

Instalación de PySpark en Windows

PySpark es una biblioteca Spark escrita en Python para ejecutar la aplicación Python usando las capacidades de Apache Spark. por lo que no hay una biblioteca de PySpark para descargar.

Sigue los pasos a continuación para instalar PySpark en Windows.

[PASO 1. Instalar la distribución Python o Anaconda](#)

[PASO 2. Instalación de PySpark en Windows](#)

[PASO 3. Instale winutils.exe en Windows](#)

[PASO 4. Ejecuta PYSPARK](#)

[Paso 5. PySpark con Jupyter notebook](#)

[Anexo: Si no tengo PIP instalado](#)

[Probar jupyter notebook con spark](#)

PASO 1. Instalar la distribución Python o Anaconda

El primer paso es descargar e instalar Python desde [Python.org](https://python.org) o la distribución Anaconda que incluye Python, Spyder IDE y Jupyter notebook.

Link anaconda: <https://docs.anaconda.com/anaconda/install/windows/>

Link Python: <https://www.python.org/downloads/>

PASO 1.2 Instalar JDK y definir variables de entorno

También debemos asegurarnos de que tenemos java instalado y debemos asegurarnos de que las variables de entorno están correctamente establecidas.

Link descarga:

<https://www.oracle.com/es/java/technologies/javase/javase8-archive-downloads.html>

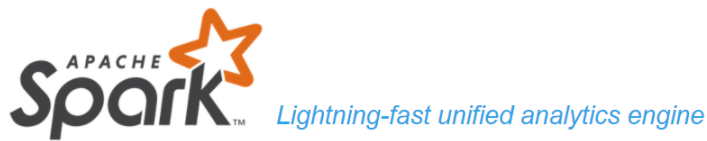
Establecemos las variables de entorno:

```
JAVA_HOME = C:\Program Files\Java\jdk1.8.0_201
```

```
PATH = %PATH%;C:\Program Files\Java\jdk1.8.0_201\bin
```

PASO 2. Instalación de PySpark en Windows

En la página de descarga de Spark(<https://spark.apache.org/downloads.html>) , seleccione el enlace "Descargar Spark (punto 3)" para descargar. Si desea utilizar una versión diferente de Spark & Hadoop, seleccione la que deseaba de los menús desplegados y el enlace en el punto 3 cambia a la versión seleccionada y le proporciona un enlace actualizado para descargar.



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

Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Download Spark: [spark-3.1.1-bin-hadoop2.7.tgz](#)
4. Verify this release using the 3.1.1 [signatures](#), [checksums](#) and [project release KEYS](#).

Note that, Spark 2.x is pre-built with Scala 2.11 except version 2.4.2, which is pre-built with Scala 2.12. Spark 3.0+ is pre-built with Scala 2.12.



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

Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Download Spark: [spark-3.1.1-bin-hadoop2.7.tgz](#)
4. Verify this release using the 3.1.1 [signatures](#), [checksums](#) and [project release KEYS](#).

Note that, Spark 2.x is pre-built with Scala 2.11 except version 2.4.2, which is pre-built with Scala 2.12. Spark 3.0+ is pre-built with Scala 2.12.

Latest Preview Release

Una vez lo hemos descargado, descomprime el binario usando 7zip o winrar y copia la carpeta de dentro spark-3.0.0-bin-hadoop2.7 en la ruta que quieras por ejemplo **c:\apps**

Ahora crearemos las variables de entorno.

```
SPARK_HOME = C:\apps\spark-3.0.0-bin-hadoop2.7
```

```
HADOOP_HOME = C:\apps\spark-3.0.0-bin-hadoop2.7
```

```
PATH=%PATH%;C:\apps\spark-3.0.0-bin-hadoop2.7\bin
```

```
conda install -c conda-forge Findspark
```

O desde una consola cmd ejecutar:

pip install findspark

Podemos trabajar desde aquí o instalar jupyter

Ejecutar para arrancar

jupyter notebook

<http://localhost:8888/>

Anexo: Si no tengo PIP instalado

Las siguientes instrucciones deberían funcionar en Windows 7, Windows 8.1 y Windows 10:

1. Descarga el script del [instalador get-pip.py](#). Si estás en Python 3.2, necesitarás esta versión de [get-pip.py](#). En caso de tener Python 3.3 o 3.4 usar estas versiones de PiP correspondientemente [Python 3.3 get-pip.py](#) o [Python 3.4 get-pip.py](#). De cualquier manera, haga clic derecho en el enlace y seleccione Guardar como y guárdelo en cualquier carpeta del pc, como su carpeta de Descargas.
2. Abra el símbolo del sistema y navegue hasta el archivo get-pip.py.
3. Ejecute el siguiente comando: `python get-pip.py`

Probar jupyter notebook con spark

Probar por ejemplo esto:

```
import findspark
findspark.init()
findspark.find()
import pyspark
```

```
from pyspark import SparkContext, SparkConf
from pyspark.sql import SparkSession
sc = SparkContext.getOrCreate()
```

```
nums = sc.parallelize([1,2,3,4])
```

nums

```
In [3]: import findspark
findspark.init()
findspark.find()
import pyspark

from pyspark import SparkContext, SparkConf
from pyspark.sql import SparkSession
sc = SparkContext.getOrCreate()

nums = sc.parallelize([1,2,3,4])
nums
```

Out[3]: ParallelCollectionRDD[2] at readRDDFromFile at PythonRDD.scala:274

In []: