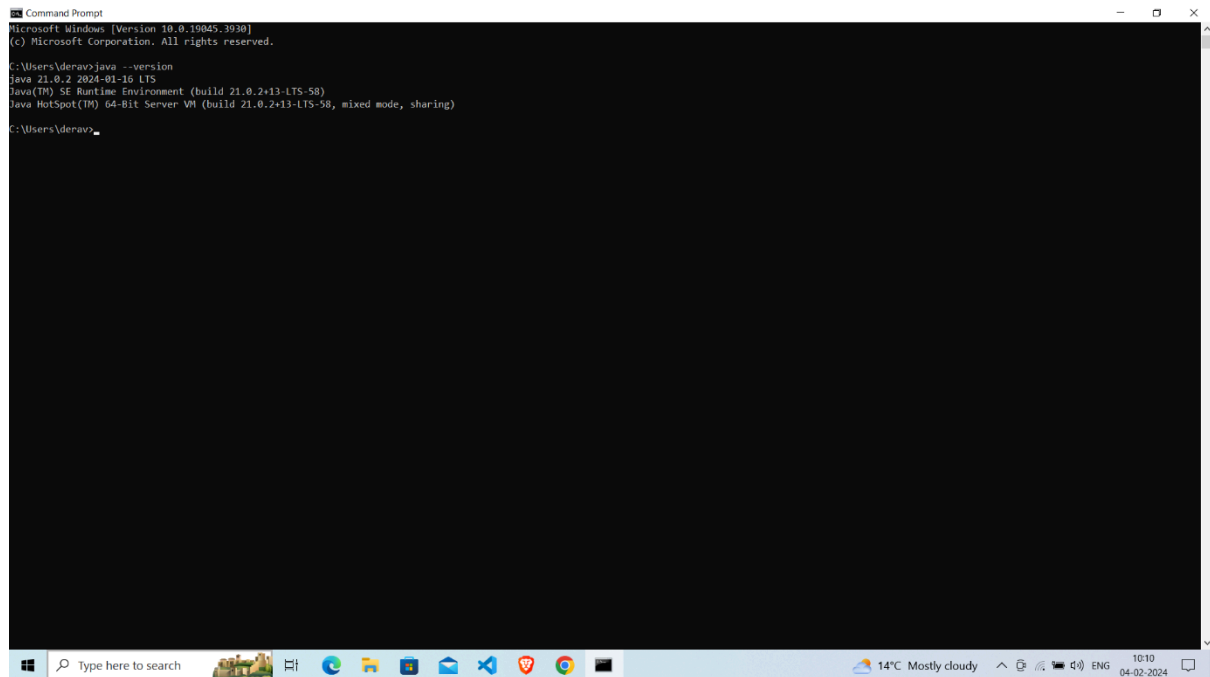


Installation of Apache Spark on my system

Step 1: Install Java 8

Start > type cmd> click Command Prompt.

java -version



```
Microsoft Windows [Version 10.0.19045.3930]
(c) Microsoft Corporation. All rights reserved.

C:\Users\derav>java -version
java 21.0.2 2024-01-16 LTS
Java(TM) SE Runtime Environment (build 21.0.2+13-LTS-58)
Java HotSpot(TM) 64-Bit Server VM (build 21.0.2+13-LTS-58, mixed mode, sharing)

C:\Users\derav>
```

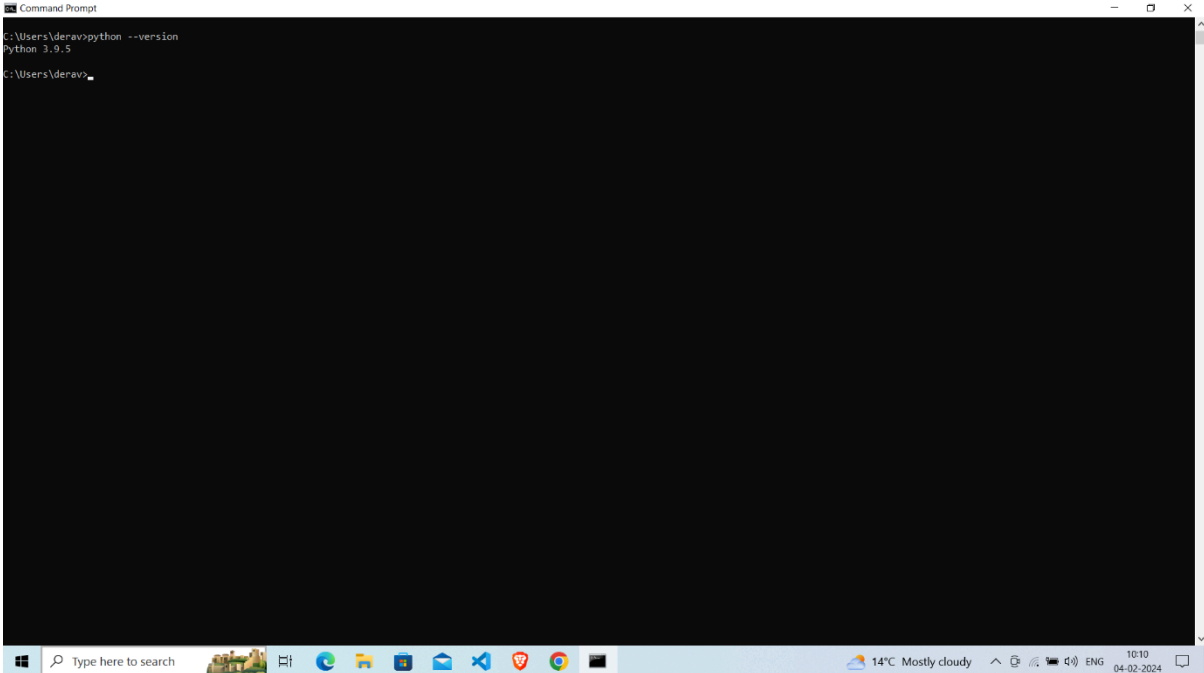
If you don't have Java installed:

1. Open a browser window, and navigate to <https://java.com/en/download/>.
2. Click the Java Download button and save the file to a location of your choice.
3. Once the download finishes double-click the file to install Java.

Step 2: Install Python

1. To install the Python package manager, navigate to <https://www.python.org/> in your web browser.
2. Mouse over the Download menu option and click Python 3.11.
3. Once the download finishes, run the file.
4. Near the bottom of the first setup dialog box, check off Add Python 3.11 to PATH. Leave the other box checked.
5. Next, click Customize Installation.
6. You can leave all boxes checked at this step, or you can uncheck the options you do not want.

7. Select the box Install for all users and leave other boxes as they are.
8. Under Customize install location, click Browse and navigate to the C drive. Add a new folder and name it Python.
9. Select that folder and click OK.
10. When the installation completes, click the Disable path length limit option at the bottom and then click Close.
11. If you have a command prompt open, restart it. Verify the installation by checking the version of Python:



```
Command Prompt
C:\Users\derav>python --version
Python 3.9.5
C:\Users\derav>
```

Step 3: Download Apache Spark

1. Open a browser and navigate to <https://spark.apache.org/downloads.html>.
2. Under the Download Apache Spark heading, there are two drop-down menus. Use the current non-preview version.
Choose a Spark release -> 3.5.0
Choose a package type -> Pre-built for Apache Hadoop 3.
3. Click the spark-3.5.0-bin-hadoop3.tgz link
4. A page with a list of mirrors loads where you can see different servers to download from. Pick any from the list and save the file to your Downloads folder.

Step 4: Add winutils.exe File

Download the winutils.exe file for the underlying Hadoop version for the Spark installation you downloaded.

1. Navigate to this URL <https://github.com/cdarlint/winutils> and inside the bin folder, locate winutils.exe, and click it.
2. Find the Download button on the right side to download the file.
3. Now, create new folders Hadoop and bin on C: using Windows Explorer or the Command Prompt.
4. Copy the winutils.exe file from the Downloads folder to C:\Hadoop\bin

Step 5: Configure Environment Variables

Configuring environment variables in Windows adds the Spark and Hadoop locations to your system PATH. It allows you to run the Spark shell directly from a command prompt window.

1. Click Start and type environment.
2. Select the result labelled Edit the system environment variables.
3. A System Properties dialog box appears. In the lower-right corner, click Environment Variables and then click New in the next window.
4. For Variable Name type SPARK_HOME.
5. For Variable Value type C:\Spark\spark-3.5.0-bin-hadoop3 and click OK. If you changed the folder path, use that one instead.
6. In the top box, click the Path entry, then click Edit. Be careful with editing the system path. Avoid deleting any entries already on the list.
7. You should see a box with entries on the left. On the right, click New.
8. The system highlights a new line. Enter the path to the Spark folder C:\Spark\spark-3.5.0-bin-hadoop3\bin. We recommend using %SPARK_HOME%\bin to avoid possible issues with the path.
9. Repeat this process for Hadoop and Java.

For Hadoop, the variable name is HADOOP_HOME and for the value use the path of the folder you created earlier: C:\Hadoop. Add C:\Hadoop\bin to the Path variable field, but we recommend using %HADOOP_HOME%\bin.

For Java, the variable name is JAVA_HOME and for the value use the path to your Java JDK directory (example, C:\Program Files\Java\<jdk_version>).

10. Click OK to close all open windows.

Step 6: Launch Spark

1. Open a new command prompt Window using the right-click and Run as administrator:

2. To start Spark, enter:

```
Command Prompt
'spark' is not recognized as an internal or external command,
operable program or batch file.

C:\Users\derav>pyspark --version
Welcome to
      ____              __
     /  _ \            /  |
    /  / \   ____     /   | _____/
   /  /_>  / __ \   /    | |
  /_____/ /___/ \_/_____|_|
Spark version 3.5.0

Using Scala version 2.12.18, Java HotSpot(TM) 64-Bit Server VM, 21.0.2
Branch HEAD
Compiled by user ubuntu on 2023-09-09T01:53:20Z
Revision ce5ddad998373636e94071e7cef2f31021add07b
Url https://github.com/apache/spark
Type --help for more information.

C:\Users\derav>
```

Inbox (1,439) - pu

pyspark - Google

sparksetup.docx

Apache spark con

WhatsApp

ln (2) Install Apache

Pandas Data Pro

Top 67 Work Fro

PySparkShell

Not secure | pushpa-pc-4040/jobs/

TCS NQT Coding Sh... Striver's SDE Sheet... gitignore.io - Creat... Linkeder - Create a... Practice WhatsApp DevOps Untitled - dbdiagra... Online Markdown E... leetcode-contest

spark 3.5.0 Jobs Stages Storage Environment Executors PySparkShell application UI

Spark Jobs ^(?)

User: derav

Total Uptime: 1.2 min

Scheduling Mode: FIFO

Event Timeline

☐ Enable zooming

Executors

Added

Removed

Jobs

Succeeded

Failed

Running

| | | | | | | | | | | | | | | | |
|--|----------|-----|----------|-----|-----|-----|-----|----------|-----|-----|-----|-----|----------|-----|-----|
| | 600 | 800 | 000 | 200 | 400 | 600 | 800 | 000 | 200 | 400 | 600 | 800 | 000 | 200 | 400 |
| | 13:33:19 | | 13:33:20 | | | | | 13:33:21 | | | | | 13:33:22 | | |

Type here to search

16°C Mostly cloudy

13:34

04-02-2024