🔒 **ey-advisory-technology-testing** / **etl-step-definitions**   Internal

<> Code    ⊙ Issues    ⁂ Pull requests    ▶ Actions    ▦ Projects    📖 **Wiki**    ⊙ Security

# Installation

Edit    New Page                                                    Jump to bottom

Gajendra Varadhan edited this page on Dec 23, 2020 · 1 revision

## Spark

For Mac OS, there is no further installations needed. For Windows 10 and above, please follow the below steps for Installing Apache Spark in your machine.

## Prerequisites

- A system running Windows 10
- A user account with administrator privileges (required to install software, modify file permissions, and modify system PATH)
- Command Prompt or Powershell
- A tool to extract .tar files, such as 7-Zip

## Step 1: Install Java 11 or above

Apache Spark requires Java 11 or above. I currently have Java 14 which works without issues. You can check to see if Java is installed using the command prompt.

Open the command line by clicking Start > type cmd > click Command Prompt.

Type the following command in the command prompt:

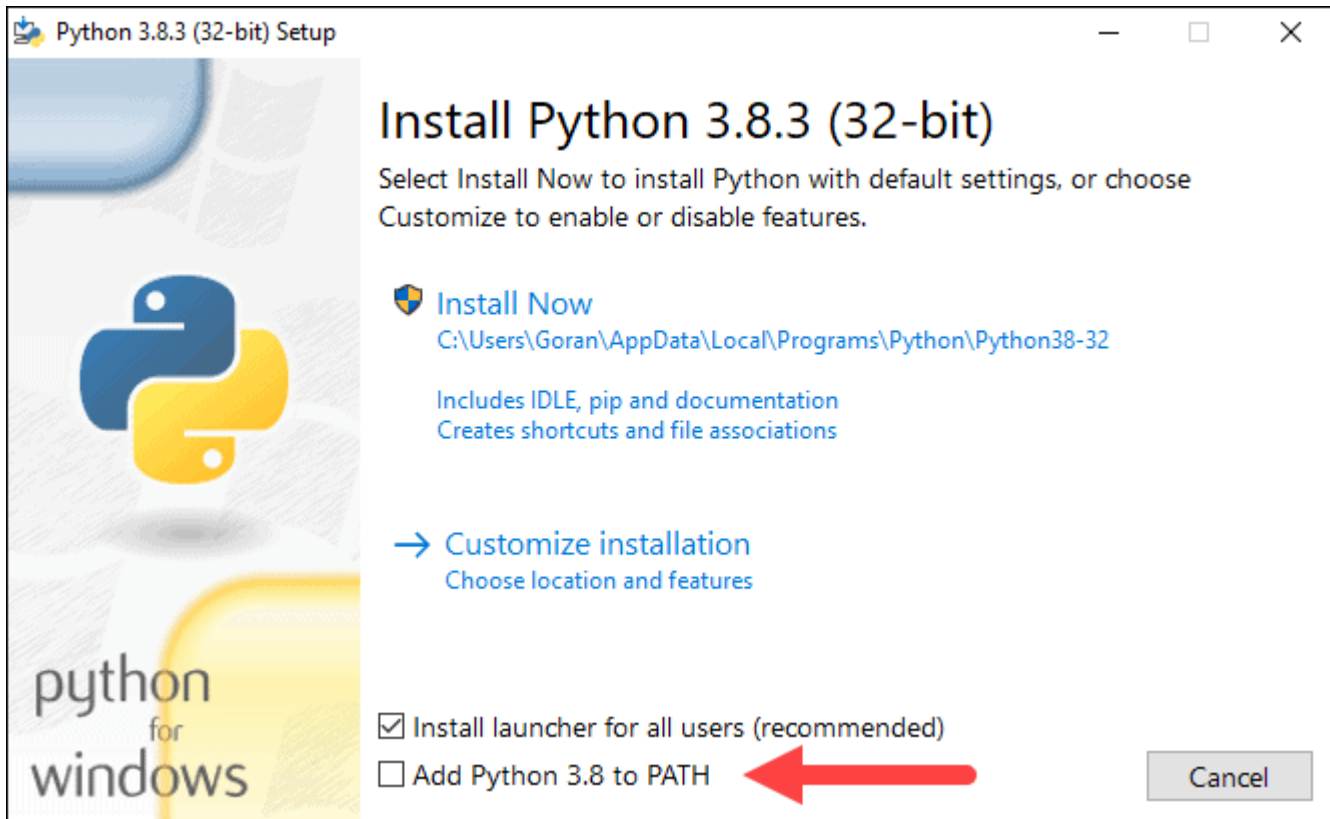```
java -version
```

If you don't have Java installed:

Download it from https://jdk.java.net/

## Step 2: Install Python (Optional)

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 3.8.3**. **3.8.3** is the latest version at the time of writing this manual.

3. Once the download finishes, run the file.



4. Near the bottom of the first setup dialog box, check off Add Python 3.8 to PATH. Leave the other box checked.

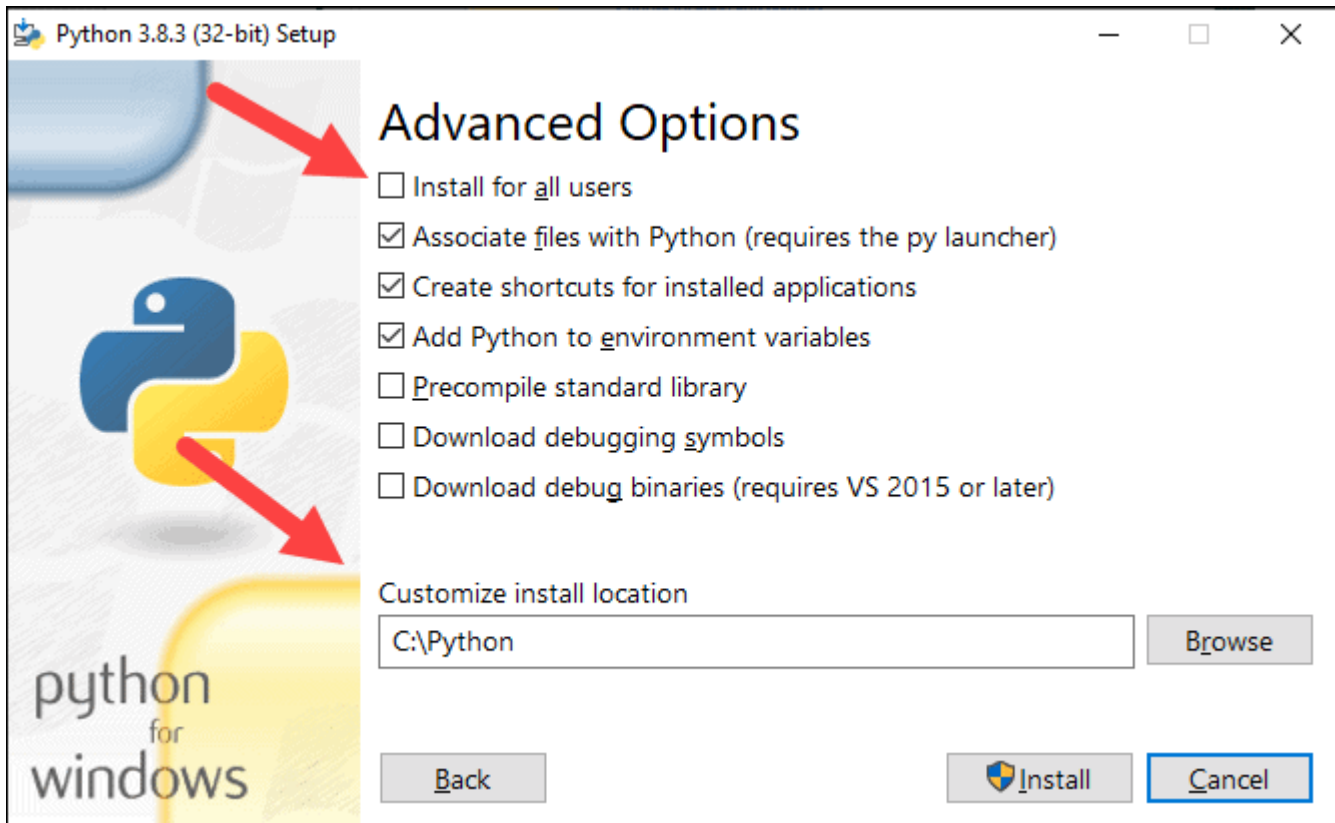5. Next, click Customize installation.

6. You can leave all boxes checked at this step, or you can uncheck the options you do not want.

7. Click Next.

8. Select the box Install for all users and leave other boxes as they are.

9. Under Customize install location, click Browse and navigate to the C drive. Add a new folder and name it Python.

10. Select that folder and click OK.

11. Click Install, and let the installation complete.

12. When the installation completes, click the Disable path length limit option at the bottom and then click Close.

13. If you have a command prompt open, restart it. Verify the installation by checking the version of Python:

```
python --version
```

The output should print `Python 3.8.3`.

## 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.

It can be any latest version. In our case, let's choose a Spark release from the drop-down menu `3.0.1 (Sep 02, 2020)`. In the second drop-down Choose a package type, select `Pre-built for Apache Hadoop 3.2 and later`. 3. Click the **spark-3.0.1-bin-hadoop3.2.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 5: Install Apache Spark

Installing Apache Spark involves extracting the downloaded file to the desired location.

1. Extract spark-3.0.1-bin-hadoop3.2.tgz to the desired location (say **C:**). Usually the content is again an archive file. So, extract again and you can see the bin folder.
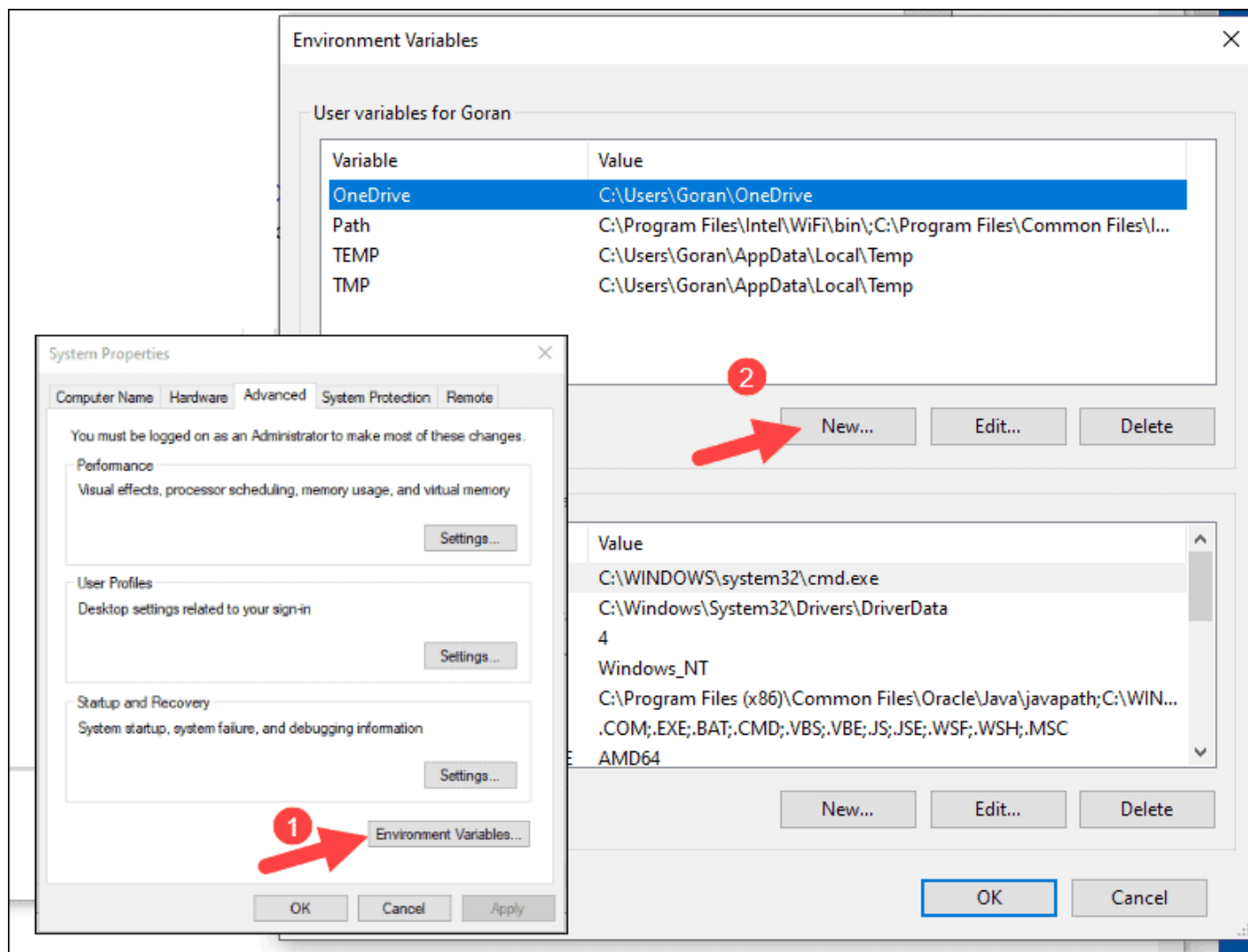
## Step 6: Add winutils.exe and hadoop.dll files

Download the winutils.exe and hadoop.dll files for the Hadoop version selected above.

1. Navigate to this URL https://github.com/cdarlint/winutils/tree/master/hadoop-3.2.1/bin, download the winutils.exe, and hadoop.dll files by clicking on them and clicking the download button on the top right.

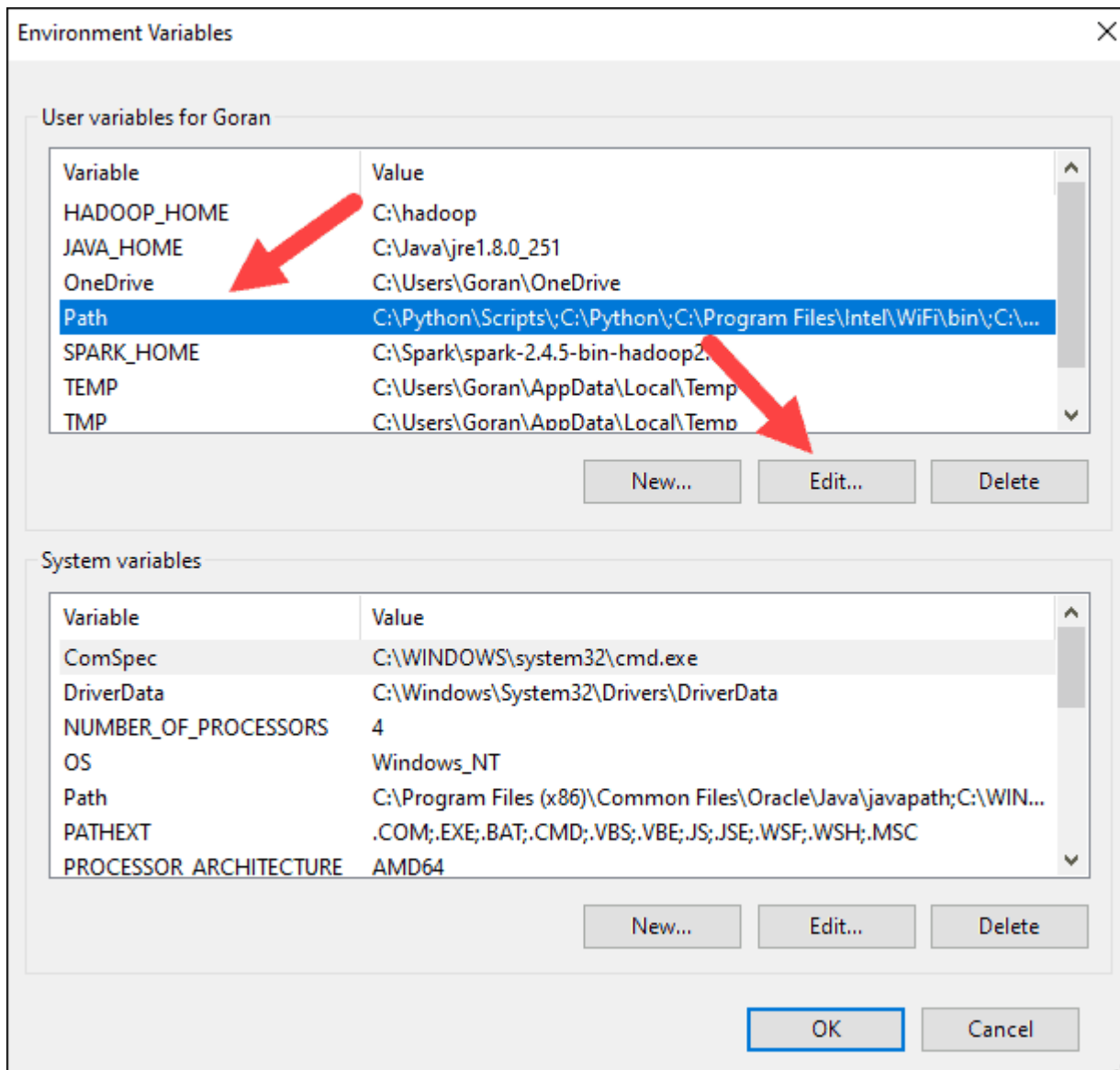2. Paste both the files under the bin folder of the spark extracted above.

## Step 7: Configure Environment Variables

This step adds the Spark and Hadoop locations to your system PATH. It allows you to run the Spark shell directly from a command prompt window.
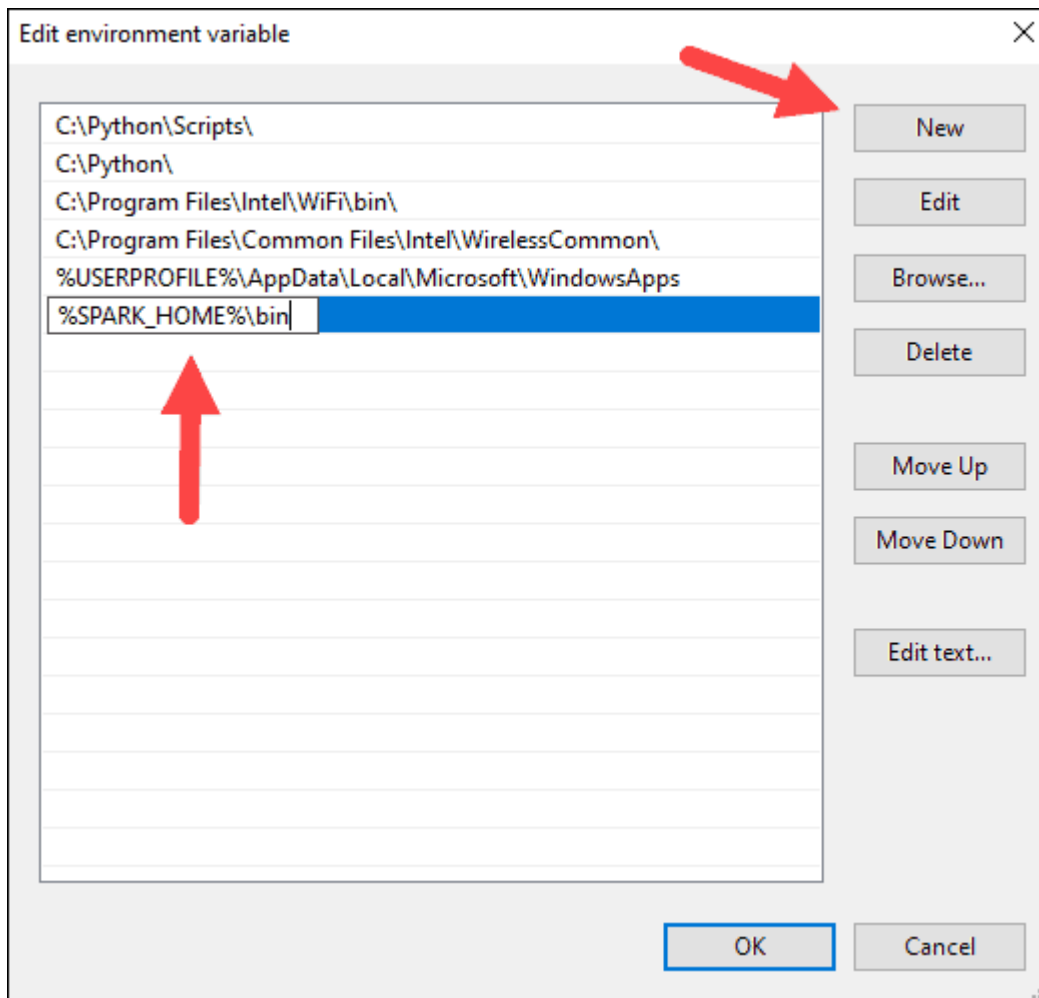
1. Click Start and type environment.

2. Select the result labeled 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. Create a new Variable with **Name **as **SPARK_HOME**.

5. For Variable **Value** type the folder where the spark bin is located, and click OK. If you changed the folder path, use that one instead.

6. Create another variable with **Name **as **HADOOP_HOME**

7. For the value, type **%SPARK_HOME%**

8. Edit the variable named Path.

9. Add an entry like **%SPARK_HOME%\bin**

10. Click **OK** to close all open windows.

11. Reboot your machine for changes to take effect.

12. Open command prompt/terminal and type `spark-shell` . It will start Spark 3.0.1 successfully

+ Add a custom footer

▾ **Pages**  3

Find a Page...

**Home**

**Installation**

**Using a GitHub hosted dependency**

+ Add a custom sidebar

## Clone this wiki locally

https://github.com/ey-advisory-technology-testing/etl-step-definitions.wiki.git