

# Tutorial - Entorno de Trabajo Curso BigData Programa Ciencia de Datos Fundatec 2010

Alonso Nuñez Sanchez

Servicio/Aplicación: **Entorno de Trabajo / BigData**  
Tipo de Documento: **Funcional**

## Información de Control del Documento.

Información Documento.

<b>Nombre Documento</b>	Tutorial Entorno de Trabajo Curso BigData
<b>Nombre proyecto</b>	BigData / Ciencia de Datos
<b>Autor</b>	Alonso Nuñez Sanchez
<b>Versión</b>	1.0

### Historial de Modificaciones.

Versión	Fecha	Agregados/Modificaciones	Editor
1.0	23/11/2019	Creación Inicial del Documento	Alonso Nuñez Sanchez

## Contenido

<b>Información General.....</b>	<b>5</b>
Objetivo .....	5
Limitaciones .....	5
Aspectos generales importantes. ....	5
<b>MiniConda y entorno de trabajo.....</b>	<b>6</b>
Descripción .....	6
Descarga.....	6
Configuración.....	6
Consideraciones. ....	9
<b>JDK .....</b>	<b>10</b>
Descripción .....	10
Descarga.....	10
Configuración.....	12
Consideraciones. ....	12
<b>Java.....</b>	<b>13</b>
Descripción .....	13
Descarga.....	13
Configuración.....	15
Consideraciones. ....	15
<b>Spark.....</b>	<b>16</b>
Descripción .....	16
Descarga.....	16
Configuración.....	17
Consideraciones. ....	17
<b>FindSpark – buscar recursos de Spark para ser utilizados en Pyton....</b>	<b>18</b>
Descripción .....	18
Descarga.....	18
Configuración.....	19
Consideraciones. ....	19
<b>PySpark - ejecutar Spark en Python .....</b>	<b>20</b>
Descripción .....	20
Descarga.....	20
Configuración.....	20
Consideraciones. ....	20

## Información General.

### Objetivo

Documentar los pasos a seguir para preparar el entorno de trabajo necesario para el curso BigData del programa Ciencia de Datos

### Necesidades

Conocer los prerequisites de software necesarios para ejecutar las tareas propias del curso  
Brindar las herramientas y el conocimiento necesario para la puesta a punto del ambiente

### Limitaciones

El presente documento se limita a la preparación e instalación de los componentes de software necesarios para funcionar en el Sistema Operativo Windows, específicamente Windows 10

### Aspectos generales importantes.

El alcance del documento es únicamente para cubrir los aspectos académicos del curso en cuestión.

## MiniConda y entorno de trabajo.

### Descripción

Conda es un sistema open source para la gestión de entornos de trabajo.

Aunque fue creado para programas en Python, permite la ejecución de múltiples lenguajes de programación.

MiniConda es una instalación reducida de Conda. Para el alcance del curso, los recursos obtenidos con miniconda son suficientes

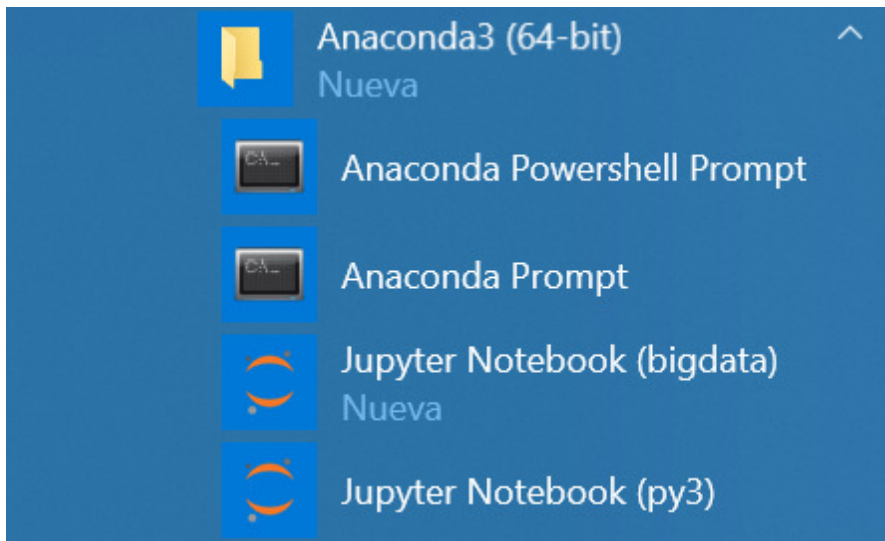
### Descarga

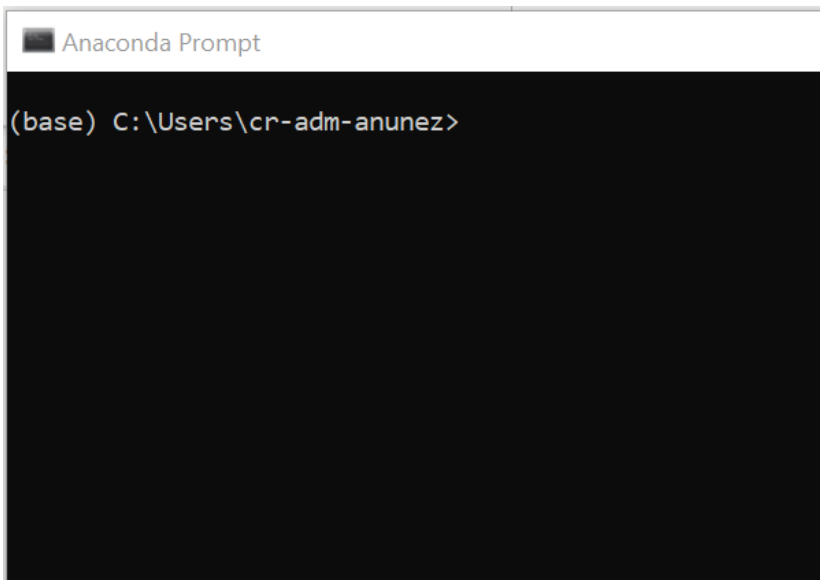
<https://conda.io/miniconda.html>

### Instalación

### Configuración

Una vez instalado, lo que tendremos es una consola que se puede encontrar en el explorador de aplicaciones, como “Anaconda Prompt”, dentro de la carpeta “Anaconda3”





Desde esta consola se ejecutan todos los comandos necesarios para la instalación y ejecución de paquetes, librerías y entornos de trabajo. Para efectos del presente documento, nos limitaremos a los comandos específicos necesarios

Lo primero será crear el ambiente de trabajo a utilizar, para esto se utiliza el comando **create**, de la siguiente forma:

```
>> conda create --name bigdata python=3
```

En este caso creamos un ambiente de trabajo de tipo python3 llamado "bigdata"

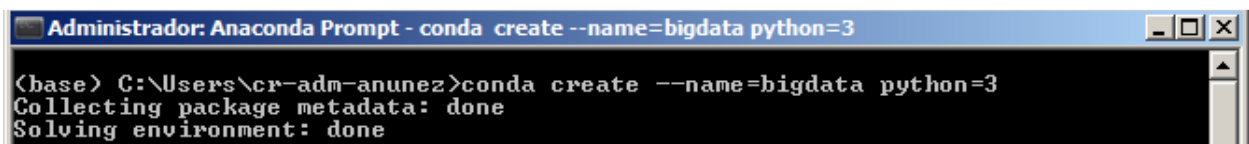
Posteriormente, se activa el ambiente recién creado, con el comando **activate**

```
>> conda activate bigdata
```

Activado el ambiente, el siguiente paso es instalar las librerías necesarias, comando **install**

```
>> conda install numpy
>> conda install matplotlib
>> conda install pandas
>> conda install jupyter notebook
>> conda install pyspark
>> conda install findspark
```

Dependiendo de la versión de conda, y los canales/repositorios configurados, pueda ser necesario un paso adicional para lograr instalar findspark, ver detalles en la sección "FinSpark" del presente documento



```

Administrador: Anaconda Prompt - conda create --name=bigdata python=3

The following packages will be downloaded:

package                                     | build                                     |
-----|-----|
ca-certificates-2019.10.16                | 0                                         | 163 KB
certifi-2019.9.11                         | py38_0                                   | 155 KB
openssl-1.1.1d                             | he774522_3                               | 5.7 MB
pip-19.3.1                                 | py38_0                                   | 1.9 MB
python-3.8.0                               | hff0d562_2                               | 19.6 MB
setuptools-41.6.0                         | py38_0                                   | 687 KB
sqlite-3.30.1                              | he774522_0                               | 962 KB
vs2015_runtime-14.16.27012                | hf0eaf9b_0                               | 2.4 MB
wheel-0.33.6                              | py38_0                                   | 53 KB
wincertstore-0.2                           | py38_0                                   | 15 KB
-----|-----|
Total:                                     |                                         | 31.6 MB

The following NEW packages will be INSTALLED:

ca-certificates    pkgs/main/win-64::ca-certificates-2019.10.16-0
certifi            pkgs/main/win-64::certifi-2019.9.11-py38_0
openssl            pkgs/main/win-64::openssl-1.1.1d-he774522_3
pip                pkgs/main/win-64::pip-19.3.1-py38_0
python             pkgs/main/win-64::python-3.8.0-hff0d562_2
setuptools         pkgs/main/win-64::setuptools-41.6.0-py38_0
sqlite             pkgs/main/win-64::sqlite-3.30.1-he774522_0
vc                 pkgs/main/win-64::vc-14.1-h0510ff6_4
vs2015_runtime     pkgs/main/win-64::vs2015_runtime-14.16.27012-hf0eaf9b_0
wheel              pkgs/main/win-64::wheel-0.33.6-py38_0
wincertstore       pkgs/main/win-64::wincertstore-0.2-py38_0

Proceed [Y/n]?

```

```

(base) C:\Users\cr-adm-anunez>conda activate bigdata

(bigdata) C:\Users\cr-adm-anunez>conda install numpy matplotlib pandas pyspark j
upyter notebook
Collecting package metadata: done
Solving environment: -

```



```

Administrador: Anaconda Prompt - conda install numpy matplotlib pandas pyspark jupyter notebook

The following packages will be downloaded:

package                                     build
-----
attrs-19.3.0                               py_0                                39 KB
bleach-3.1.0                               py_0                               111 KB
certifi-2019.9.11                         py37_0                             155 KB
colorama-0.4.1                             py_0                                17 KB
decorator-4.4.1                           py_0                                13 KB
importlib_metadata-0.23                   py37_0                              44 KB
ipykernel-5.1.3                           py37h39e3cac_0                     168 KB
ipython-7.9.0                             py37h39e3cac_0                      1.1 MB
ipywidgets-7.5.1                          py_0                               107 KB
jedi-0.15.1                               py37_0                             715 KB
jinja2-2.10.3                             py_0                                95 KB
jsonschema-3.2.0                          py37_0                             112 KB
jupyter_client-5.3.4                      py37_0                             160 KB
jupyter_core-4.6.1                        py37_0                              97 KB
matplotlib-3.1.1                         py37hc8f65d3_0                      6.6 MB
mkl-service-2.3.0                        py37hb782905_0                     200 KB
mkl_fft-1.0.15                           py37h14836fe_0                      137 KB
mkl_random-1.1.0                         py37h675688f_0                     270 KB
more-itertools-7.2.0                     py37_0                              99 KB
nbconvert-5.6.1                          py37_0                             512 KB
notebook-6.0.2                           py37_0                              7.7 MB
numpy-1.17.3                             py37h4ceb530_0                      5 KB
numpy-base-1.17.3                       py37hc3f5095_0                     4.8 MB
pandas-0.25.3                             py37ha925a31_0                     9.8 MB
parso-0.5.1                              py_0                                68 KB
pip-19.3.1                               py37_0                             1.9 MB
prompt_toolkit-2.0.10                    py_0                                227 KB
py4j-0.10.7                              py37_0                             251 KB
pyparsing-2.4.5                          py_0                                62 KB
pyspark-2.4.4                             py37he774522_0                      95 KB
python-3.7.5                             h8c8aaf0_0                         205.0 MB
python-dateutil-2.8.1                    py_0                                18.6 MB
pytz-2019.3                              py_0                                224 KB
pyzmq-18.1.0                             py_0                                231 KB
qtconsole-4.6.0                           py37ha925a31_0                     442 KB
setuptools-41.6.0                         py_0                                95 KB
six-1.13.0                               py37_0                             687 KB
terminado-0.8.3                           py37_0                             27 KB
testpath-0.4.4                            py_0                                25 KB
traitlets-4.3.3                           py37_0                              88 KB
wheel-0.33.6                              py37_0                             138 KB
widgetsnbextension-3.5.1                  py37_0                             58 KB
zipp-0.6.0                               py37_0                             1.8 MB
Total:                                     262.9 MB

```

Consideraciones.

## JDK

### Descripción

JDK (Java Development Kit) son un conjunto de herramientas que permiten el desarrollo de programas en el lenguaje de programación Java, provee diversas librerías para tal efecto.

### Descarga

<https://www.oracle.com/technetwork/java/javase/downloads/jdk13-downloads-5672538.html>

### Java SE Development Kit 13.0.1

You must accept the **Oracle Technology Network License Agreement for Oracle Java SE** to download this software.

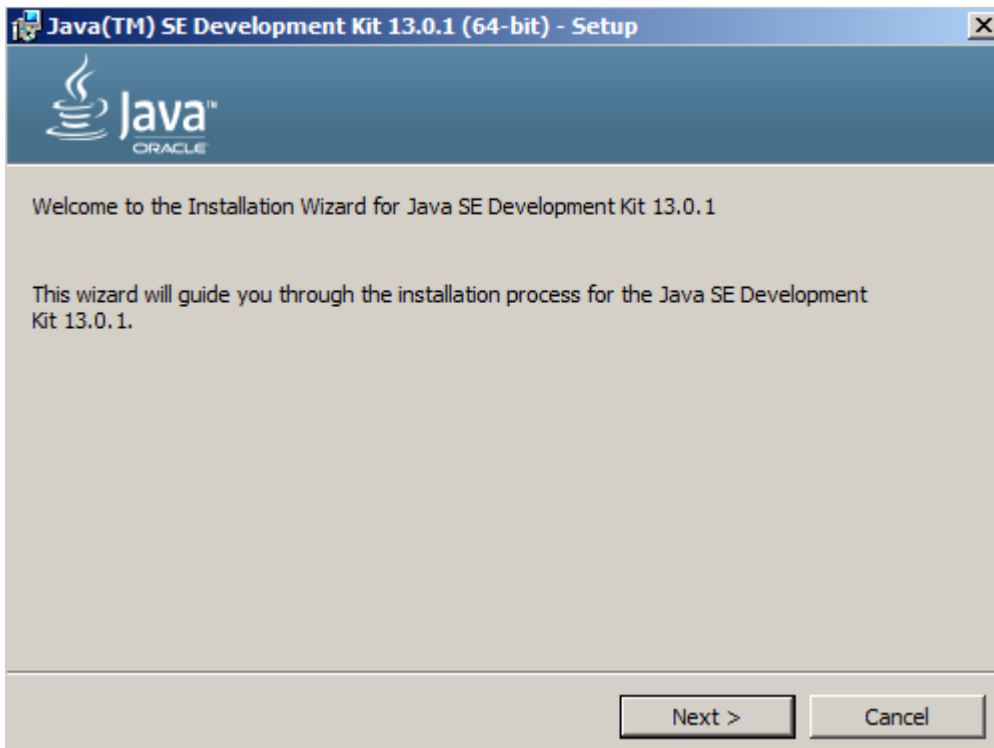
☒ **Accept License Agreement**
☐ Decline License Agreement

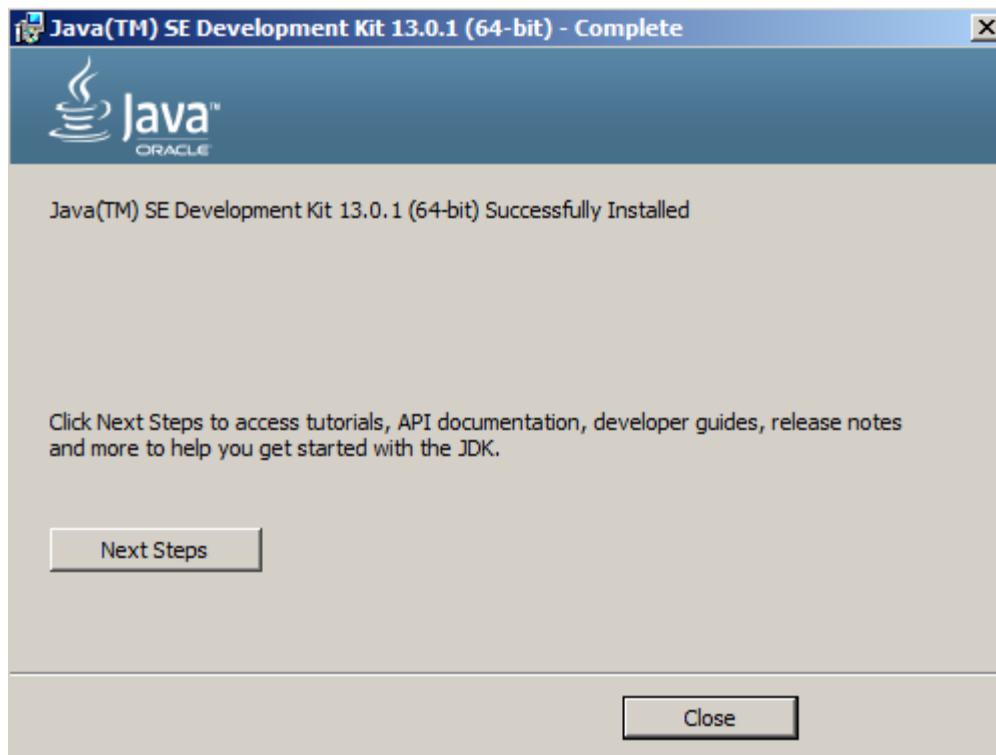
Product / File Description	File Size	Download
Linux	155.88 MB	<a href="#">jdk-13.0.1_linux-x64_bin.deb</a>
Linux	163.17 MB	<a href="#">jdk-13.0.1_linux-x64_bin.rpm</a>
Linux	180 MB	<a href="#">jdk-13.0.1_linux-x64_bin.tar.gz</a>
macOS	172.78 MB	<a href="#">jdk-13.0.1_osx-x64_bin.dmg</a>
macOS	173.11 MB	<a href="#">jdk-13.0.1_osx-x64_bin.tar.gz</a>
Windows	159.84 MB	<a href="#">jdk-13.0.1_windows-x64_bin.exe</a>
Windows	178.99 MB	<a href="#">jdk-13.0.1_windows-x64_bin.zip</a>

1. Ingresar al link anterior
2. Marcar "Accept Licence Agreement"
3. Seleccionar el Sistema Operativo y arquitectura, en este caso "Windows x64", [jdk-13.0.1\\_windows-x64\\_bin.exe](#)

## Instalación

Ejecutar el instalador y completar el wizard sin cambiar configuraciones





## Configuración

No requerida

Consideraciones.

# Java.

## Descripción

## Descarga

<https://www.java.com/es/download/win10.jsp>

Recursos de ayuda

- » [¿Qué es Java?](#)
- » [Eliminar versiones anteriores de Java](#)
- » [Desactivar Java](#)
- » [Mensajes de error](#)
- » [Solucionar problemas de Java](#)
- » [Otra ayuda](#)

Usuarios de Windows de 64 bits

¿Utiliza exploradores de 32 y 64 bits?

- » [Preguntas frecuentes sobre Java de 64 bits para Windows](#)

Instalación fuera de línea

¿Problemas al descargar? Intente con el [installer fuera de línea](#).

## Descargar Java para Windows

**Recomendado Version 8 Update 231 (Tamaño de archivo: 1.97 MB)**

Fecha de versión: 15 de octubre de 2019



### Actualización importante de la licencia de Oracle Java

**La licencia de Oracle Java ha cambiado para las versiones publicadas a partir del 16 de abril de 2019.**

El nuevo [acuerdo de licencia de Oracle Technology Network para Oracle Java SE](#) es sustancialmente diferente a las licencias de Oracle Java anteriores. La nueva licencia permite ciertos usos, como el uso personal y de desarrollo, sin coste alguno (aunque podría haber otros usos autorizados en licencias de Oracle Java anteriores que ya no estén disponibles). Revise las condiciones con atención antes de descargar y utilizar este producto. Puede consultar las preguntas frecuentes [aquí](#).

La licencia comercial y el soporte están disponibles con una [suscripción de Java SE](#) de bajo coste.

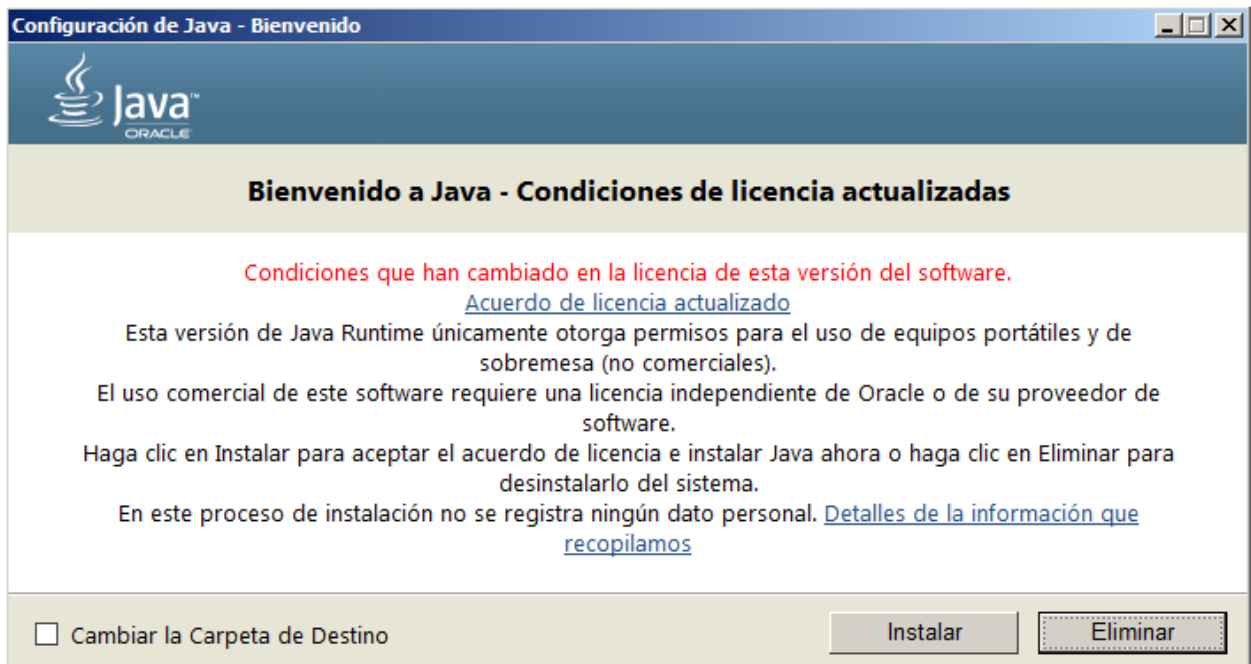
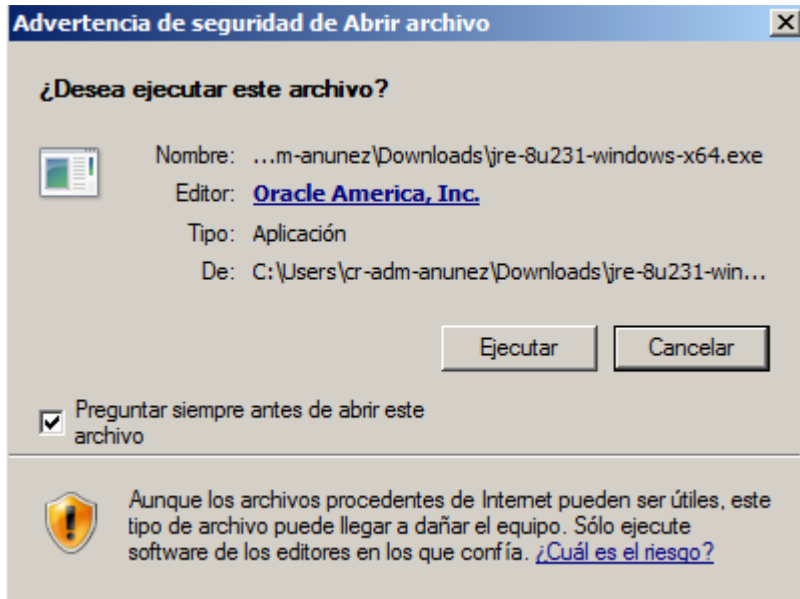
Oracle también ofrece la última versión de OpenJDK con la [licencia pública general](#) de código abierto en [jdk.java.net](#).

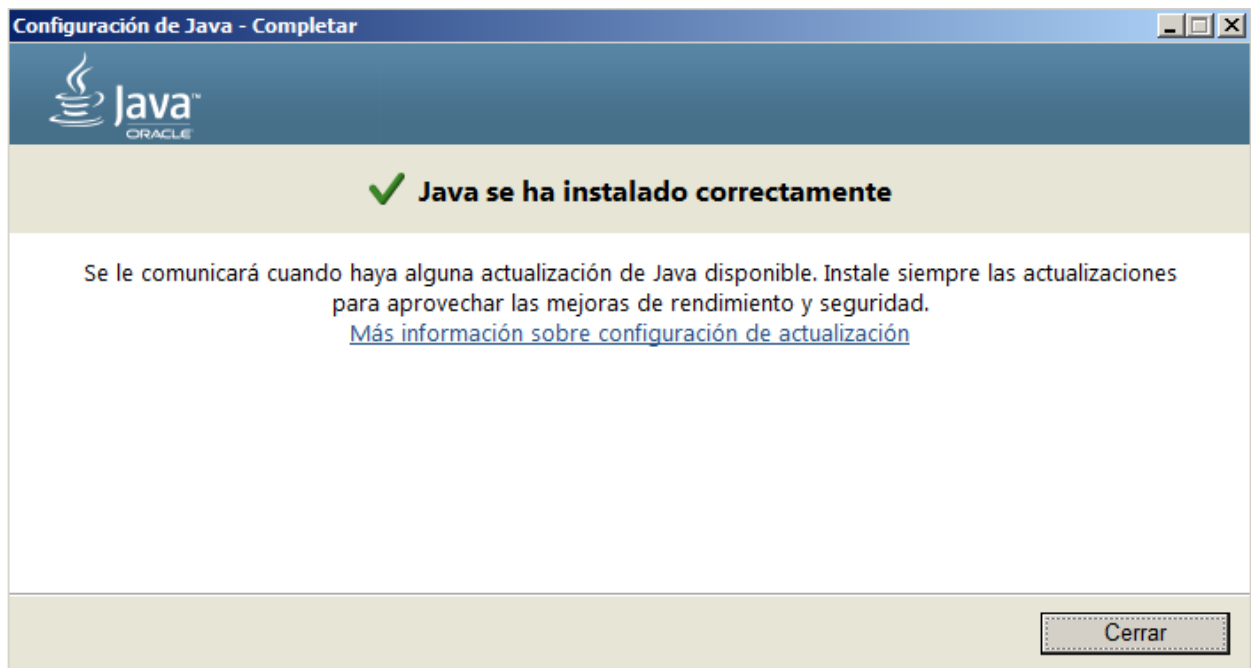
**Aceptar e iniciar descarga gratuita**

Al descargar Java, confirma que ha leído y acepta las condiciones del [acuerdo de licencia de Oracle Technology Network para Oracle Java SE](#)

## Instalación

Ejecutar el instalador y completar el wizard sin cambiar configuraciones





Configuración

No requerida

Consideraciones.

# Spark.

## Descripción

Es un framework enfocado en el manejo de grandes volúmenes de datos y la ejecución de cómputo intensivo sobre estos datos.

Es la evolución de lo que fue Hadoop, entre sus ventajas respecto a Hadoop está la reducción en tiempos de ejecución

## Descarga

<https://spark.apache.org/downloads.html>

Usaremos la versión 2.4.4, paquete “Pre-built for Apache Hadoop 2.7”

## Download Apache Spark™

1. Choose a Spark release:
2. Choose a package type:
3. Download Spark: [spark-2.4.4-bin-hadoop2.7.tgz](#)
4. Verify this release using the 2.4.4 [signatures](#), [checksums](#) and [project release KEYS](#).

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

## Instalación

No requiere instalación, solamente copiar el contenido de la carpeta descargada en una ubicación a elegir, en este caso será c:\spark\

- bin
- conf
- data
- examples
- jars
- kubernetes
- licenses
- python
- R
- sbin
- yarn
- LICENSE
- NOTICE
- README.md
- RELEASE



## Configuración

No requiere ninguna configuración, simplemente la copia del contenido descargado y el directorio elegido

## Consideraciones.

Las pruebas realizadas y primeros códigos ejecutados funcionan correctamente con la versión indicada. Esto no elimina la posibilidad de probar con la versión más reciente

## FindSpark – buscar recursos de Spark para ser utilizados en Python

### Descripción

FindSpark es una librería de Python para “buscar” y utilizar los recursos de Spark

Básicamente se usa para indicarle a Python la ruta donde se copiaron los archivos de Spark (ver sección “Spark” del presente documento)

### Descarga

Se descarga desde conda cuando se ejecuta el comando install

### Instalación

Se instala desde conda, en el ambiente de trabajo activo, con el comando install

```
>> conda install findspark
```

Dependiendo de la versión de conda, o de los canales/repositorios disponibles, pueda ser necesario indicar un canal distinto, esto se logra con el comando **config --add channels** indicando el canal deseado, en este caso será conda-forge, y volviendo a ejecutar el comando **install**

```
>> conda config --add channels conda-forge
```

```
<bigdata> C:\Users\cr-adm-anunez>conda install findspark
Collecting package metadata: done
Solving environment: failed

PackagesNotFoundError: The following packages are not available from current channels:

- findspark

Current channels:

- https://repo.anaconda.com/pkgs/main/win-64
- https://repo.anaconda.com/pkgs/main/noarch
- https://repo.anaconda.com/pkgs/free/win-64
- https://repo.anaconda.com/pkgs/free/noarch
- https://repo.anaconda.com/pkgs/r/win-64
- https://repo.anaconda.com/pkgs/r/noarch
- https://repo.anaconda.com/pkgs/msys2/win-64
- https://repo.anaconda.com/pkgs/msys2/noarch

To search for alternate channels that may provide the conda package you're
looking for, navigate to

    https://anaconda.org

and use the search bar at the top of the page.
```

```
<bigdata> C:\Users\cr-adm-anunez>conda config --add channels conda-forge
<bigdata> C:\Users\cr-adm-anunez>conda install findspark
Collecting package metadata: done
Solving environment: done
```

The following packages will be downloaded:

package	build		
ca-certificates-2019.9.11	hecc5488_0	181 KB	conda-forge
certifi-2019.9.11	py37_0	147 KB	conda-forge
findspark-1.3.0	py_1	6 KB	conda-forge
openssl-1.1.1d	hfa6e2cd_0	4.7 MB	conda-forge
Total:		5.0 MB	

The following NEW packages will be INSTALLED:

findspark conda-forge/noarch::findspark-1.3.0-py\_1

The following packages will be SUPERSEDED by a higher-priority channel:

ca-certificates pkgs/main::ca-certificates-2019.10.16~ --> conda-forge::ca-certificates-2019.9.11-hecc5488\_0  
 certifi pkgs/main --> conda-forge  
 openssl pkgs/main::openssl-1.1.1d-he774522\_3 --> conda-forge::openssl-1.1.1d-hfa6e2cd\_0

Proceed <[y]/n>?

## Configuración

No requerida

Consideraciones.

## PySpark - ejecutar Spark en Python

### Descripción

PySpark es un API de Python para poder ejecutar Spark

### Descarga

Se descarga desde conda cuando se ejecuta el comando install

### Instalación

Se instala desde conda, en el ambiente de trabajo activo, con el comando install

>> **conda install pyspark**

### Configuración

No requerida

### Consideraciones.