

Notes for this computer

Windows path for R has been set for R version 3.2.1 as “C:\Program Files\R\R-3.2.1\bin\x64”

Spark is installed in “C:\Apache\spark-1.4.1-bin-hadoop2.6” (NOTE- use “/” instead of “\” when writing the path inside R.

Installing and Starting SparkR Locally on Windows OS and RStudio

July 26, 2015

By [emaasit](#)

(This article was first published on [Emaasit's Blog » R](#), and kindly contributed to [R-bloggers](#))

Introduction

With the recent release of Apache Spark 1.4.1 on July 15th, 2015, I wanted to write a step-by-step guide to help new users get up and running with SparkR locally on a Windows machine using command shell and RStudio. SparkR provides an R frontend to Apache Spark and using Spark’s distributed computation engine allows R-Users to run large scale data analysis from the R shell. The steps listed here are also documented in my online book title “[Getting Started with SparkR for Big Data Analysis](#)” which can be accessed at: <http://www.danielemaasit.com/getting-started-with-sparkr/>. These steps will get you up and running in less than 5 mins.

Prerequisite

Make sure you have Java 6+ installed on your computer and the system environments set.

Step 1: Download Spark

Open your web browser and open this web page: <http://spark.apache.org/>. This is the official website for the Apache Spark project. You should see a large green button to the right of the page that reads “Download Spark”, as shown in Figure 1. Click the green button.



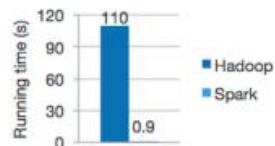
[Download](#) [Libraries](#) [Documentation](#) [Examples](#) [Community](#) [FAQ](#)

Apache Spark™ is a fast and general engine for large-scale data processing.

Speed

Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

Spark has an advanced DAG execution engine that supports cyclic data flow and in-memory computing.



Logistic regression in Hadoop and Spark

Ease of Use

Write applications quickly in Java, Scala, Python, R.

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python and R shells.

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

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

Word count in Spark's Python API

Latest News

[Spark 1.4.1 released](#) (Jul 15, 2015)
[Spark Summit 2015 Videos Posted](#) (Jun 29, 2015)
[Spark 1.4.0 released](#) (Jun 11, 2015)
[One month to Spark Summit 2015 in San Francisco](#) (May 15, 2015)

[Archive](#)

[Download Spark](#)

Built-in Libraries:

[Spark SQL](#)
[Spark Streaming](#)
[MLlib \(machine learning\)](#)
[GraphX \(graph\)](#)
[Third-Party Packages](#)

Figure 1: Apache Spark home page

Clicking the green button will take you to the download page as shown in Figure 2 below.

Download Spark

The latest release of Spark is Spark 1.4.1, released on July 15, 2015 ([release notes](#)) ([git tag](#))

1. Choose a Spark release:
2. Choose a package type:
3. Choose a download type:
4. Download Spark: [spark-1.4.1.tgz](#)
5. Verify this release using the [1.4.1 signatures and checksums](#).

Note: Scala 2.11 users should download the Spark source package and build with Scala 2.11 support.

Link with Spark

Spark artifacts are [hosted in Maven Central](#). You can add a Maven dependency with the following coordinates:

```
groupId: org.apache.spark
artifactId: spark-core_2.10
version: 1.4.1
```

Spark Source Code Management

If you are interested in working with the newest under-development code or contributing to Spark development, you can also check out the master branch from Git:

```
# Master development branch
```

Latest News

- [Spark 1.4.1 released](#) (Jul 15, 2015)
- [Spark Summit 2015 Videos Posted](#) (Jun 29, 2015)
- [Spark 1.4.0 released](#) (Jun 11, 2015)
- [One month to Spark Summit 2015 in San Francisco](#) (May 15, 2015)

[Archive](#)

[Download Spark](#)

Built-in Libraries:

- [Spark SQL](#)
- [Spark Streaming](#)
- [MLlib \(machine learning\)](#)
- [GraphX \(graph\)](#)

Third-Party Packages

Figure 2: The download page

You should follow the steps 1 to 3 to create a download link for a Spark Package of your choice. On the “2. Choose a package type” option, select any pre-built package type from the drop-down list (Figure 3). Since we want to experiment locally on windows, a pre-built package for Hadoop 2.6 and later will suffice.

Download Spark

The latest release of Spark is Spark 1.4.1, released on July 15, 2015 ([release notes](#)) ([git tag](#))

1. Choose a Spark release:
2. Choose a package type:
3. Choose a download type:
4. Download Spark: [spark-1.4.1.tgz](#)
5. Verify this release using the [1.4.1 signatures and checksums](#).

Note: Scala 2.11 users should c

- Source Code [can build several Hadoop versions]
- Source Code [can build several Hadoop versions]
- Pre-built with user-provided Hadoop [can use with most Hadoop distributions]
- Pre-built for Hadoop 2.6 and later**
- Pre-built for Hadoop 2.4 and later
- Pre-built for Hadoop 2.3
- Pre-built for Hadoop 1.X
- Pre-built for CDH 4

Figure 3: Choose a package type

On the “3. Choose a download type” option, select “Direct Download” from the drop-down list (Figure 4).

Download Spark

The latest release of Spark is Spark 1.4.1, released on July 15, 2015 ([release notes](#)) ([git tag](#))

1. Choose a Spark release: 1.4.1 (Jul 15 2015) ▾
2. Choose a package type: Pre-built for Hadoop 2.6 and later ▾
3. Choose a download type: Select Apache Mirror ▾
Select Apache Mirror
Direct Download
4. Download Spark: [spark-1.4.1-bin-hadoop2.tgz](#)
5. Verify this release using the [1.4.1 signatures and checksums](#).

Note: Scala 2.11 users should download the Spark source package and build with Scala 2.11 support.

Figure 4: Choose a download type

After selecting the download type, a link is created next to the option “4. Download Spark” (Figure 5). Click this link to download Spark.

Download Spark

The latest release of Spark is Spark 1.4.1, released on July 15, 2015 ([release notes](#)) ([git tag](#))

1. Choose a Spark release: 1.4.1 (Jul 15 2015) ▾
2. Choose a package type: Pre-built for Hadoop 2.6 and later ▾
3. Choose a download type: Direct Download ▾
4. Download Spark: [spark-1.4.1-bin-hadoop2.tgz](#)
5. Verify this release using the [1.4.1 signatures and checksums](#).

Note: Scala 2.11 users should download the Spark source package and build with Scala 2.11 support.

Figure 5: Click the download link

Save the zipped file to your computer (Figure 6).

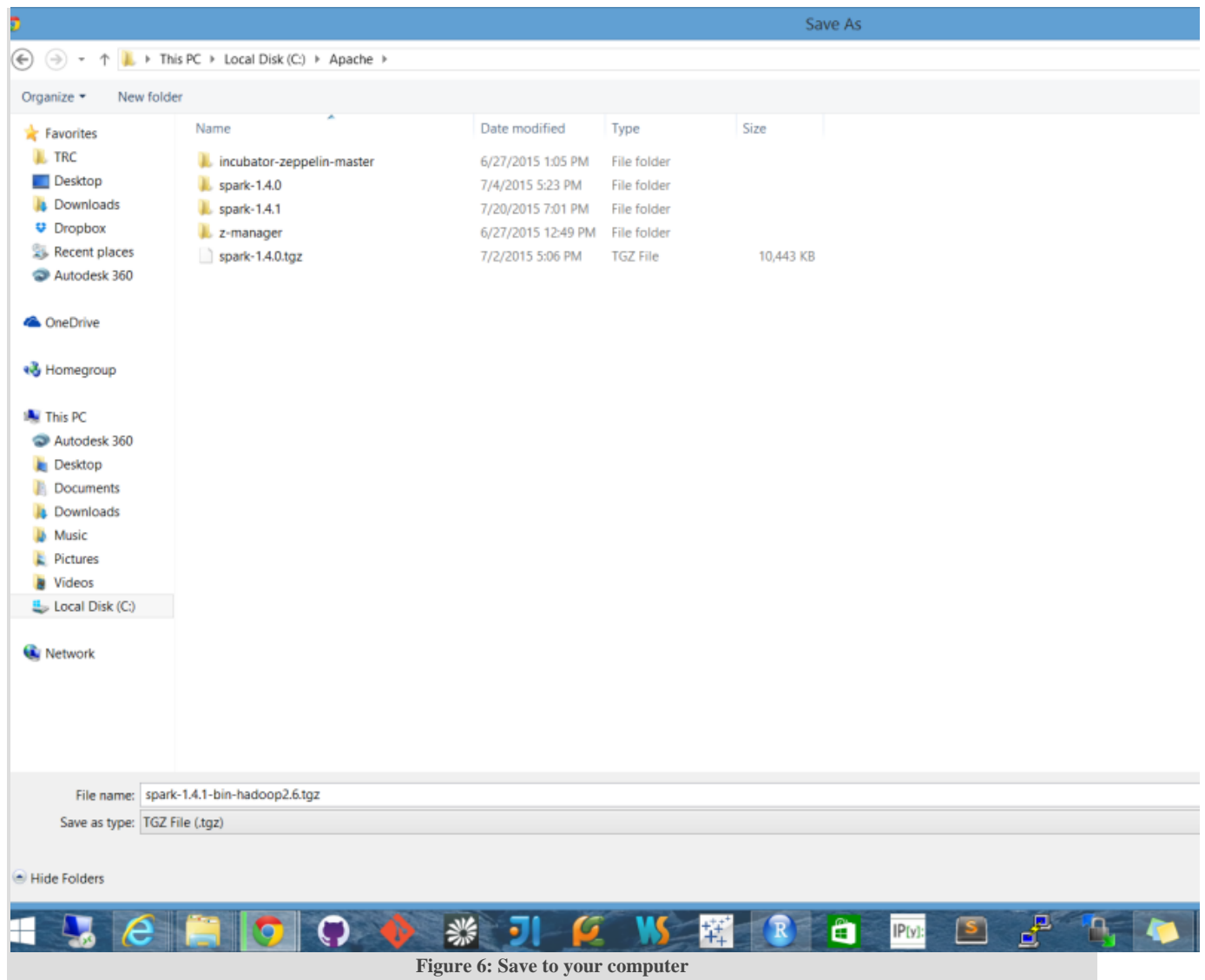


Figure 6: Save to your computer

Step 2: Unzip Built Package

Unzip and save the files to a directory folder of your choice. In Figure 7 below, I chose to save to *"C:/Apache/Spark-1.4.1"*.

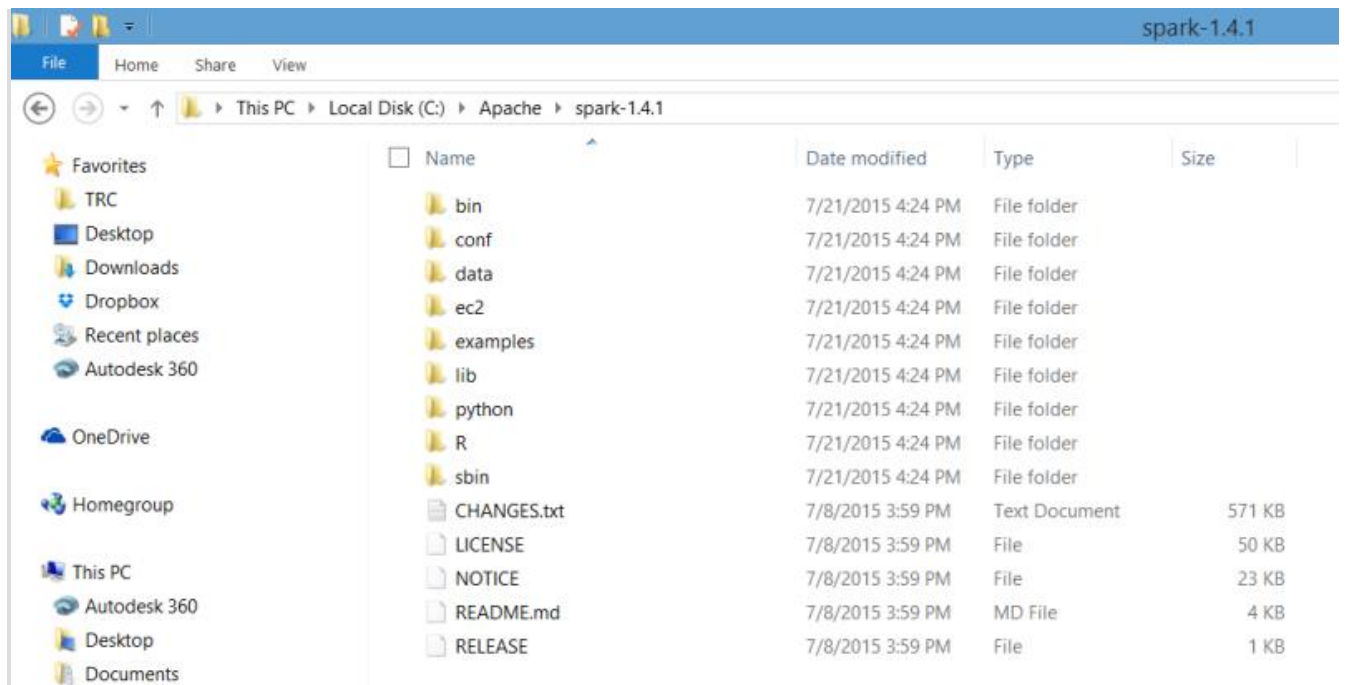


Figure 7: List of Files in unzipped folder

Step 3: Run in Command Prompt

Now start your favorite command shell and change directory to your Spark folder as shown in Figure 8.

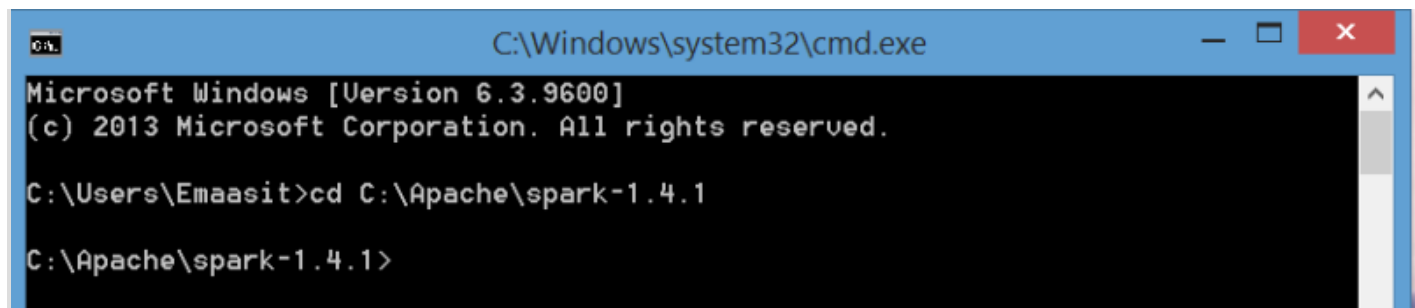


Figure 8: Start command prompt and change directory

To start SparkR, simply run the command ".binsparkR" on the top-level Spark directory as shown in Figure 9 below.

```
C:\Apache\spark-1.4.1>.\bin\sparkR
```

Figure 9: Start SparkR

You will see logs on your screen that should take at most 15 seconds to launch SparkR. If everything ran smoothly you should see a welcome message that reads "Welcome to SparkR!" as shown in Figure 10.

```
15/07/21 17:01:20 INFO SparkUI: Started SparkUI at http://192.168.56.1:4040
15/07/21 17:01:20 INFO Executor: Starting executor ID driver on host localhost
15/07/21 17:01:20 INFO Utils: Successfully started service 'org.apache.spark.net
work.netty.NettyBlockTransferService' on port 62447.
15/07/21 17:01:20 INFO NettyBlockTransferService: Server created on 62447
15/07/21 17:01:20 INFO BlockManagerMaster: Trying to register BlockManager
15/07/21 17:01:20 INFO BlockManagerMasterEndpoint: Registering block manager loc
alhost:62447 with 265.4 MB RAM, BlockManagerId(driver, localhost, 62447)
15/07/21 17:01:20 INFO BlockManagerMaster: Registered BlockManager

Welcome to SparkR!
Spark context is available as sc, SQL context is available as sqlContext
>
>
```

Figure 10: Welcome to SparkR!

At this point you are ready to start prototyping with SparkR on the command shell. Note that a Spark context and a SQL Context have been initialized for you as “*sc*” and “*sqlContext*” respectively. You can now start experimenting using the example shown in **Step 4.5**.

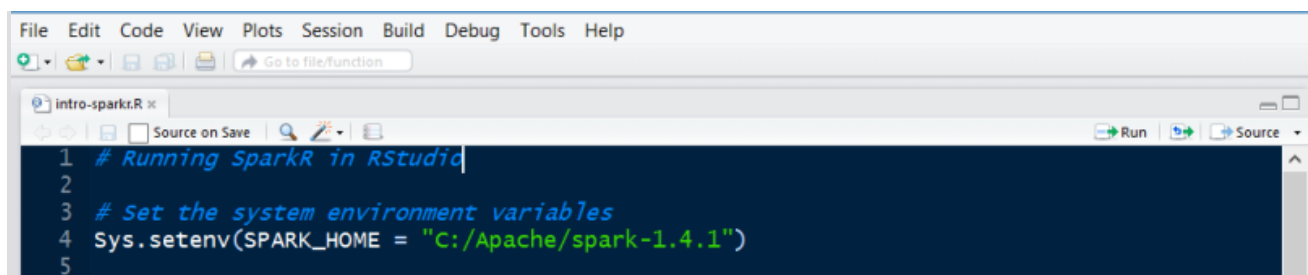
Running in RStudio

While using SparkR in the command shell is good for quickly getting started, most R users typically use an Integrated Development Environment (IDE) like RStudio for development and running production ready code. Step 4 below will guide you to get started using SparkR in RStudio.

Step 4: Run in RStudio

- **Step 4.1: Set System Environment**

Once you have opened RStudio, you need to set the system environment first. You have to point your R session to the installed version of SparkR. Use the code shown in Figure 11 below but replace the *SPARK_HOME* variable using the path to your Spark folder. Mine is “C:/Apache/Spark-1.4.1”.

The image shows the RStudio interface with a menu bar (File, Edit, Code, View, Plots, Session, Build, Debug, Tools, Help) and a toolbar. A script editor window titled 'intro-spark.R' is open, displaying the following R code:

```
1 # Running SparkR in RStudio
2
3 # Set the system environment variables
4 Sys.setenv(SPARK_HOME = "C:/Apache/spark-1.4.1")
5
```

Figure 11: Set System Environment Variable

- **Step 4.2: Set the Library Paths**

Second, you have to set the library path for Spark as shown in Figure 12 below.


```

5
6
7 # Set the library path
8 .libPaths(c(file.path(Sys.getenv("SPARK_HOME"), "R", "lib"), .libPaths()))
9
10

```

Figure 12: Set the Library Path

- **Step 4.3: Load SparkR Library**

Next, you can now load SparkR just as you would any other R library using the `library()` command as shown in Figure 13.

```

13
14
15 library(sparkR)
16 #load the sparkR library
17
18

```

Figure 13: Load the SparkR library

- **Step 4.4: Initialize Spark Context and SQL Context**

Initialize SparkR by creating a Spark context using the command `sparkR.init()`. The argument in this command is `master = "local[N]"`, where N stands for the number of threads that you want to use.

Also, you need to create a SQL context to be able to work with DataFrames (the main abstraction in SparkR). Use the command `sparkRSQL.init()` to create a SQL context from your Spark context as shown in Figure 14.

```

14
15 # Create a spark context and a SQL context
16 sc <- sparkR.init(master = "local")
17 sqlContext <- sparkRSQL.init(sc)
18
19

```

Figure 14: Initialize Spark Context and SQL Context

When you run the above commands (From step 4.1 to 4.4), this invokes the “spark-submit” script that launches java, as shown in Figure 15. If this runs successfully, your Spark context and SQL context should be created and at this stage you should be able to start experimenting with SparkR.


```

> #load the sparkr library
> library(SparkR)

Attaching package: 'SparkR'

The following object is masked from 'package:stats':

  filter

The following objects are masked from 'package:base':

  intersect, sample, table

> # Create a spark context and a SQL context
> sc <- sparkR.init(master = "local")
Launching java with spark-submit command C:/Apache/spark-1.4.1/bin/spark-submit.cmd
sparkr-shell C:\Users\Emaasit\AppData\Local\Temp\RtmpQdRVBW\backend_port20886132bb9
> sqlContext <- sparkRSQL.init(sc)
> |

```

Figure 15: All set

- **Step 4.5: A Quick Example**

You can start experimenting with SparkR on the command shell and in RStudio using the example provided below. You can monitor your Spark jobs using the Spark UI at localhost:4040

```

# Set the system environment variables
Sys.setenv(SPARK_HOME = "C:/Apache/spark-1.4.1")
.libPaths(c(file.path(Sys.getenv("SPARK_HOME"), "R", "lib"), .libPaths()))
#load the Sparkr library
library(SparkR)
# Create a spark context and a SQL context
sc <- sparkR.init(master = "local")
sqlContext <- sparkRSQL.init(sc)
#create a sparkR DataFrame
DF <- createDataFrame(sqlContext, faithful)
head(DF)

# Create a simple local data.frame
localDF <- data.frame(name=c("John", "Smith", "Sarah"), age=c(19, 23, 18))
# Convert local data frame to a SparkR DataFrame
df <- createDataFrame(sqlContext, localDF)
# Print its schema
printSchema(df)
# root
# |-- name: string (nullable = true)
# |-- age: double (nullable = true)
# Create a DataFrame from a JSON file
path <- file.path(Sys.getenv("SPARK_HOME"), "examples/src/main/resources/people.json")
peopleDF <- jsonFile(sqlContext, path)
printSchema(peopleDF)
# Register this DataFrame as a table.
registerTempTable(peopleDF, "people")
# SQL statements can be run by using the sql methods provided by sqlContext
teenagers <- sql(sqlContext, "SELECT name FROM people WHERE age >= 13 AND age <= 19")

```

```
# Call collect to get a local data.frame
teenagersLocalDF <- collect(teenagers)
# Print the teenagers in our dataset
print(teenagersLocalDF)
# Stop the SparkContext now
sparkR.stop()
```



view `rawintro-sparkr.R` hosted with

by **GitHub**

Final Remarks

The purpose of this blog post was to get you up and running quickly with SparkR locally on a personal computer. In the next blog post, I will show you how to use SparkR on a cloud computing framework like [Amazon Elastic Compute Cloud](#) (EC2) to manipulate large datasets with millions of records.

Tagged: [Apache Spark](#), [Big Data](#), [R](#), [RStudio](#), [SparkR](#) [Comments: 8](#)