

ILLINOIS DATA SCIENCE INITIATIVE

TECHNICAL REPORTS

---

## Setting up Spark for Python

---

*Author:*

Sameet Sapra & Joshua Chang & Professor Brunner

February 25, 2017

# Setting up Spark for Python

**SAMEET SAPRA<sup>1</sup>, JOSHUA CHANG<sup>2</sup>, AND PROFESSOR BRUNNER<sup>3</sup>**

<sup>1</sup>National Center For Supercomputing Applications (NCSA)

<sup>2</sup>Laboratory for Computation, Data, and Machine Learning

<sup>3</sup>Illinois Data Science Initiative

Compiled February 25, 2017

## An introduction to setting up Apache Spark with Python

<https://github.com/lcdm-uiuc>

### 1. WHAT IS APACHE SPARK?

Apache Spark is a fast cluster computing system. Spark can be configured with multiple cluster managers like YARN, Mesos, Amazon EC2, or standalone mode. Along with that it can be configured in local mode and standalone mode. This technical report will install Spark on top of YARN.

Spark supports many high level tools like GraphX and MLlib, for graphs processing and machine learning. It also has many APIs in Java, Scala, and Python, and we will go over PySpark, Python's API for Spark.

### 2. WHY PYSPARK?

PySpark is an API that interfaces with RDD's in Python. It is built on top of Spark's Java API and exposes the Spark programming model to Python. PySpark makes use of a library called Py4J, which enables Python programs to dynamically access Java objects in a Java Virtual Machine. This allows data to be processed in Python and cached in the JVM.

### 3. PYSPARK PREREQUISITES

Assume that the environment is a cluster with Hadoop installed. Verify that Java, Scala, and Yarn are installed by checking the versions. If they are installed, skip to the Spark configuration steps, otherwise follow the installation instructions given below.

### 4. INSTALLATION

Install Java on the system with:

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

Install Spark on the system by downloading the rpm. Here we will install spark-1.6:

```
> wget http://apache.mirrors.ionfish.org/spark
  /spark-1.6.0/spark-1.6.0-bin-hadoop2.6.tgz
> tar -zxvf spark-1.6.0-bin-hadoop2.6.tgz
> mv spark-1.6.0-bin-hadoop2.6 spark
```

To install yarn:

```
> sudo wget https://dl.yarnpkg.com/rpm/yarn.
  repo -O /etc/yum.repos.d/yarn.repo
> sudo yum install yarn
```

Let's configure Spark to only show errors on startup, rather than all info. This will make the Spark shell less cluttered when it is opened on startup.

```
> cd /usr/hdp/2.3.6.0-3796/spark/conf
> vi log4j.properties
```

Find a line that describe the rootCategory of log:

```
log4j.rootCategory=INFO, console
```

Replace INFO with ERROR:

```
log4j.rootCategory=ERROR, console
```

If all that worked, you should see the Spark shell start up when you type

```
> spark-shell
```

### 5. SETTING UP PYSPARK

Now let's ensure that python 2.7 is installed and configured.

```
> yum install -y centos-release-SCL
> yum install -y python27
> yum -y install python-pip
```

Finally, let's set up the environment variables for Spark and Python.

```
> vi $HOME/.bashrc
> export SPARK_HOME=/usr/hdp/2.3.6.0-3796/spark
> export PYTHONPATH=$SPARK_HOME/python
> export SPARK_HIVE=true
```

Then, reload the environment:

```
> source $HOME/.bashrc
```

Again, let's make sure Python is setup correctly. Type

```
> python
```

You should see the Python REPL come up. Finally, we can start writing code in Python.

## 6. WRITING OUR FIRST LINES OF PYSPARK

This example is taken from the Apache Spark website. It runs a parallelized operation to compute the value of  $\pi$  and is a good example to see the benefits of Spark.

```
def inside(p):  
    x, y = random.random(), random.random()  
    return x*x + y*y < 1  
  
count = sc.parallelize(xrange(0, NUM_SAMPLES)) \  
    .filter(inside).count()  
print "Pi is roughly %f" % (4.0 * count / NUM_SAMPLES)
```

Once you've written your Python code, you can compile and deploy it with:

```
> spark-submit --master yarn-cluster MY_PYTHON_FILE.py
```

The `'-num-executors 10'` flag (arbitrary number) may be added to specify how many executors, the objects responsible for executing tasks, are to be used. Using as many executors as data nodes is recommended to decrease minimize runtime.

## 7. CONCLUSION

This technical report serves as a guide to set up an environment to run Spark on HDFS and write some simple Python code to take advantage of Spark using PySpark.