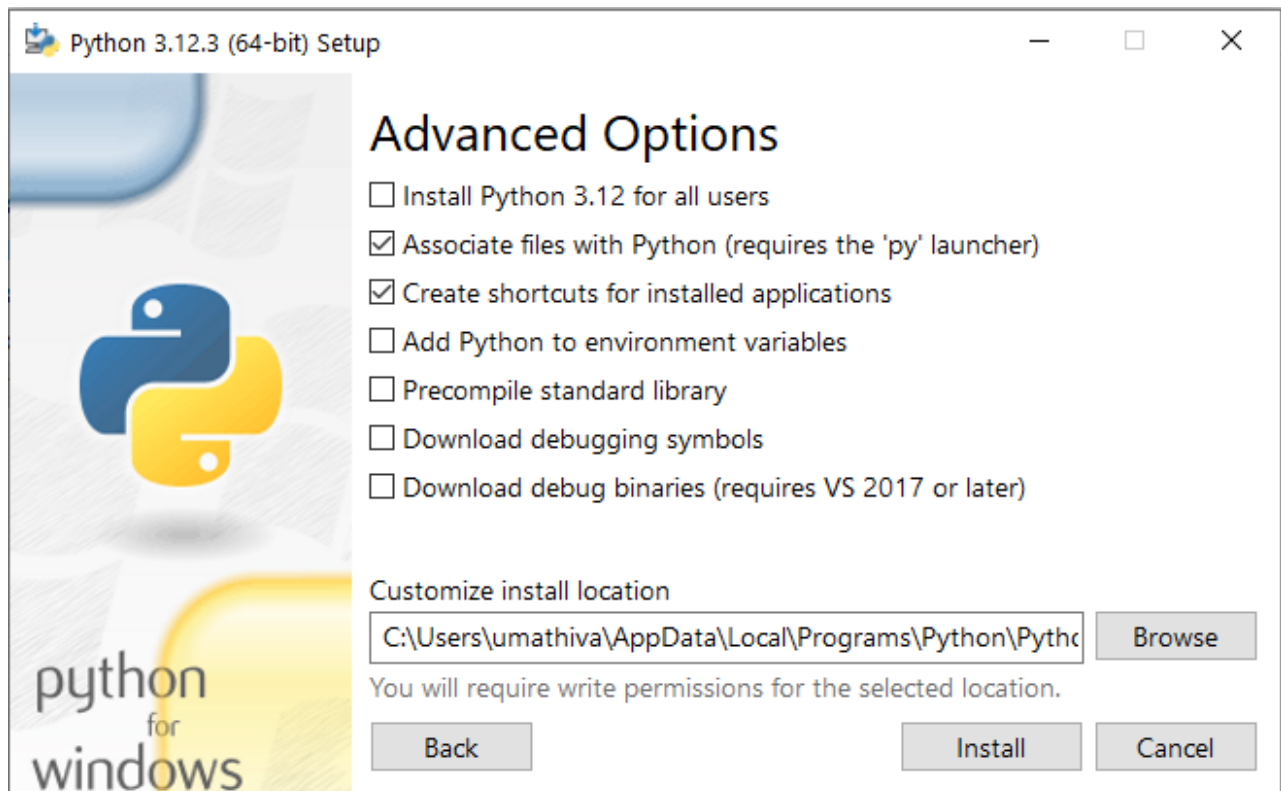# Apache Spark Deployment

# Introduction

The spark applications can be developed and executed on a spark cluster, to run/execute the spark applications first you must download, extract and copy the Apache Spark, and Hadoop to your local development system.  Please follow the below instructions to run Spark applications on your local machine.

# Apache Spark Installation and Configuration

**Pre-requisite:** Ensure to download and install the latest version of python on your local machine (make sure to choose/add Python related variables)
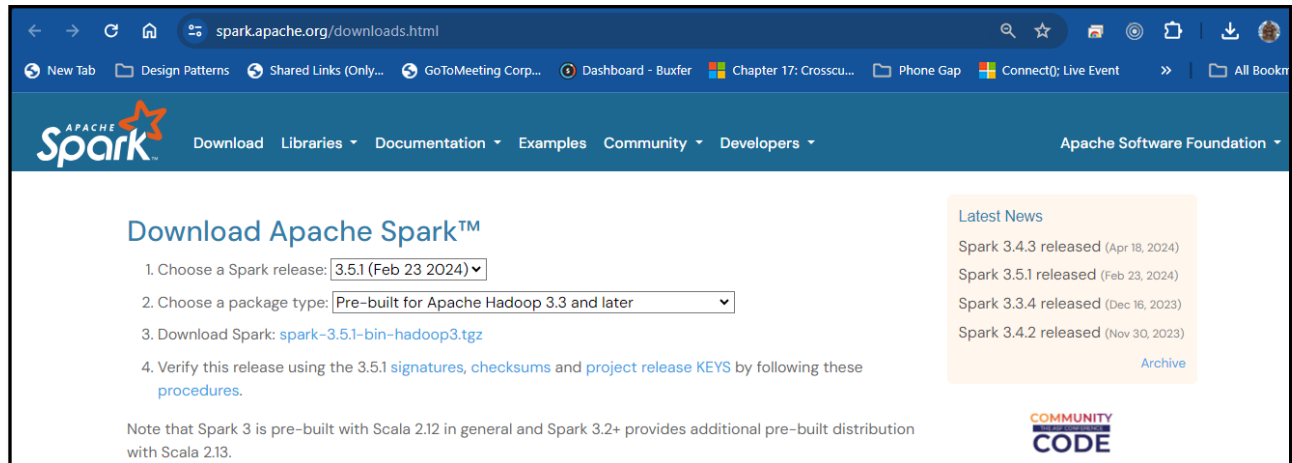
[https://www.python.org/ftp/python/3.12.3/python-3.12.3-amd64.exe](https://www.python.org/ftp/python/3.12.3/python-3.12.3-amd64.exe) (Make sure to choose custom installation to get this page where you can choose the '**Add Python to environment variables**)

**Step 1:**

Download the Apache spark (version 3.5.1) from the below location.

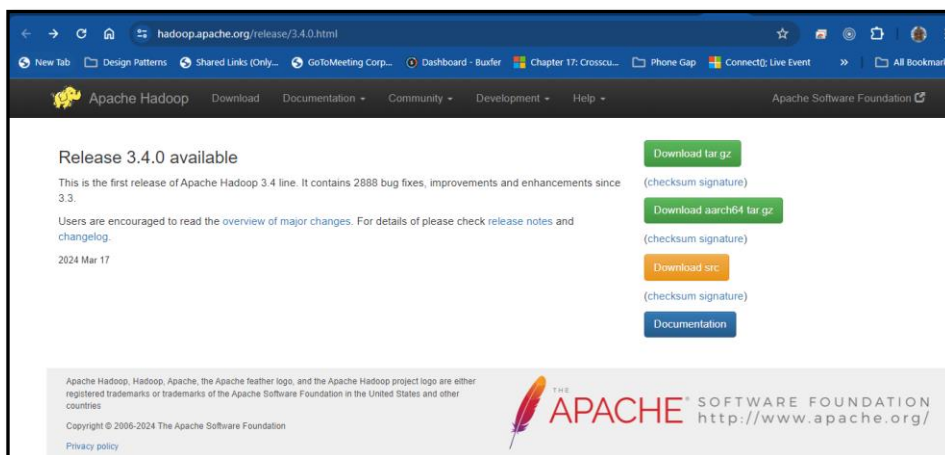https://spark.apache.org/downloads.html



https://www.apache.org/dyn/closer.lua/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz

Extract the **.tgz** file using (winrar utility) into a folder something like **C:\Spark**

**Step 2:**

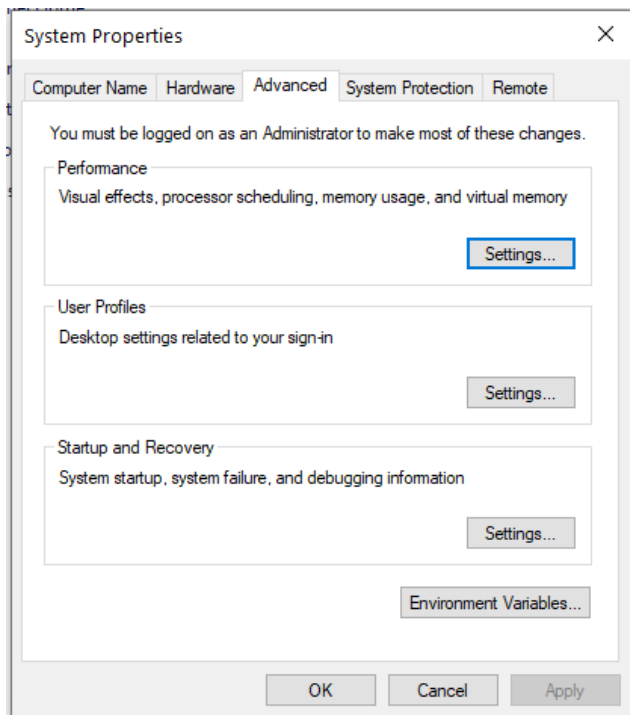Download Hadoop (version 3.4.0) from the below location.

https://hadoop.apache.org/release/3.4.0.html

https://archive.apache.org/dist/hadoop/common/hadoop-3.4.0/hadoop-3.4.0.tar.gz

Extract the **.tgz** file using (winrar utility) into a folder like **C:\Hadoop**

**Step 3:**

Now add the below environment variables (**user variables**) to your system



**HADOOP_HOME**: Despite Spark can run without Hadoop, the version we downloaded is prebuilt for Hadoop 3.3 and looks in the code for it (**C:\Hadoop**). To fix this inconvenience we need to set this variable to the folder containing the winutils.exe file.

HADOOP_HOME = C:\Spark\hadoop

**JAVA_HOME:** usually this variable should be already set when you install java, but it is better to verify that exists and is correct.

JAVA_HOME = C:\Program Files\Java\jdk-12.0.2

**Trainer:** Udayakumar M                                                                                  Data Engineering

**SPARK_HOME**: This variable should contain the path of the bin folder of where you uncompressed Spark .tgz file which is **C:\Spark**

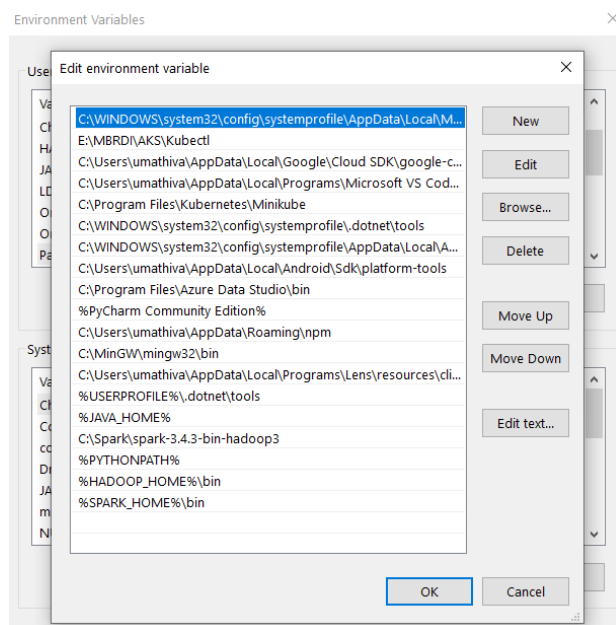SPARK_HOME = C:\Spark\spark-3.4.3-bin-hadoop3

After the above variables are set, now add the below paths to PATH variable.

%JAVA_HOME%

%PYTHONPATH%

%HADOOP_HOME%\bin

%SPARK_HOME%\bin



Now to develop and run your spark programs from VS Code Jupyter Notebook add the below environment variables

**PYSPARK_PYTHON** – Specify the path of the Python location where you have installed.

**PYSPARK_DRIVER_PYTHON –** Specify the value '**jupyter'** here.

**PYSPARK_DRIVER_PYTHON_OPTS –** Specify the value '**notebook'** here.

We are finally done and could start the spark-shell which is an interactive way to analyze data using Scala or Python. In this way we are going also to test our Spark installation.

From the command prompt go to **the Spark f**older and type the following command to run the Scala shell:



Notice the Spark Context and Spark Session got created successfully, optionally you shall navigate to http://127.0.0.1:4040 where Spark context web UI will be available.



Now run few scala commands to check the spark results

scala>println("Welcome to Dell")

scala>val data = sc.parallelize(List(30,40,50))

scala>data.collect

Hurray, you have successfully deployed Spark on your machine and ready to run spark applications in your local machine.