

Install Apache Spark Stand-alone

When you need to scale up your machine learning abilities, you will need a distributed computation. PySpark interface to Spark is a good option. Here is a simple guide, on installation of Apache Spark with PySpark, alongside your anaconda, on your windows machine.

Install Java 8

Before you can start with spark, you need to make sure you have java 8 installed, or to install it. First of all, check if JAVA is installed: open cmd (windows command prompt) , or anaconda prompt, from start menu and run:

```
C:\>java -version
```

You Should get something like:

```
java version "1.8.0_144"  
Java(TM) SE Runtime Environment (build 1.8.0_144-b01)  
Java HotSpot(TM) Client VM (build 25.144-b01, mixed mode, sharing)
```

If you don't have Java 8 installed, you need to install it.

Go to [Java's official](#) download website, accept Oracle license and download Java JDK 8, suitable to your system.

Java SE Development Kit 8u201		
You must accept the Oracle Binary Code License Agreement for Java SE to download this software.		
Thank you for accepting the Oracle Binary Code License Agreement for Java SE; you may now download this software.		
Product / File Description	File Size	Download
Linux ARM 32 Hard Float ABI	72.98 MB	jdk-8u201-linux-arm32-vfp-hflt.tar.gz
Linux ARM 64 Hard Float ABI	69.92 MB	jdk-8u201-linux-arm64-vfp-hflt.tar.gz
Windows x86	197.66 MB	jdk-8u201-windows-i586.exe
Windows x64	207.46 MB	jdk-8u201-windows-x64.exe

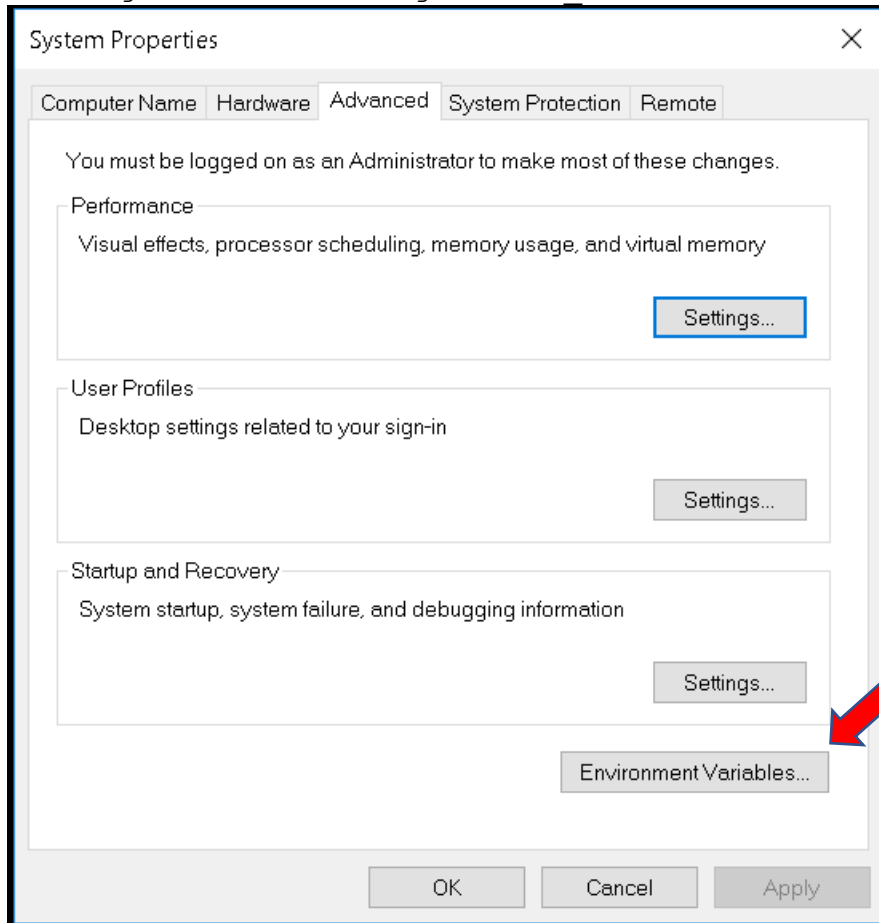
Run the executable, and JAVA by default will be installed in C:\Program Files\Java\jdk1.8.0_201

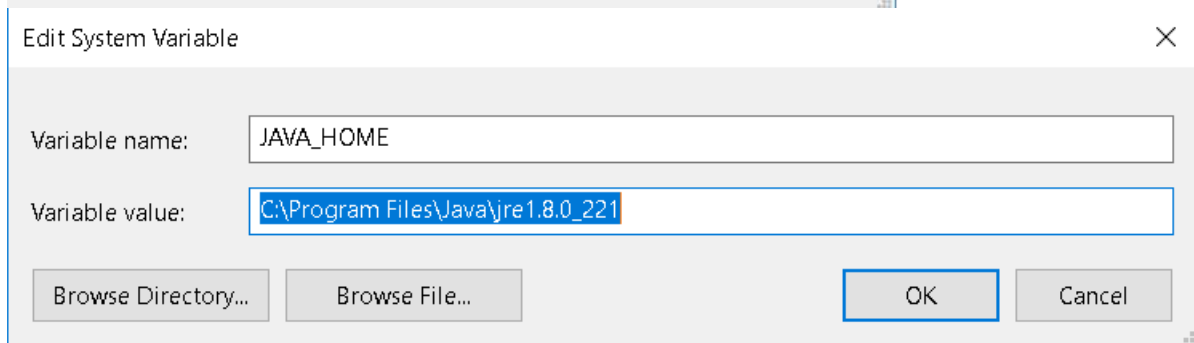
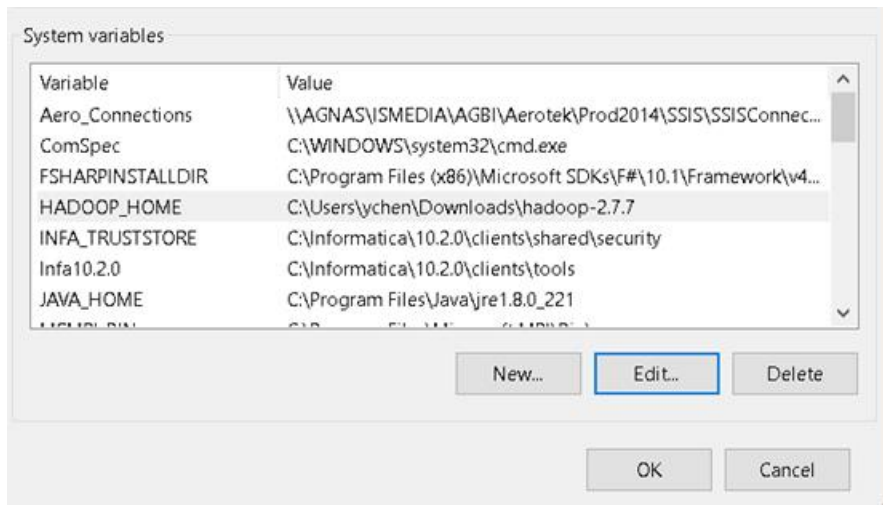
Add the following environment variable:

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

Add to PATH variable the following directory:

C:\Program Files\Java\jdk1.8.0_201\bin





Download and Install Spark

Go to [Spark home page](#), and download the **.tgz** file from 2.4.3 version.

Download Apache Spark™

1. Choose a Spark release: **2.4.3 (May 07 2019)** ▼
2. Choose a package type: **Pre-built for Apache Hadoop 2.7 and later** ▼
3. Download Spark: [spark-2.4.3-bin-hadoop2.7.tgz](#) 
4. Verify this release using the 2.4.3 [signatures](#), [checksums](#) and [project release KEYS](#).

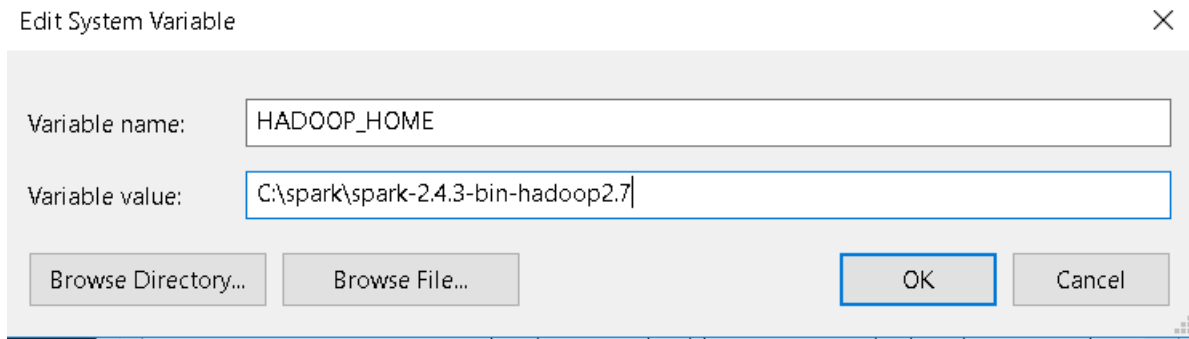
Note that, Spark is pre-built with Scala 2.11 except version 2.4.2, which is pre-built with Scala 2.12.

Extract the file to your chosen directory (7z can open tgz). In my case, it was C:\spark. There is another compressed directory in the tar, extract it (into here) as well.

Setup the environment variables

SPARK_HOME = C:\spark\spark-2.4.3-bin-hadoop2.7

HADOOP_HOME = C:\spark\spark-2.4.3-bin-hadoop2.7



Add the following path to PATH environment variable:
C:\spark\spark-2.4.3-bin-hadoop2.7\bin

Check PySpark Installation

In your anaconda prompt, or any python supporting cmd, type pyspark, to enter pyspark shell. You supposed to see the following:

```
Command Prompt - pyspark
Microsoft Windows [Version 10.0.16299.1217]
(c) 2017 Microsoft Corporation. All rights reserved.

C:\Users\ychen>pyspark
Python 3.7.3 (default, Mar 27 2019, 17:13:21) [MSC v.1915 64 bit (AMD64)] :: Anaconda, Inc. on win32

Warning:
This Python interpreter is in a conda environment, but the environment has
not been activated. Libraries may fail to load. To activate this environment
please see https://conda.io/activation

Type "help", "copyright", "credits" or "license" for more information.
19/07/22 23:35:02 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java cl
asses where applicable
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
Welcome to

  ____      _
 / ___|  _ \| | | |
 \___ \| |_) | |_| |
  ___) | |_) | | | |
 |____|_|_|\___|_|_|_|

 version 2.4.1

Using Python version 3.7.3 (default, Mar 27 2019 17:13:21)
SparkSession available as 'spark'.
>>>
```

To exit pyspark shell, type Ctrl-z and enter. Or the python command exit()

Setup Hadoop with winutils.exe

In hadoop binaries repository of <https://github.com/steveloughran/winutils>, choose your hadoop version, then go to bin, and download the winutils.exe file. In my case: <https://github.com/steveloughran/winutils/blob/master/hadoop-2.7.1/bin/winutils.exe>

Save winutils.exe in to bin directory of your spark installation, **SPARK_HOME\bin** directory. In my case: C:\spark\spark-2.3.2-bin-hadoop2.7\bin. Now the trick. It's not a must, things did not work well for me without it.

1. Create the folder **C:\tmp\hive**
2. Execute the following command in **cmd** started using the option **Run as administrator or Run Elevated**.

```
winutils.exe chmod -R 777 C:\tmp\hive
winutils.exe ls -F C:\tmp\hive
```

The output is something of the sort:
drwxrwxrwx|1|LAPTOP-.....

PySpark with Jupyter notebook

Install conda findspark to access spark instance from jupyter notebook. By writing:

```
C:\>conda install -c conda-forge findspark
```

Open your python jupyter notebook, and write inside:

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

Last line will output SPARK_HOME path. The following lines will create a simple lambda function that can run parallelly on spark context.

```
from pyspark import SparkContext, SparkConf
from pyspark.sql import SparkSessionconf =
pyspark.SparkConf().setAppName('appName').setMaster('local')
sc = pyspark.SparkContext(conf=conf)
spark = SparkSession(sc)
nums = sc.parallelize([1,2,3,4])
nums.map(lambda x: x*x).collect()
sc.stop()
```

The following lines will test/show some important environment setup to use pyspark on jupyter notebook.

```
import os
print(os.environ['SPARK_HOME'])
print(os.environ['JAVA_HOME'])
print(os.environ['PATH'])
```