and Mesos respectively. In each case, you would need to establish a working YARN or Mesos cluster prior to installing and configuring Spark. In the case of Spark on YARN, this typically involves deploying Spark to an existing Hadoop cluster.

I will cover Spark Standalone and Spark on YARN installation examples in this hour because these are the most common deployment modes in use today.

# Preparing to Install Spark

Spark is a cross-platform application that can be deployed on

- ▶ Linux (all distributions)
- ▶ Windows
- ▶ Mac OS X

Although there are no specific hardware requirements, general Spark instance hardware recommendations are

- ▶ 8 GB or more memory
- ▶ Eight or more CPU cores
- ▶ 10 gigabit or greater network speed
- ▶ Four or more disks in *JBOD* configuration (JBOD stands for "Just a Bunch of Disks," referring to independent hard disks not in a RAID—or Redundant Array of Independent Disks—configuration)

Spark is written in *Scala* with programming interfaces in Python (PySpark) and Scala. The following are software prerequisites for installing and running Spark:

- ▶ Java
- ▶ Python (if you intend to use PySpark)

If you wish to use Spark with R (as I will discuss in **Hour 15, "Getting Started with Spark and R"**), you will need to install R as well. Git, Maven, or SBT may be useful as well if you intend on building Spark from source or compiling Spark programs.

If you are deploying Spark on YARN or Mesos, of course, you need to have a functioning YARN or Mesos cluster before deploying and configuring Spark to work with these platforms.

I will cover installing Spark in Standalone mode on a single machine on each type of platform, including satisfying all of the dependencies and prerequisites.

# Installing Spark in Standalone Mode

In this section I will cover deploying Spark in Standalone mode on a single machine using various platforms. Feel free to choose the platform that is most relevant to you to install Spark on.

## Getting Spark

In the installation steps for Linux and Mac OS X, I will use pre-built releases of Spark. You could also download the source code for Spark and build it yourself for your target platform using the build instructions provided on the official Spark website. I will use the latest Spark binary release in my examples. In either case, your first step, regardless of the intended installation platform, is to download either the release or source from: **http://spark.apache.org/downloads.html**

This page will allow you to download the latest release of Spark. In this example, the latest release is 1.5.2, your release will likely be greater than this (e.g. 1.6.x or 2.x.x).



**FIGURE 3.1**
The Apache Spark downloads page.

NOTE

The Spark releases do not actually include Hadoop as the names may imply. They simply include libraries to integrate with the Hadoop clusters and distributions listed. Many of the Hadoop classes are required regardless of whether you are using Hadoop. I will use the `spark-1.5.2-bin-hadoop2.6.tgz` package for this installation.

CAUTION

### Using the "Without Hadoop" Builds

You may be tempted to download the "without Hadoop" or `spark-x.x.x-bin-without-hadoop.tgz` options if you are installing in Standalone mode and not using Hadoop.

The nomenclature can be confusing, but this build is expecting many of the required classes that are implemented in Hadoop to be present on the system. Select this option only if you have Hadoop installed on the system already. Otherwise, as I have done in my case, use one of the `spark-x.x.x-bin-hadoopx.x` builds.

▼ TRY IT YOURSELF

### Install Spark on Red Hat/Centos

In this example, I'm installing Spark on a Red Hat Enterprise Linux 7.1 instance. However, the same installation steps would apply to Centos distributions as well.

1. As shown in Figure 3.1, download the `spark-1.5.2-bin-hadoop2.6.tgz` package from your local mirror into your home directory using `wget` or `curl`.

2. If Java 1.7 or higher is not installed, install the Java 1.7 runtime and development environments using the OpenJDK `yum` packages (alternatively, you could use the Oracle JDK instead):

   ```
   sudo yum install java-1.7.0-openjdk java-1.7.0-openjdk-devel
   ```

3. Confirm Java was successfully installed:

   ```
   $ java -version
   java version "1.7.0_91"
   OpenJDK Runtime Environment (rhel-2.6.2.3.el7-x86_64 u91-b00)
   OpenJDK 64-Bit Server VM (build 24.91-b01, mixed mode)
   ```

4. Extract the Spark package and create SPARK_HOME:

   ```
   tar -xzf spark-1.5.2-bin-hadoop2.6.tgz
   sudo mv spark-1.5.2-bin-hadoop2.6 /opt/spark
   export SPARK_HOME=/opt/spark
   export PATH=$SPARK_HOME/bin:$PATH
   ```

▼

The `SPARK_HOME` environment variable could also be set using the `.bashrc` file or similar user or system profile scripts. You need to do this if you wish to persist the `SPARK_HOME` variable beyond the current session.

5. Open the `PySpark` shell by running the `pyspark` command from any directory (as you've added the Spark `bin` directory to the PATH). If Spark has been successfully installed, you should see the following output (with informational logging messages omitted for brevity):

```
Welcome to

      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 1.5.2
      /_/
Using Python version 2.7.5 (default, Feb 11 2014 07:46:25)
SparkContext available as sc, HiveContext available as sqlContext.
>>>
```

6. You should see a similar result by running the `spark-shell` command from any directory.

7. Run the included Pi Estimator example by executing the following command:

```
spark-submit --class org.apache.spark.examples.SparkPi \
--master local \
$SPARK_HOME/lib/spark-examples*.jar 10
```

8. If the installation was successful, you should see something similar to the following result (omitting the informational log messages). Note, this is an estimator program, so the actual result may vary:

```
Pi is roughly 3.140576
```

## NOTE

Most of the popular Linux distributions include Python 2.*x* with the `python` binary in the system path, so you normally don't need to explicitly install Python; in fact, the `yum` program itself is implemented in Python.

You may also have wondered why you did not have to install Scala as a prerequisite. The Scala binaries are included in the assembly when you build or download a pre-built release of Spark.

▼ TRY IT YOURSELF

## Install Spark on Ubuntu/Debian Linux

In this example, I'm installing Spark on an Ubuntu 14.04 LTS Linux distribution.

As with the Red Hat example, Python 2. 7 is already installed with the operating system, so we do not need to install Python.

1. As shown in Figure 3.1, download the `spark-1.5.2-bin-hadoop2.6.tgz` package from your local mirror into your home directory using `wget` or `curl`.

2. If Java 1.7 or higher is not installed, install the Java 1.7 runtime and development environments using Ubuntu's APT (Advanced Packaging Tool). Alternatively, you could use the Oracle JDK instead:

   ```
   sudo apt-get update
   sudo apt-get install openjdk-7-jre
   sudo apt-get install openjdk-7-jdk
   ```

3. Confirm Java was successfully installed:

   ```
   $ java -version
   java version "1.7.0_91"
   OpenJDK Runtime Environment (IcedTea 2.6.3) (7u91-2.6.3-0ubuntu0.14.04.1)
   OpenJDK 64-Bit Server VM (build 24.91-b01, mixed mode)
   ```

4. Extract the Spark package and create `SPARK_HOME`:

   ```
   tar -xzf spark-1.5.2-bin-hadoop2.6.tgz
   sudo mv spark-1.5.2-bin-hadoop2.6 /opt/spark
   export SPARK_HOME=/opt/spark
   export PATH=$SPARK_HOME/bin:$PATH
   ```

   The `SPARK_HOME` environment variable could also be set using the `.bashrc` file or similar user or system profile scripts. You will need to do this if you wish to persist the `SPARK_HOME` variable beyond the current session.

5. Open the `PySpark` shell by running the `pyspark` command from any directory. If Spark has been successfully installed, you should see the following output:

   ```
   Welcome to
         ____              __
        / __/__  ___ _____/ /__
       _\ \/ _ \/ _ `/ __/  '_/
      /__ / .__/\_,_/_/ /_/\_\   version 1.5.2
         /_/
   Using Python version 2.7.6 (default, Mar 22 2014 22:59:56)
   SparkContext available as sc, HiveContext available as sqlContext.
   >>>
   ```

▼

6. You should see a similar result by running the `spark-shell` command from any directory.

7. Run the included Pi Estimator example by executing the following command:

```
spark-submit --class org.apache.spark.examples.SparkPi \
--master local \
$SPARK_HOME/lib/spark-examples*.jar 10
```

8. If the installation was successful, you should see something similar to the following result (omitting the informational log messages). Note, this is an estimator program, so the actual result may vary:

```
Pi is roughly 3.140576
```

TRY IT YOURSELF ▼

## Install Spark on Mac OS X

In this example, I install Spark on OS X Mavericks (10.9.5).

Mavericks includes installed versions of Python (2.7.5) and Java (1.8), so I don't need to install them.

1. As shown in Figure 3.1, download the `spark-1.5.2-bin-hadoop2.6.tgz` package from your local mirror into your home directory using curl.

2. Extract the Spark package and create `SPARK_HOME`:

```
tar -xzf spark-1.5.2-bin-hadoop2.6.tgz
sudo mv spark-1.5.2-bin-hadoop2.6 /opt/spark
export SPARK_HOME=/opt/spark
export PATH=$SPARK_HOME/bin:$PATH
```

3. Open the `PySpark` shell by running the `pyspark` command in the Terminal from any directory. If Spark has been successfully installed, you should see the following output:

```
Welcome to
      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 1.5.2
      /_/
Using Python version 2.7.5 (default, Feb 11 2014 07:46:25)
SparkContext available as sc, HiveContext available as sqlContext.
>>>
```

The `SPARK_HOME` environment variable could also be set using the `.profile` file or similar user or system profile scripts.

▼

4. You should see a similar result by running the `spark-shell` command in the terminal from any directory.

5. Run the included Pi Estimator example by executing the following command:

```
spark-submit --class org.apache.spark.examples.SparkPi \
--master local \
$SPARK_HOME/lib/spark-examples*.jar 10
```

6. If the installation was successful, you should see something similar to the following result (omitting the informational log messages). Note, this is an estimator program, so the actual result may vary:

```
Pi is roughly 3.140576
```

▼ TRY IT YOURSELF

## Install Spark on Microsoft Windows

Installing Spark on Windows can be more involved than installing it on Linux or Mac OS X because many of the dependencies (such as Python and Java) need to be addressed first.

This example uses a Windows Server 2012, the server version of Windows 8.

1. You will need a decompression utility capable of extracting `.tar.gz` and `.gz` archives because Windows does not have native support for these archives. 7-zip is a suitable program for this. You can obtain it from **http://7-zip.org/download.html**.

2. As shown in Figure 3.1, download the `spark-1.5.2-bin-hadoop2.6.tgz` package from your local mirror and extract the contents of this archive to a new directory called `C:\Spark`.

3. Install Java using the Oracle JDK Version 1.7, which you can obtain from the Oracle website. In this example, I download and install the `jdk-7u79-windows-x64.exe` package.

4. Disable IPv6 for Java applications by running the following command as an administrator from the Windows command prompt :

```
setx /M _JAVA_OPTIONS "-Djava.net.preferIPv4Stack=true"
```

5. Python is not included with Windows, so you will need to download and install it. You can obtain a Windows installer for Python from **https://www.python.org/getit/**. I use Python 2.7.10 in this example. Install Python into `C:\Python27`.

6. Download the Hadoop common binaries necessary to run Spark compiled for Windows x64 from `hadoop-common-bin`. Extract these files to a new directory called `C:\Hadoop`.

7. Set an environment variable at the machine level for HADOOP_HOME by running the following command as an administrator from the Windows command prompt:

```
setx /M HADOOP_HOME C:\Hadoop
```

8. Update the system path by running the following command as an administrator from the Windows command prompt:

```
setx /M path "%path%;C:\Python27;%PROGRAMFILES%\Java\jdk1.7.0_79\bin;C:\Hadoop"
```

9. Make a temporary directory, C:\tmp\hive, to enable the HiveContext in Spark. Set permission to this file using the winutils.exe program included with the Hadoop common binaries by running the following commands as an administrator from the Windows command prompt:

```
mkdir C:\tmp\hive
C:\Hadoop\bin\winutils.exe chmod 777 /tmp/hive
```

10. Test the Spark interactive shell in Python by running the following command:

```
C:\Spark\bin\pyspark
```

You should see the output shown in Figure 3.2.



**FIGURE 3.2**
The PySpark shell in Windows.

11. You should get a similar result by running the following command to open an interactive Scala shell:

```
C:\Spark\bin\spark-shell
```

12. Run the included Pi Estimator example by executing the following command:

```
C:\Spark\bin\spark-submit --class org.apache.spark.examples.SparkPi --master
local C:\Spark\lib\spark-examples*.jar 10
```

**13.** If the installation was successful, you should see something similar to the following result shown in Figure 3.3. Note, this is an estimator program, so the actual result may vary:



**FIGURE 3.3**
The results of the SparkPi example program in Windows.

## Installing a Multi-node Spark Standalone Cluster

Using the steps outlined in this section for your preferred target platform, you will have installed a single node Spark Standalone cluster. I will discuss Spark's cluster architecture in more detail in **Hour 4, "Understanding the Spark Runtime Architecture."** However, to create a multi-node cluster from a single node system, you would need to do the following:

▶ Ensure all cluster nodes can resolve hostnames of other cluster members and are routable to one another (typically, nodes are on the same private subnet).

▶ Enable passwordless SSH (Secure Shell) for the Spark master to the Spark slaves (this step is only required to enable remote login for the slave daemon startup and shutdown actions).

▶ Configure the `spark-defaults.conf` file on all nodes with the URL of the Spark master node.

▶ Configure the `spark-env.sh` file on all nodes with the hostname or IP address of the Spark master node.

▶ Run the `start-master.sh` script from the `sbin` directory on the Spark master node.

▶ Run the `start-slave.sh` script from the `sbin` directory on all of the Spark slave nodes.

▶ Check the Spark master UI. You should see each slave node in the `Workers` section.

▶ Run a test Spark job.

TRY IT YOURSELF ▼

## Configuring and Testing a Multinode Spark Cluster

Take your single node Spark system and create a basic two-node Spark cluster with a master node and a worker node.

In this example, I use two Linux instances with Spark installed in the same relative paths: one with a hostname of `sparkmaster`, and the other with a hostname of `sparkslave`.

1. Ensure that each node can resolve the other. The `ping` command can be used for this. For example, from `sparkmaster`:

```
ping sparkslave
```

2. Ensure the firewall rules of network ACLs will allow traffic on multiple ports between cluster instances because cluster nodes will communicate using various TCP ports (normally not a concern if all cluster nodes are on the same subnet).

3. Create and configure the `spark-defaults.conf` file on all nodes. Run the following commands on the `sparkmaster` and `sparkslave` hosts:

```
cd $SPARK_HOME/conf
sudo cp spark-defaults.conf.template spark-defaults.conf
sudo sed -i "\$aspark.master\tspark://sparkmaster:7077" spark-defaults.conf
```

4. Create and configure the `spark-env.sh` file on all nodes. Complete the following tasks on the `sparkmaster` and `sparkslave` hosts:

```
cd $SPARK_HOME/conf
sudo cp spark-env.sh.template spark-env.sh
sudo sed -i "\$aSPARK_MASTER_IP=sparkmaster" spark-env.sh
```

5. On the `sparkmaster` host, run the following command:

```
sudo $SPARK_HOME/sbin/start-master.sh
```

6. On the `sparkslave` host, run the following command:

```
sudo $SPARK_HOME/sbin/start-slave.sh spark://sparkmaster:7077
```

7. Check the Spark master web user interface (UI) at http://sparkmaster:8080/.

8. Check the Spark worker web UI at http://sparkslave:8081/.

9. Run the built-in Pi Estimator example from the terminal of either node:

```
spark-submit --class org.apache.spark.examples.SparkPi \
--master spark://sparkmaster:7077 \
--driver-memory 512m \
--executor-memory 512m \
--executor-cores 1 \
$SPARK_HOME/lib/spark-examples*.jar 10
```

▼    **10.** If the application completes successfully, you should see something like the following (omitting informational log messages). Note, this is an estimator program, so the actual result may vary:

```
Pi is roughly 3.140576
```

This is a simple example. If it was a production cluster, I would set up passwordless SSH to enable the `start-all.sh` and `stop-all.sh` shell scripts. I would also consider modifying additional configuration parameters for optimization.

CAUTION

### Spark Master Is a Single Point of Failure in Standalone Mode

Without implementing *High Availability (HA),* the Spark Master node is a *single point of failure (SPOF)* for the Spark cluster. This means that if the Spark Master node goes down, the Spark cluster would stop functioning, all currently submitted or running applications would fail, and no new applications could be submitted.

High Availability can be configured using *Apache Zookeeper,* a highly reliable distributed coordination service. You can also configure HA using the filesystem instead of Zookeeper; however, this is not recommended for production systems.

# Exploring the Spark Install

Now that you have Spark up and running, let's take a closer look at the install and its various components.

If you followed the instructions in the previous section, "Installing Spark in Standalone Mode," you should be able to browse the contents of $SPARK_HOME.

In Table 3.1, I describe each subdirectory of the Spark installation.

**TABLE 3.1  Spark Installation Subdirectories**

| Directory | Description |
|---|---|
| bin | Contains all of the commands/scripts to run Spark applications interactively through shell programs such as `pyspark`, `spark-shell`, `spark-sql` and `sparkR`, or in batch mode using `spark-submit`. |
| conf | Contains templates for Spark configuration files, which can be used to set Spark environment variables (`spark-env.sh`) or set default master, slave, or client configuration parameters (`spark-defaults.conf`). There are also configuration templates to control logging (`log4j.properties`), metrics collection (`metrics.properties`), as well as a template for the `slaves` file, which controls which slave nodes can join the Spark cluster. |

| Directory | Description |
|-----------|-------------|
| ec2 | Contains scripts to deploy Spark nodes and clusters on Amazon Web Services (AWS) Elastic Compute Cloud (EC2). I will cover deploying Spark in EC2 in **Hour 5, "Deploying Spark in the Cloud."** |
| lib | Contains the main assemblies for Spark including the main library (`spark-assembly-x.x.x-hadoopx.x.x.jar`) and included example programs (`spark-examples-x.x.x-hadoopx.x.x.jar`), of which we have already run one, SparkPi, to verify the installation in the previous section. |
| licenses | Includes license files covering other included projects such as Scala and JQuery. These files are for legal compliance purposes only and are not required to run Spark. |
| python | Contains all of the Python libraries required to run PySpark. You will generally not need to access these files directly. |
| sbin | Contains administrative scripts to start and stop master and slave services (locally or remotely) as well as start processes related to YARN and Mesos. I used the `start-master.sh` and `start-slave.sh` scripts when I covered how to install a multi-node cluster in the previous section. |
| data | Contains sample data sets used for testing mllib (which we will discuss in more detail in **Hour 16, "Machine Learning with Spark"**). |
| examples | Contains the source code for all of the examples included in `lib/spark-examples-x.x.x-hadoopx.x.x.jar`. Example programs are included in Java, Python, R, and Scala. You can also find the latest code for the included examples at **https://github.com/apache/spark/tree/master/examples**. |
| R | Contains the `SparkR` package and associated libraries and documentation. I will discuss SparkR in **Hour 15, "Getting Started with Spark and R"** |

# Deploying Spark on Hadoop

As discussed previously, deploying Spark with Hadoop is a popular option for many users because Spark can read from and write to the data in Hadoop (in HDFS) and can leverage Hadoop's process scheduling subsystem, YARN.

## Using a Management Console or Interface

If you are using a commercial distribution of Hadoop such as Cloudera or Hortonworks, you can often deploy Spark using the management console provided with each respective platform: for example, Cloudera Manager for Cloudera or Ambari for Hortonworks.

If you are using the management facilities of a commercial distribution, the version of Spark deployed may lag the latest stable Apache release because Hadoop vendors typically update their software stacks with their respective major and minor release schedules.

## Installing Manually

Installing Spark on a YARN cluster manually (that is, not using a management interface such as Cloudera Manager or Ambari) is quite straightforward to do.

▼ TRY IT YOURSELF

### Installing Spark on Hadoop Manually

1. Follow the steps outlined for your target platform (for example, Red Hat Linux, Windows, and so on) in the earlier section "Installing Spark in Standalone Mode."

2. Ensure that the system you are installing on is a Hadoop client with configuration files pointing to a Hadoop cluster. You can do this as shown:

   ```
   hadoop fs -ls
   ```

   This lists the contents of your user directory in HDFS. You could instead use the path in HDFS where your input data resides, such as

   ```
   hadoop fs -ls /path/to/my/data
   ```

   If you see an error such as `hadoop: command not found`, you need to make sure a correctly configured Hadoop client is installed on the system before continuing.

3. Set either the `HADOOP_CONF_DIR` or `YARN_CONF_DIR` environment variable as shown:

   ```
   export HADOOP_CONF_DIR=/etc/hadoop/conf
   # or
   export YARN_CONF_DIR=/etc/hadoop/conf
   ```

   As with `SPARK_HOME`, these variables could be set using the `.bashrc` or similar profile script sourced automatically.

4. Execute the following command to test Spark on YARN:

   ```
   spark-submit --class org.apache.spark.examples.SparkPi \
   --master yarn-cluster \
   $SPARK_HOME/lib/spark-examples*.jar 10
   ```

**5.** If you have access to the YARN Resource Manager UI, you can see the Spark job running in YARN as shown in Figure 3.4:



**FIGURE 3.4**
The YARN ResourceManager UI showing the Spark application running.

**6.** Clicking the **ApplicationsMaster** link in the **ResourceManager UI** will redirect you to the Spark UI for the application:



**FIGURE 3.5**
The Spark UI.

Submitting Spark applications using YARN can be done in two submission modes: `yarn-cluster` or `yarn-client`.

Using the `yarn-cluster` option, the Spark Driver and Spark Context, ApplicationsMaster, and all executors run on YARN NodeManagers. These are all concepts we will explore in detail in **Hour 4, "Understanding the Spark Runtime Architecture."** The `yarn-cluster` submission mode is intended for production or non interactive/batch Spark applications. You cannot use

`yarn-cluster` as an option for any of the interactive Spark shells. For instance, running the following command:

```
spark-shell --master yarn-cluster
```

will result in this error:

```
Error: Cluster deploy mode is not applicable to Spark shells.
```

Using the `yarn-client` option, the Spark Driver runs on the client (the host where you ran the Spark application). All of the tasks and the ApplicationsMaster run on the YARN NodeManagers however unlike `yarn-cluster` mode, the Driver does not run on the ApplicationsMaster. The `yarn-client` submission mode is intended to run interactive applications such as `pyspark` or `spark-shell`.

---

CAUTION

## Running Incompatible Workloads Alongside Spark May Cause Issues

Spark is a memory-intensive processing engine. Using Spark on YARN will allocate containers, associated CPU, and memory resources to applications such as Spark as required. If you have other memory-intensive workloads, such as Impala, Presto, or HAWQ running on the cluster, you need to ensure that these workloads can coexist with Spark and that neither compromises the other. Generally, this can be accomplished through application, YARN cluster, scheduler, or application queue configuration and, in extreme cases, operating system cgroups (on Linux, for instance).

---

# Summary

In this hour, I have covered the different deployment modes for Spark: Spark Standalone, Spark on Mesos, and Spark on YARN.

Spark Standalone refers to the built-in process scheduler it uses as opposed to using a preexisting external scheduler such as Mesos or YARN. A Spark Standalone cluster could have any number of nodes, so the term "Standalone" could be a misnomer if taken out of context. I have showed you how to install Spark both in Standalone mode (as a single node or multi-node cluster) and how to install Spark on an existing YARN (Hadoop) cluster.

I have also explored the components included with Spark, many of which you will have used by the end of this book.

You're now up and running with Spark. You can use your Spark installation for most of the exercises throughout this book.

# Q&A

**Q.** What are the factors involved in selecting a specific deployment mode for Spark?

**A.** The choice of deployment mode for Spark is primarily dependent upon the environment you are running in and the availability of external scheduling frameworks such as YARN or Mesos. For instance, if you are using Spark with Hadoop and you have an existing YARN infrastructure, Spark on YARN is a logical deployment choice. However, if you are running Spark independent of Hadoop (for instance sourcing data from S3 or a local filesystem), Spark Standalone may be a better deployment method.

**Q.** What is the difference between the `yarn-client` and the `yarn-cluster` options of the `--master` argument using `spark-submit`?

**A.** Both the `yarn-client` and `yarn-cluster` options execute the program in the Hadoop cluster using YARN as the scheduler; however, the `yarn-client` option uses the client host as the driver for the program and is designed for testing as well as interactive shell usage.

# Workshop

The workshop contains quiz questions and exercises to help you solidify your understanding of the material covered. Try to answer all questions before looking at the "Answers" section that follows.

## Quiz

1. **True or false:** A Spark Standalone cluster consists of a single node.

2. Which component is not a prerequisite for installing Spark?

    **A.** Scala

    **B.** Python

    **C.** Java

3. Which of the following subdirectories contained in the Spark installation contains scripts to start and stop master and slave node Spark services?

    **A.** `bin`

    **B.** `sbin`

    **C.** `lib`

4. Which of the following environment variables are required to run Spark on Hadoop/YARN?

    **A.** `HADOOP_CONF_DIR`

    **B.** `YARN_CONF_DIR`

    **C.** Either `HADOOP_CONF_DIR` or `YARN_CONF_DIR` will work.

## Answers

1. **False.** Standalone refers to the independent process scheduler for Spark, which could be deployed on a cluster of one-to-many nodes.

2. **A.** The Scala assembly is included with Spark; however, Java and Python must exist on the system prior to installation.

3. **B.** `sbin` contains administrative scripts to start and stop Spark services.

4. **C.** Either the `HADOOP_CONF_DIR` or `YARN_CONF_DIR` environment variable must be set for Spark to use YARN.

## Exercises

1. Using your Spark Standalone installation, execute `pyspark` to open a PySpark interactive shell.

2. Open a browser and navigate to the SparkUI at http://localhost:4040.

3. Click the Environment top menu link or navigate to Environment page directly using the url: http://localhost:4040/environment/.

4. Note some of the various environment settings and configuration parameters set. I will explain many of these in greater detail throughout the book.

# Index

## N

## U