CSCI316 Software and Programming Packages Installation Guide to install and Run PySpark in Jupyter Notebook on Windows (Adopted from instructions provided by Dr Guoxin Su)

(SJ: 25 June 2021)

Installing and Integrating PySpark with Jupyter Notebook on Windows

**A. Downloading PySpark, Anaconda, 7-Zip, and Java JDK.**

Before setting up PySpark, we need to have the following packages installed in the System:

- **Download PySpark package**
    i.   https://spark.apache.org/downloads.html



    ii.  For the package type, please choose **'Pre-built for Apache Hadoop 2.7'**
    iii. Click the label **'spark-3.1.2-bin-hadoop2.7.tgz'** to start downloading the package.
    iv.  In the next screen, choose the mirror site that you want to download the package from, and double-click on the link to start downloading.



    v.   Save the downloaded package to a working directory.

- **Installing 7-zip**
    i.   If you already have 7-zip installed on your Windows, you can skip this step, otherwise, download the 7-zip installer and install the 7-zip application.
    ii.  https://www.7-zip.org/download.html

- Choose the version of installer that fits your system and click the 'Download' link to start downloading the installer. For example, for my system, I choose the **64-bit x64** version.



- Once the downloading is completed, and the installer is checked, execute the installer to install the 7-Zip application. You can follow all the default setting unless you want to install the 7-Zip differently.

- **Installing Java JDK**
  i. If your system has no Java installed or the Java version is 7.x or less, please download and install Java from Oracle: https://www.oracle.com/java/technologies/javase/javase-jdk8-downloads.html
  ii. For my system, I downloaded and install Java SE Development Kit (Java JDK 8) (Note: So far this version works, I have tried installing the latest JDK and it did not work.)

✕

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

☑ I reviewed and accept the Oracle Technology Network License Agreement for Oracle Java SE
Required.

*You will be redirected to the login screen in order to download the file.*

Download jdk-8u291-windows-x64.exe ⬇

o   Save to a directory of your choice.



o   Next, you can proceed to install the Java JDK, but please do not use the default setting from the installer.

o   Installing Java JDK:

i.   Create a new folder in your local drive; it can be in 'C:\' or 'D:\'. For me, I install it in my 'C:\' drive, and name the folder 'Java'. Note: Please use a name that does not have a space in between.

ii.   Execute your java JDK installer. From the directory where you save the Java JDK installer, **double-click** on the installer to start the installation.

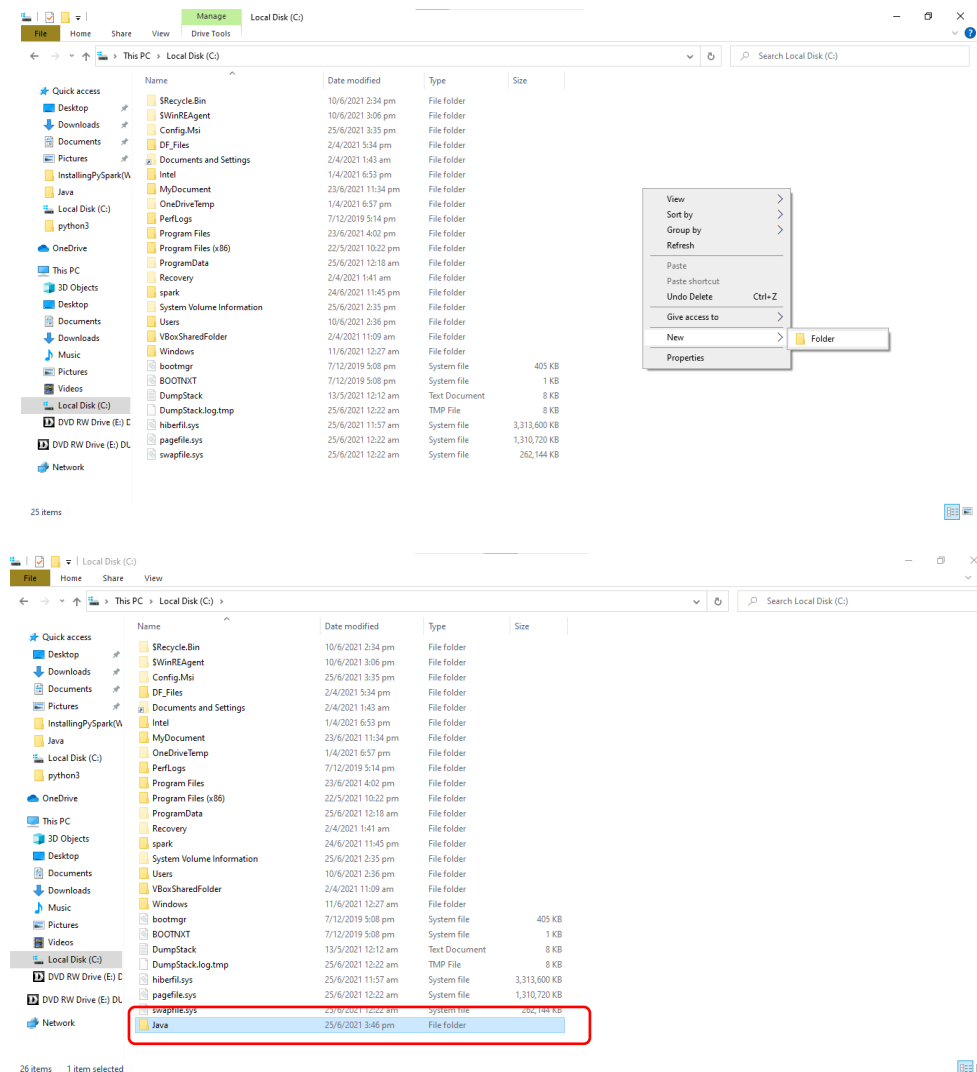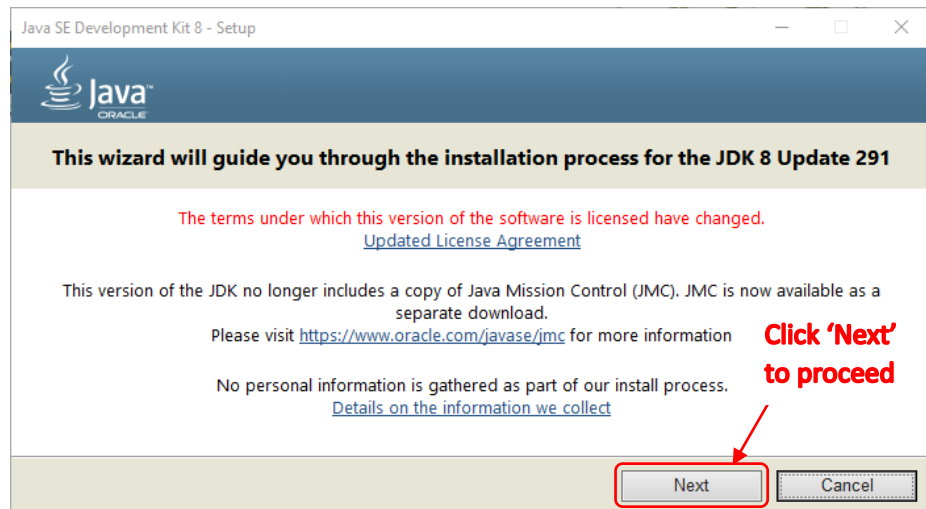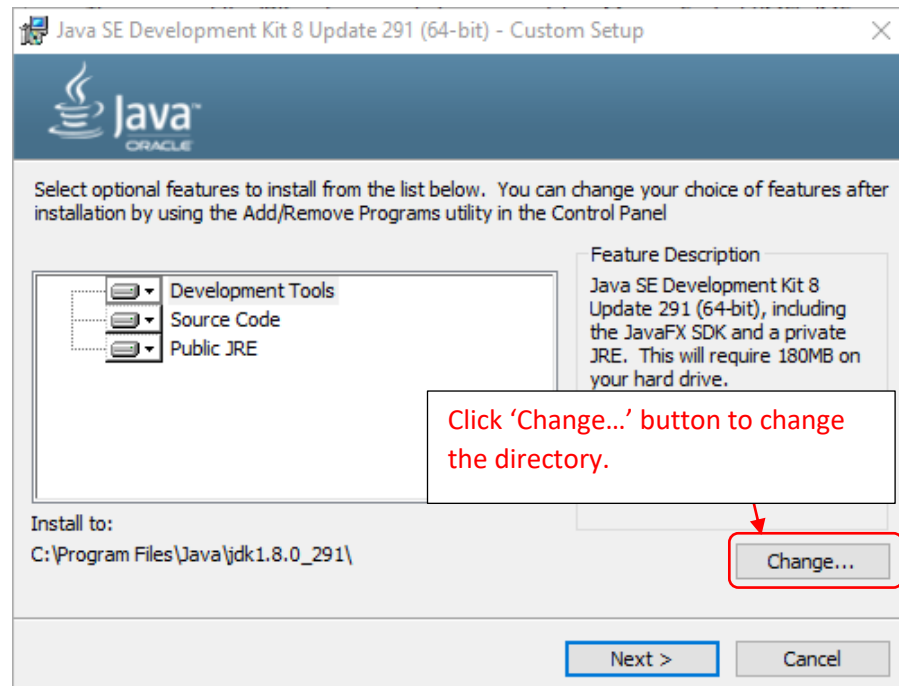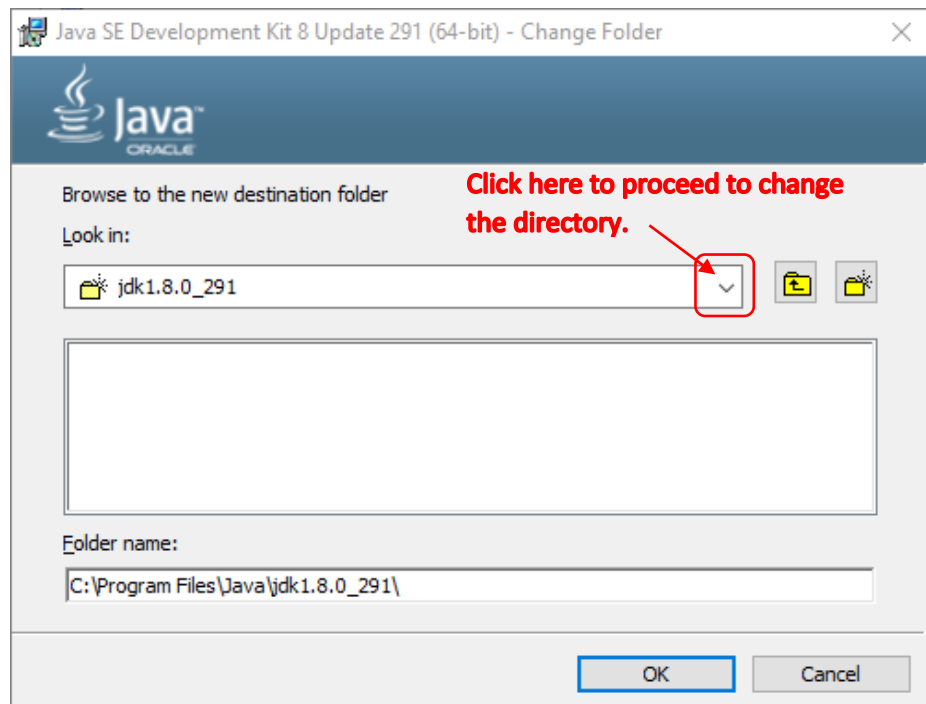| Name | Date | Type | Size | Tags |
|------|------|------|------|------|
| jdk-8u291-windows-x64 | 25/6/2021 4:39 pm | Application | 172,731 KB | |
| pipInstallFindSpark | 23/6/2021 4:50 pm | PNG File | 17 KB | |

Quick access
Desktop
Downloads

iii.   In the pop-up menu, click 'Next' to proceed.

Java SE Development Kit 8 - Setup

**This wizard will guide you through the installation process for the JDK 8 Update 291**

The terms under which this version of the software is licensed have changed.
Updated License Agreement

This version of the JDK no longer includes a copy of Java Mission Control (JMC). JMC is now available as a separate download.
Please visit https://www.oracle.com/javase/jmc for more information

No personal information is gathered as part of our install process.
Details on the information we collect

Click 'Next' to proceed

Next          Cancel

iv.   Click the 'Change...' button to change the directory where you intend to install your Java JDK.

Java SE Development Kit 8 Update 291 (64-bit) - Custom Setup

Select optional features to install from the list below. You can change your choice of features after installation by using the Add/Remove Programs utility in the Control Panel

Development Tools
Source Code
Public JRE

Feature Description

Java SE Development Kit 8 Update 291 (64-bit), including the JavaFX SDK and a private JRE. This will require 180MB on your hard drive.

Click 'Change...' button to change the directory.

Install to:
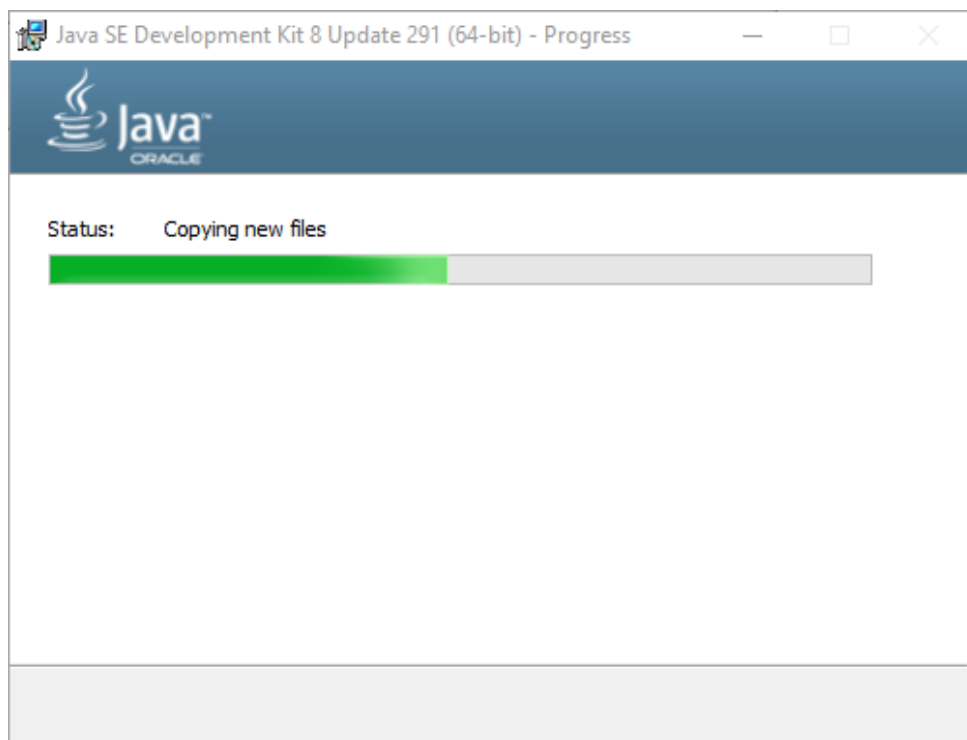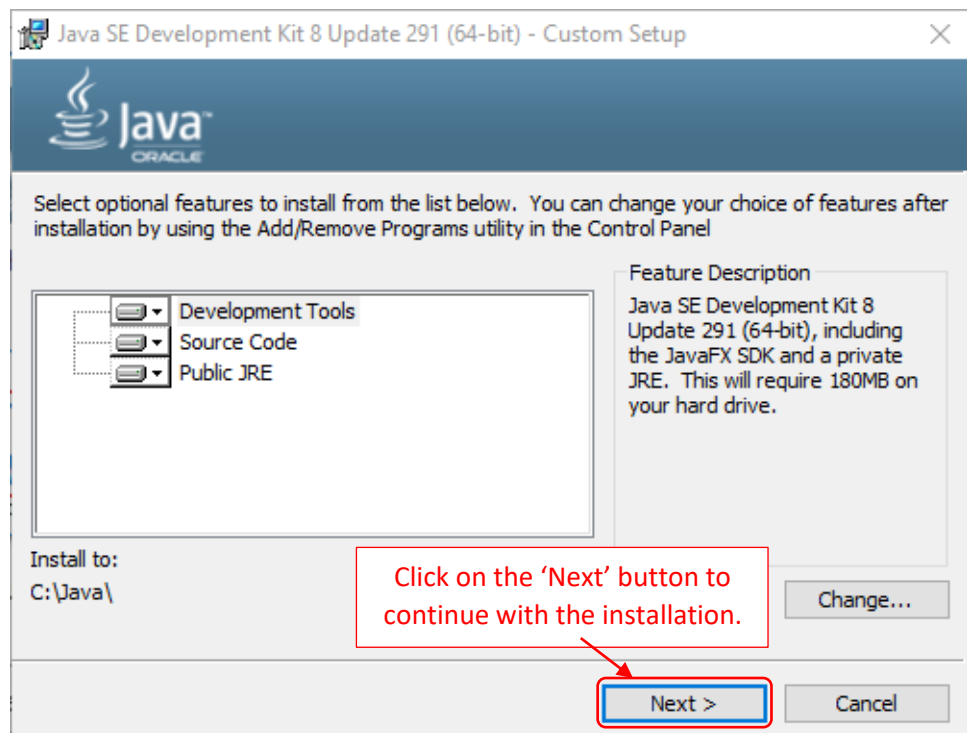C:\Program Files\Java\jdk1.8.0_291\            Change...

Next >          Cancel

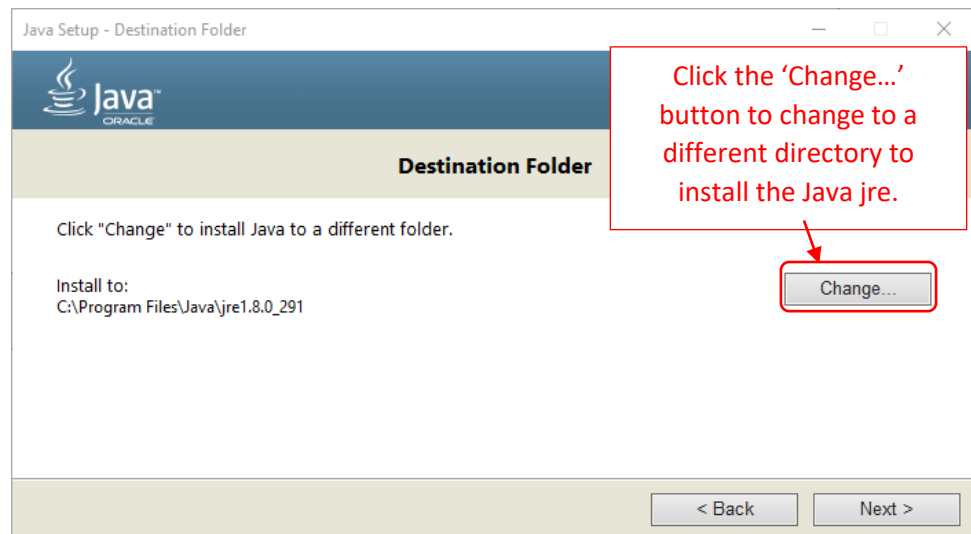v.  In the next window, click on the ⌄ icon to proceed to change the directory.



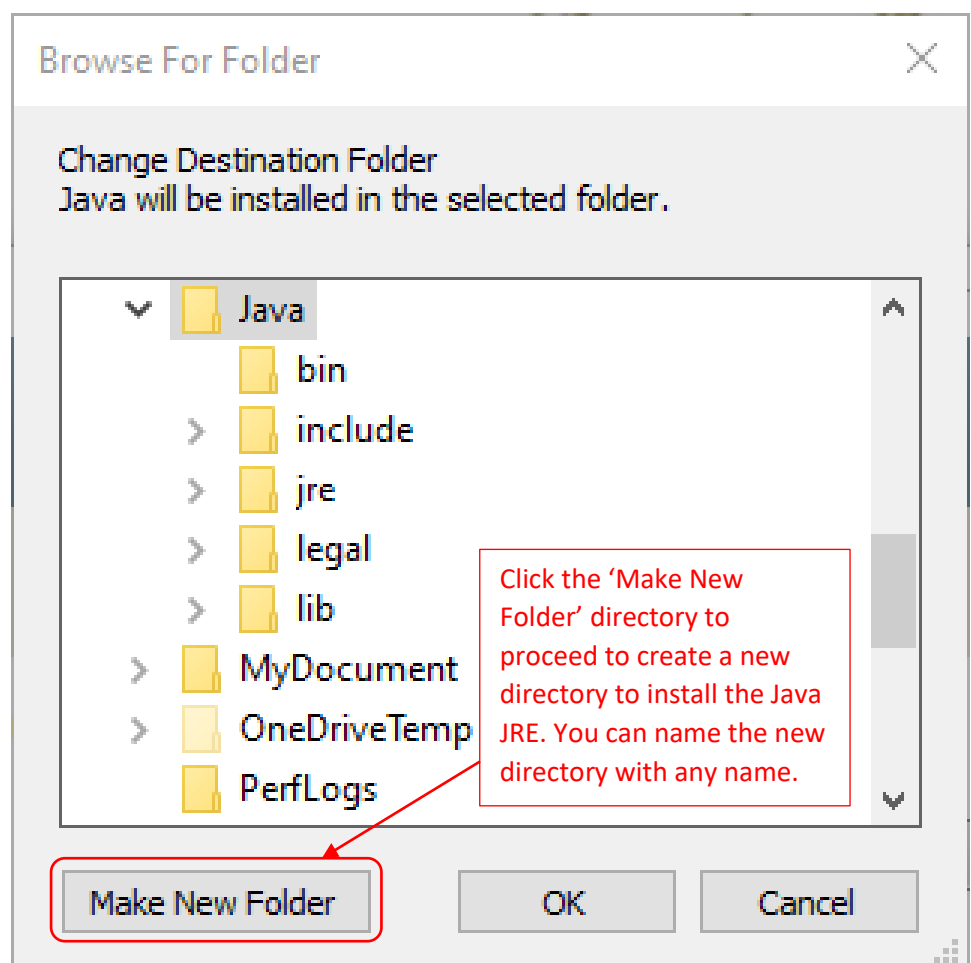vi. Change to the directory where you want to install you Java JDK and click the 'OK' button.

vii.    Click on the 'Next' button to continue the installation.

viii.    In the next window, similarly, click the 'Change…' button to change to a different
        directory to install the Java jre.
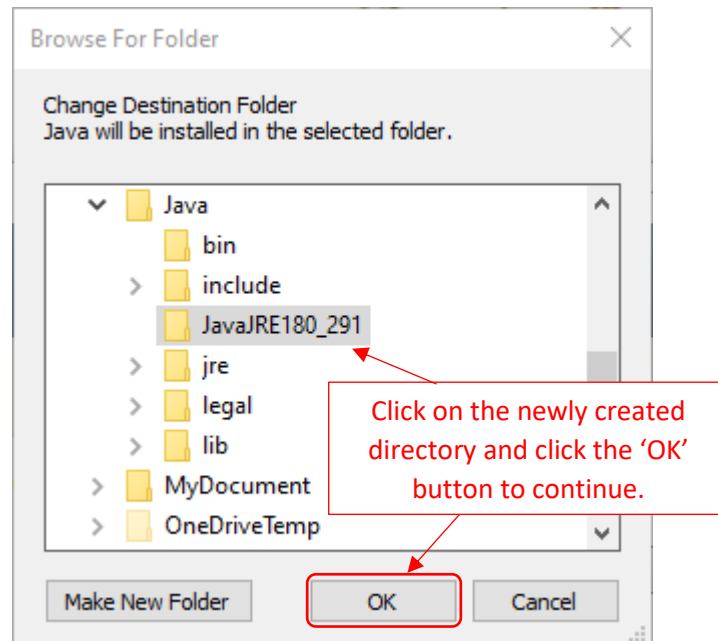


ix.    Change to the directory where you want to install you Java JRE. You may want to
       create (make) a new directory for that.
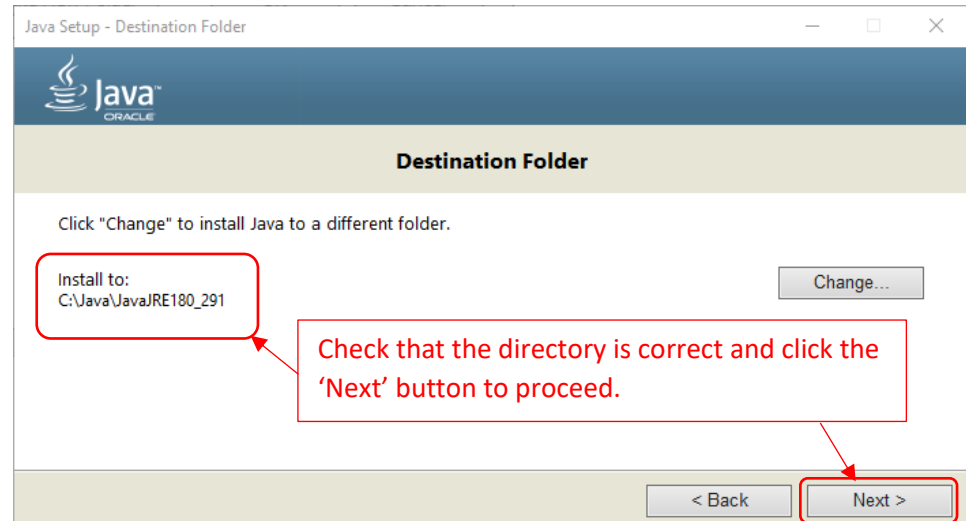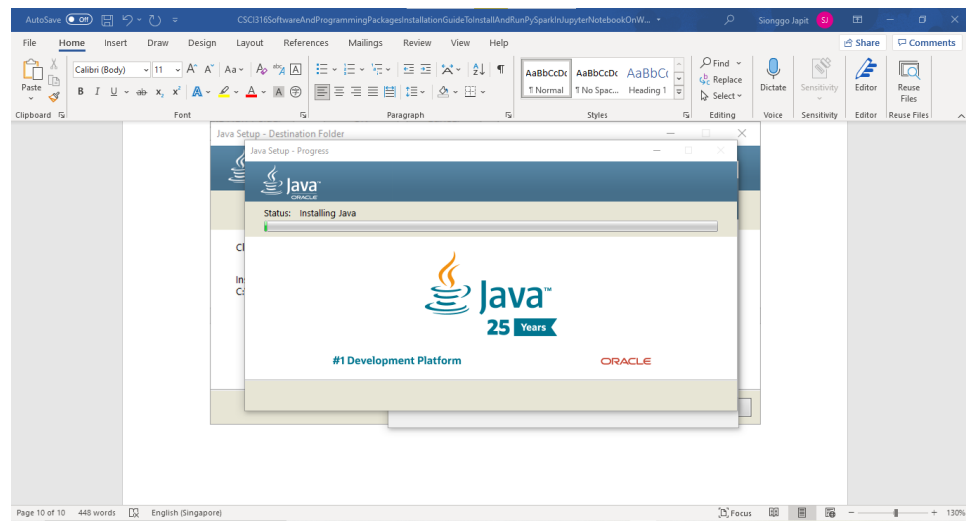
x.   Click on the newly created directory (in my example, JavaJRE180_291) and click the 'OK' button to continue.



xi.   Check that the directory is correct and click the 'Next' button to proceed with the installation.

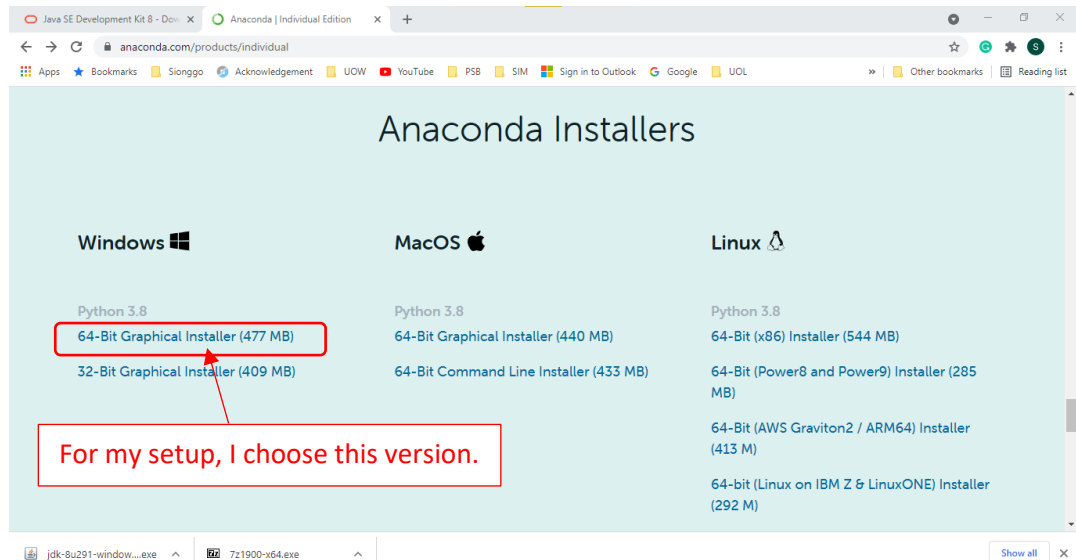xii.    Wait for the installation to complete….
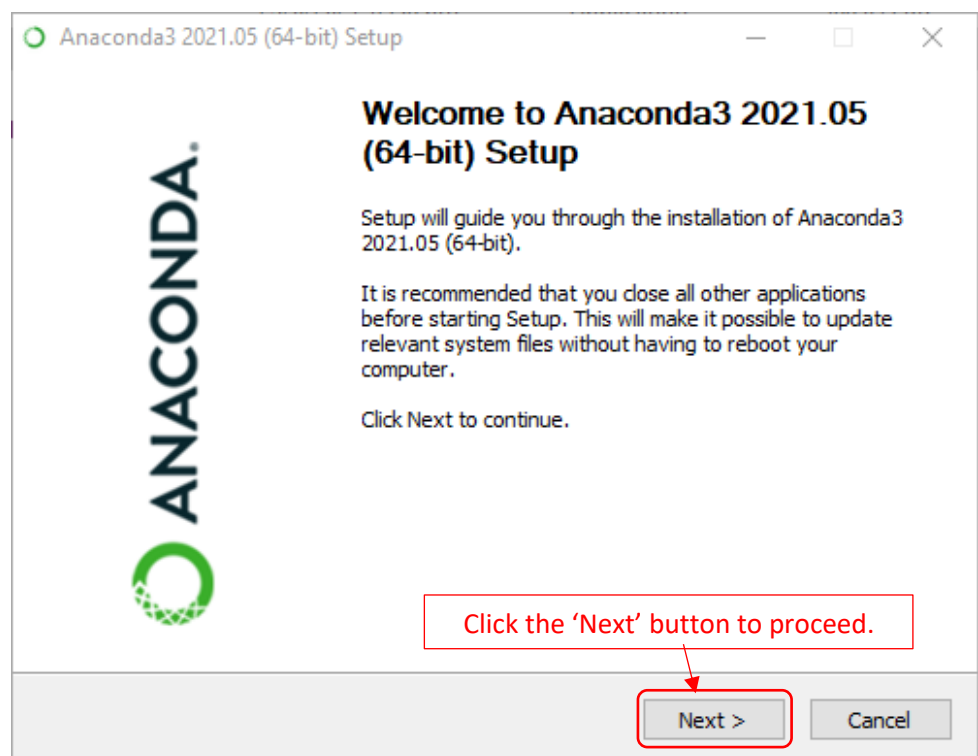


xiii.   Click the 'Close' button when done.
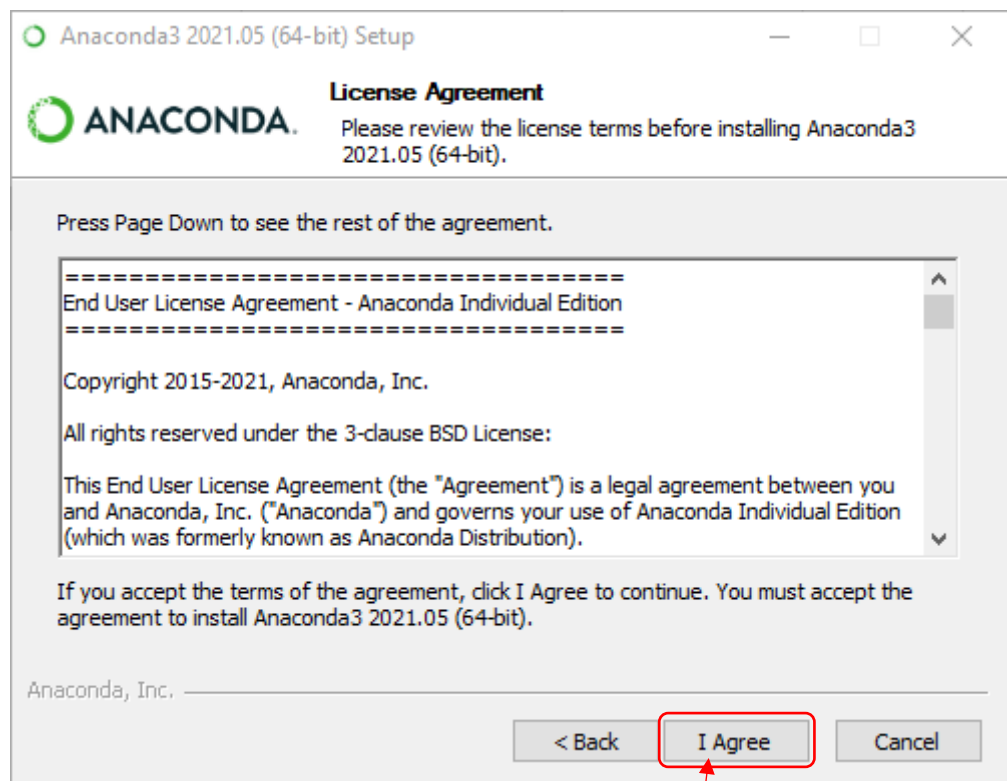
- **Installing Python and Jupyter Notebook**
    i. You can get both Python and Jupyter by installing the Python 3.x version of Anaconda distribution.
    ii. https://www.anaconda.com/products/individual
    iii. Scroll all the way down to the bottom of the page and choose the version of the installer that suit your system to start the download.



    iv. Click the 'Download' button to start the download.
    v. Installing Anaconda:
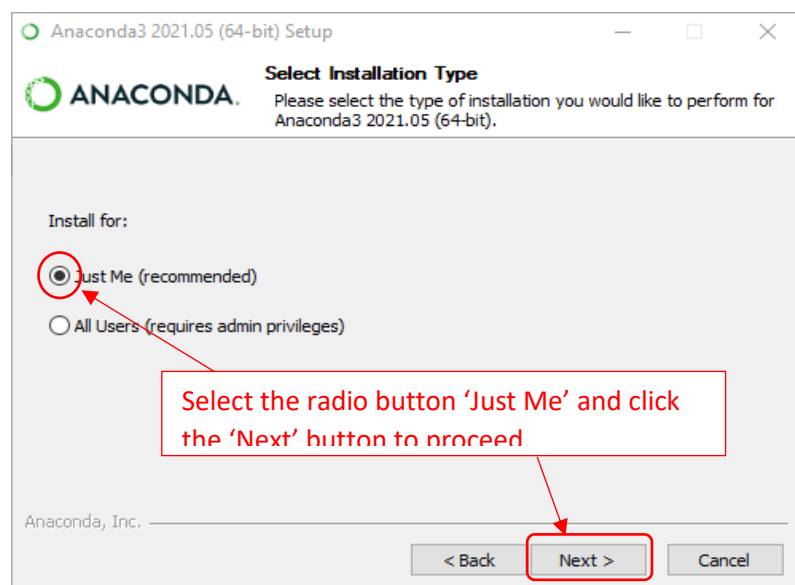        o After the downloading is complete, you can proceed to install the Anaconda distribution.

o   Read the 'License Agreement' and when finished reading, click the 'I Agree' button to proceed with the installation if you agree with the terms and condition specified in the License Agreement. 😊.
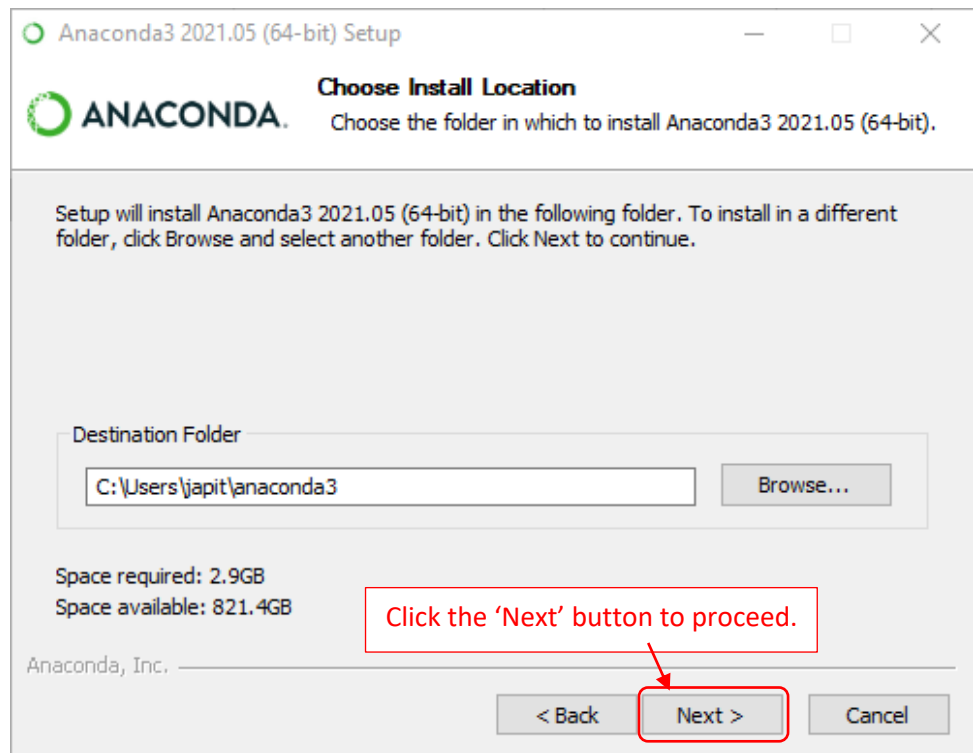


Click the 'I Agree' button to proceed with the installation if you agree with the terms and condition specified in the License Agreement.
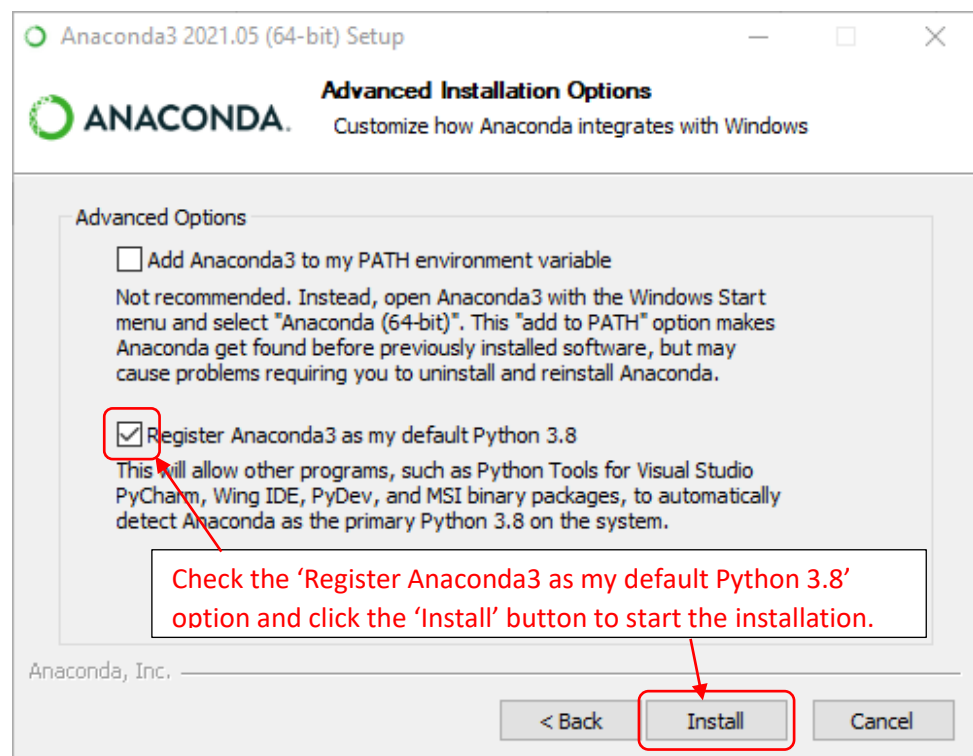
o   Select the radio button 'Just Me' and click the 'Next' button to proceed.



Select the radio button 'Just Me' and click the 'Next' button to proceed.

- You can accept the proposed destination folder. You can also change to a different directory if you want to. Click the 'Next' button to proceed.



Click the 'Next' button to proceed.

- Check the 'Register Anaconda3 as my default Python 3.8' option and click the 'Install' button to start the installation. This installation will take a while. Just relax….



Check the 'Register Anaconda3 as my default Python 3.8' option and click the 'Install' button to start the installation.
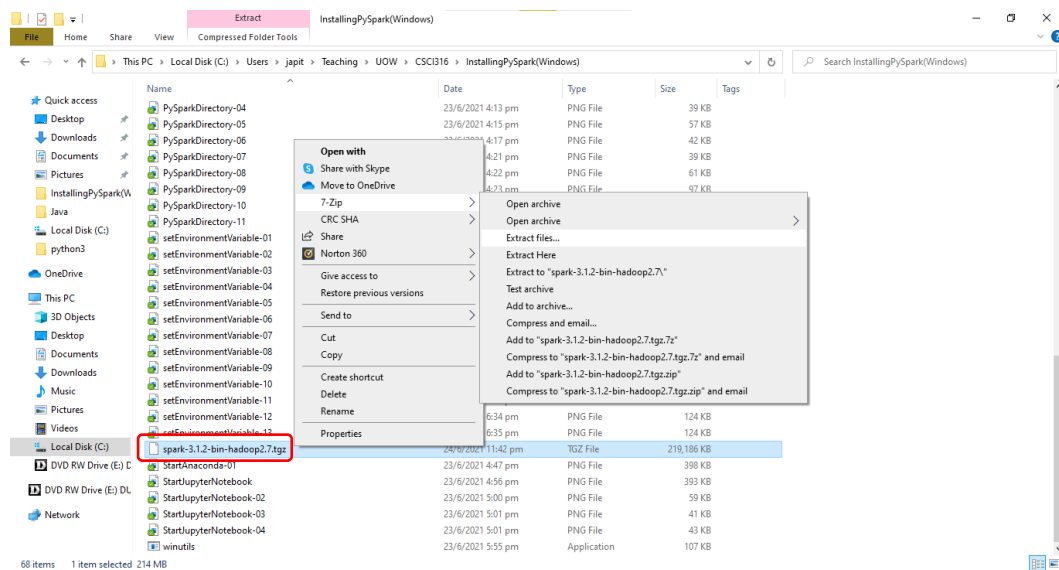
**B. Installing PySpark**

- While waiting for the Anaconda to finish the installation, we can proceed to un-zip the PySpark package for the PySpark setup.

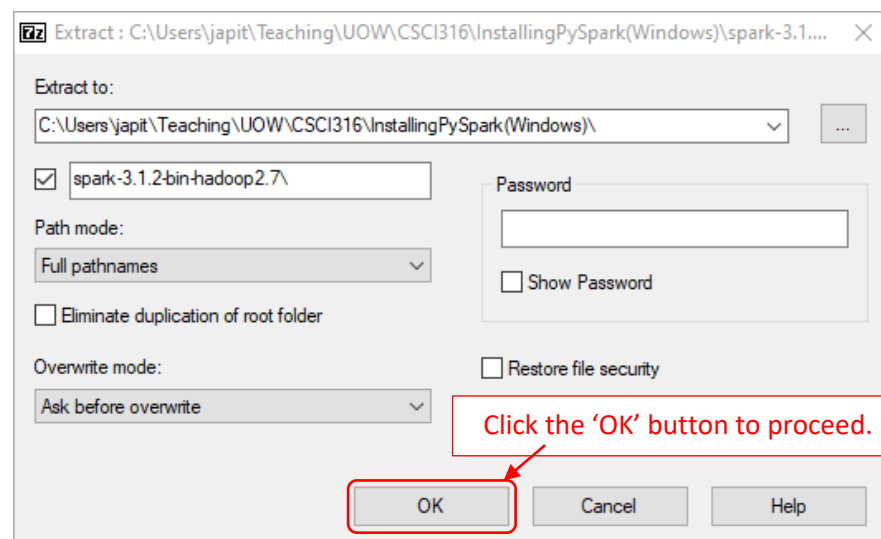  i. Unpack the spark package 'spark-3.1.2-bin-hadoop2.7.tgz' that was downloaded earlier. I save the package in my drive
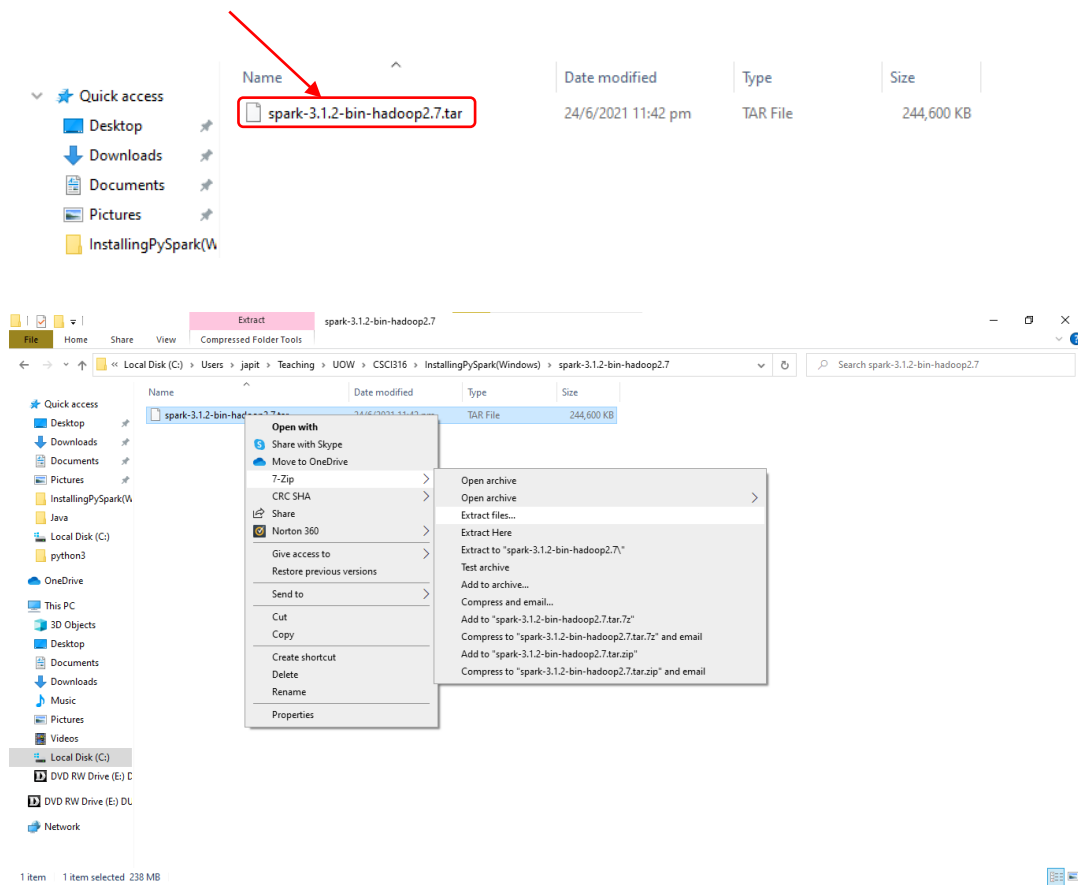  C:\Users\japit\Teaching\UOW\CSCI316\InstallingPySpark(Windows)\ spark-3.1.2-bin-hadoop2.7.tgz.

  This package is rather special. It has been packed (zipped) two times. Hence, you need to unpack (unzip) the package two times.
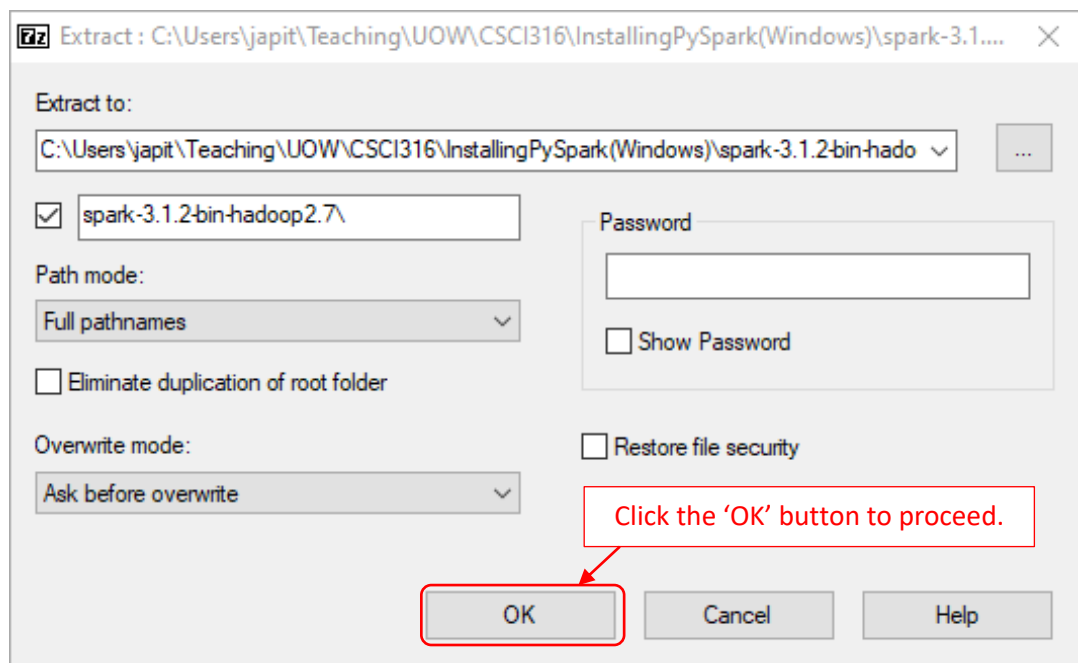


Click the 'OK' button to proceed.
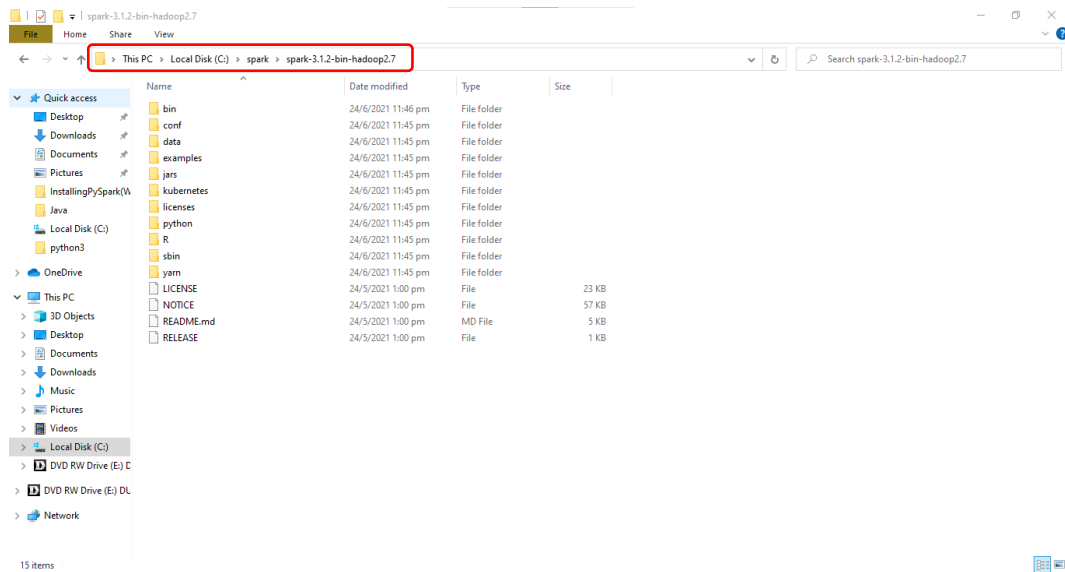


Click the 'OK' button to proceed.

After the unpacking, navigate to the folder of the unpacked file, and do another unpacking process.





Click the 'OK' button to proceed with the unpacking process.

ii.  Put the unpacked package to a directory of your choice. I put mine under C:\spark\spark-3.1.2-bin-hadoop2.7.



iii.  'winutils' is a collection of useful TCL commands that access some part of the Win32 API. This enables the user to use Windows specific services. Winutils is required when installing Hadoop on Windows environment. Winutils can be downloaded from Steve Loughran's GitHub repo.

- https://github.com/steveloughran/winutils/
- Go to the corresponding Hadoop version in the Spark distribution and find winutils.exe under /bin. Note that when we download the PySpark package in the earlier step, we choose the 'Pre-built for Apache Hadoop 2.7'.
- Click the link 'hadoop-2.7.1' to navigate to the next page.

steveloughran add 2.6.4 and 2.7.1 windows binaries    ☐1   7665f01 on Feb 13, 2016   ⏱ History

..

📁 bin     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 README.md     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

**README.md**

These are actually the binaries off HDP 2.3.0; they should be interchangeble with the ASF 2.7.1 release

(usual disclaimers etc: if you want the artifacts direct you can download the whole install from hortonworks.com.

This is just what ended up in my hdp/bin dir after the installation -though I have deleted the pdb files.

- o Click the link 'bin' to navigate to the next page.

📄 snappy-stubs-internal.obj     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 snappy.dll     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 snappy.dll.intermediate.manifest     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 snappy.exp     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 snappy.lastbuildstate     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 snappy.lib     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 snappy.obj     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 snappy.write.1.tlog     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 timelineserver.exe     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 timelineserver.xml     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 winutils.exe     add 2.6.4 and 2.7.1 windows binaries     5 years ago

📄 yarn     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

📄 yarn.cmd     Add HDP2.3.0/Hadoop 2.7.1 windows binary artifacts     6 years ago

- o Click the link 'winutils.exe' to navigate to the next page.

ℙ master ▾    winutils / hadoop-2.7.1 / bin / winutils.exe      Go to file   ⋯

steveloughran add 2.6.4 and 2.7.1 windows binaries     Latest commit 7665f01 on Feb 13, 2016   ⏱ History
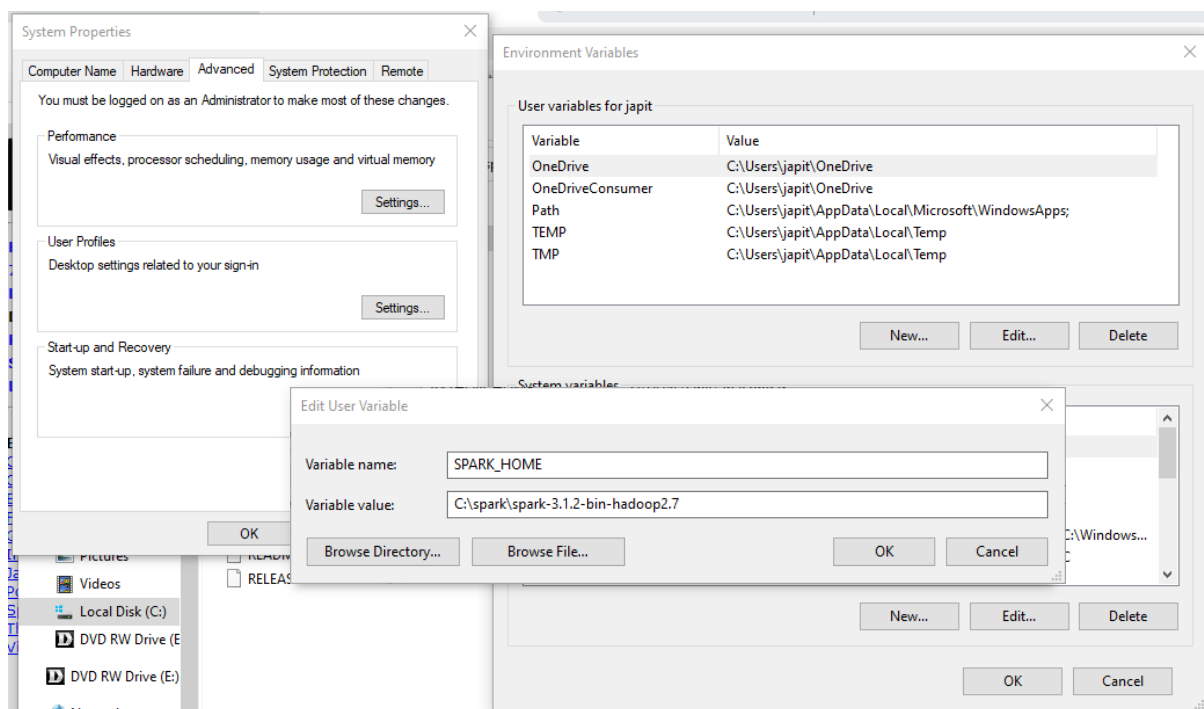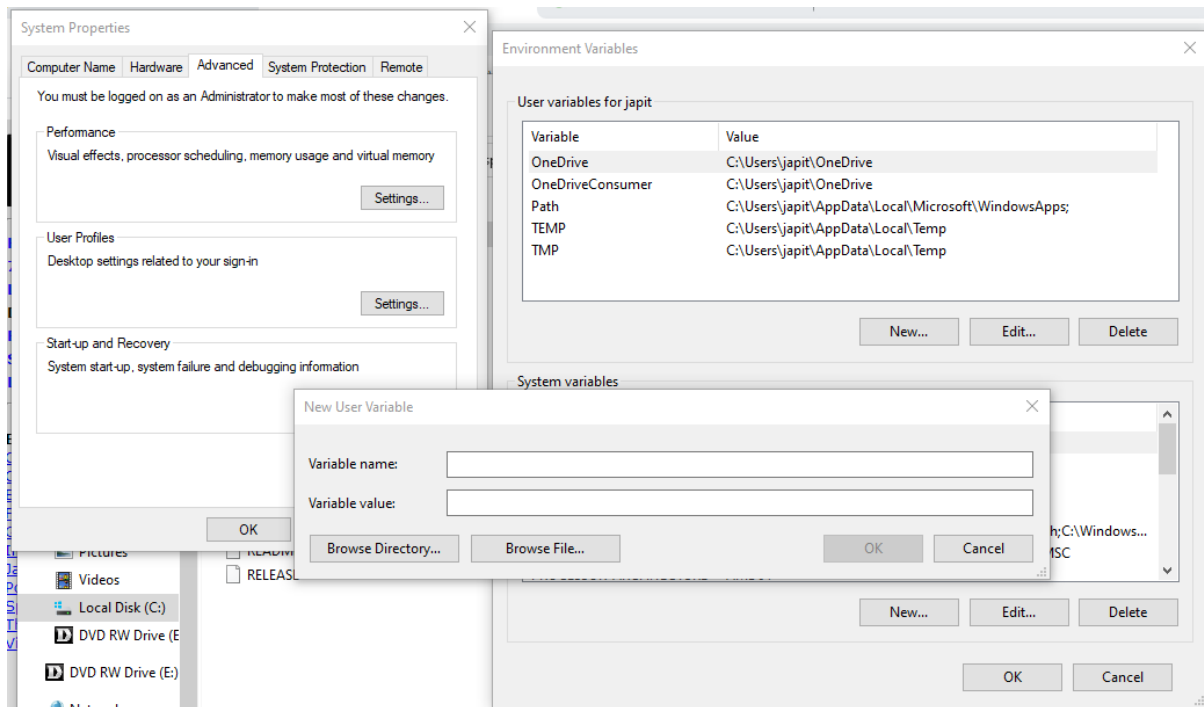
👥 1 contributor

107 KB          Download   🖥 🗑

View raw

- o Click the button 'Download' to start the download.
- o The downloaded 'winutils.exe' need to be placed in the bin directory of the spark folders 'C:\spark\spark-3.1.2-bin-hadoop2.7.'.

o The 'winutils.exe' is placed here.



- **Setting of user's environment variable**
- In order for the system to know where the various components (packages) are installed, we need to add the path (directory) of each component to the user's environment variables. These variables let Windows find where the files are when we start the PySpark kernel.
- You can find the environment variable settings by typing "environment…" in the search box. When the System Properties windows is shown, click on the 'Environment Variables…' button to proceed.

- Do the same for the rest of the components shown below. Note that the values shown here are according to my setup (the directories where I place the packages.) You need to follow according to your setup.

| Name | Values |
| --- | --- |
| SPARK_HOME | C:\spark\spark-3.1.2-bin-hadoop2.7 |
| HADOOP_HOME | C:\spark\spark-3.1.2-bin-hadoop2.7 |
| PYSPARK_DRIVER_PYTHON | Jupyter |
| PYSPARK_DRIVER_PYTHON_OPTS | notebook |

- You also need to set the path where you unpacked your spark to the system's path variable. In the same environment variable settings window, look for the 'Path' variable and click the 'Edit…' button to proceed.



- In the next 'Edit environment variables' pop-up window, click the 'New' button, and enter/add the path where you unpacked your spark packages follows by a '\bin'. For example, I unpacked my spark in C:\spark\spark-3.1.2-bin-hadoop2.7. In this path, there is one directory where the binaries are stored. I will add that directory to the path, hence, my path will be 'C:\spark\spark-3.1.2-bin-hadoop2.7\bin'.

Click the 'OK' button to save your entry.

- By now, I suppose the installation of Anaconda may have been completed. Check that the Anaconda installation is indeed completed.
- Next, we need install findspark, pyspark, and other required libraries for CSCI316. We will use Anaconda to do so:
    i.  Open the Anaconda prompt.



ii. Installing findspark
   ➢ At the prompt, type **'pip install findspark'** followed with a 'return' key.



Since I have installed my 'findpark' earlier, the system will indicate that the requirement has already been satisfied, otherwise, the system will install the 'findpark' for you.

iii.     Installing pyspark
➢ At the prompt, type '**pip install pyspark**' followed with a 'return' key.



Similarly, since I have installed my 'pyspark' earlier, the system will indicate that the requirement has already been satisfied, otherwise, the system will install the 'pypark' for you.

iv.     Installing tensorflow
➢ At the prompt, type '**pip install –upgrade tensorflow**' followed with a 'return' key.

v. Installing pyarrow

➢ At the prompt, type '**conda install -c conda-forge pyarrow**' followed with a 'return' key.

```
Anaconda Prompt (anaconda3)                                                — □ ×
zope                          1.0                          py38_1
zope.event                    4.5.0                        py38_0
zope.interface                5.3.0                 py38h2bbff1b_0
zstd                          1.4.5                        h04227a9_0

(base) C:\Users\japit>conda install -c conda-forge pyarrow
Collecting package metadata (current_repodata.json): done
Solving environment: -
The environment is inconsistent, please check the package plan carefully
The following packages are causing the inconsistency:

  - defaults/win-64::anaconda==2021.05=py38_0
  - defaults/win-64::astropy==4.2.1=py38h2bbff1b_1
  - defaults/win-64::bkcharts==0.2=py38_0
  - defaults/win-64::bokeh==2.3.2=py38haa95532_0
  - defaults/win-64::bottleneck==1.3.2=py38h2a96729_1
  - defaults/noarch::dask==2021.4.0=pyhd3eb1b0_0
  - defaults/win-64::imagecodecs==2021.3.31=py38h5da4933_0
  - defaults/noarch::imageio==2.9.0=pyhd3eb1b0_0
  - defaults/win-64::matplotlib==3.3.4=py38haa95532_0
  - defaults/win-64::matplotlib-base==3.3.4=py38h49ac443_0
  - defaults/win-64::mkl_fft==1.3.0=py38h277e83a_2
  - defaults/win-64::mkl_random==1.2.1=py38hf11a4ad_2
  - defaults/win-64::numba==0.53.1=py38hf11a4ad_0
  - defaults/win-64::numexpr==2.7.3=py38hb80d3ca_1
  - defaults/win-64::numpy==1.20.1=py38h34a8a5c_0
  - defaults/win-64::pandas==1.2.4=py38hd77b12b_0
  - defaults/win-64::patsy==0.5.1=py38_0
  - defaults/win-64::pyerfa==1.7.3=py38h2bbff1b_0
  - defaults/win-64::pytables==3.6.1=py38ha5be198_0
  - defaults/win-64::pywavelets==1.1.1=py38he774522_2
```

```
Anaconda Prompt (anaconda3)                                                — □ ×
python_abi-3.8         | 4 KB    | ################################################## | 100%
boost-cpp-1.69.0       | 31.9 MB | ################################################## | 100%
aws-sdk-cpp-1.8.185    | 2.5 MB  | ################################################## | 100%
arrow-cpp-4.0.1        | 4.3 MB  | ################################################## | 100%
aws-c-event-stream-0   | 26 KB   | ################################################## | 100%
glog-0.5.0             | 90 KB   | ################################################## | 100%
aws-checksums-0.1.9    | 51 KB   | ################################################## | 100%
libllvm9-9.0.1         | 48 KB   | ################################################## | 100%
pyarrow-4.0.1          | 1.8 MB  | ################################################## | 100%
gmpy2-2.1.0b5          | 190 KB  | ################################################## | 100%
utf8proc-2.6.1         | 312 KB  | ################################################## | 100%
ca-certificates-2021   | 171 KB  | ################################################## | 100%
abseil-cpp-20200225.   | 1.9 MB  | ################################################## | 100%
pathtools-0.1.2        | 8 KB    | ################################################## | 100%
gflags-2.2.2           | 80 KB   | ################################################## | 100%
certifi-2021.5.30      | 142 KB  | ################################################## | 100%
libboost-1.73.0        | 20.3 MB | ################################################## | 100%
mpfr-4.0.2             | 1.9 MB  | ################################################## | 100%
anaconda-custom        | 36 KB   | ################################################## | 100%
aws-c-common-0.4.57    | 150 KB  | ################################################## | 100%
double-conversion-3.   | 101 KB  | ################################################## | 100%
uriparser-0.9.3        | 47 KB   | ################################################## | 100%
openssl-1.1.1k         | 5.7 MB  | ################################################## | 100%
mpir-3.0.0             | 3.0 MB  | ################################################## | 100%
re2-2021.06.01         | 467 KB  | ################################################## | 100%
mpc-1.1.0              | 322 KB  | ################################################## | 100%
Preparing transaction: done
Verifying transaction: done
Executing transaction: done

(base) C:\Users\japit>conda list
```

vi. Installing cairocffi

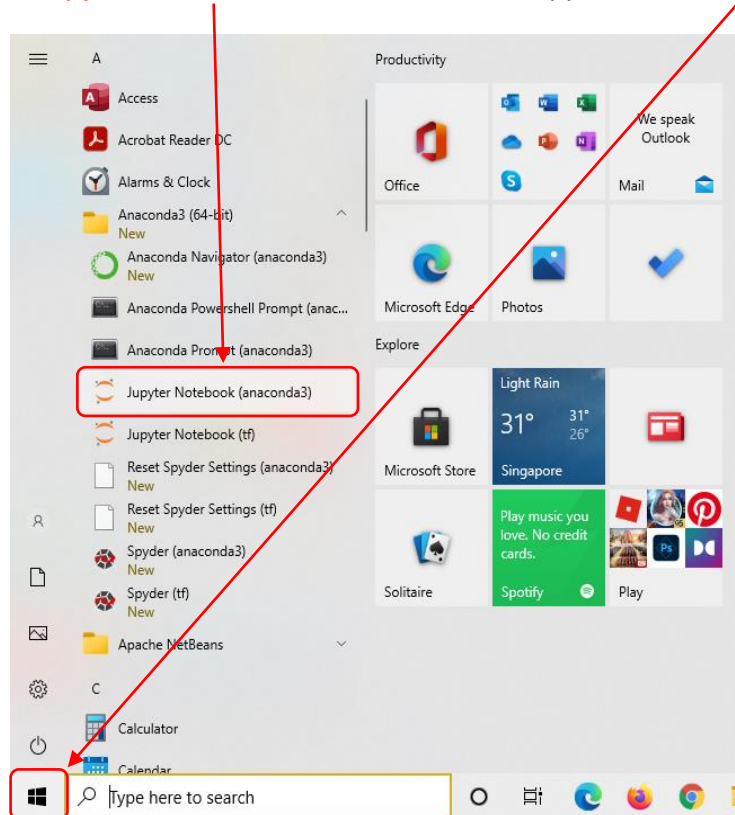➤ At the prompt, type '**conda install -c conda-forge/label/cf202003 cairocffi**' followed with a 'return' key.



vii. We can now check/verify that all the required software/libraries are installed. To do that at the prompt, type '**conda list**' followed with a 'return' key.

**C. Running PySpark in Jupyter Notebook**

    i.    To start Jupyter Notebook, open Jupyter Notebook via the windows '**Start**' icon. Click on the lable '**Jupyter Notebook (anaconda3)**' to start Jupyter Notebook.



The Jupyter Notebook server is started. Leave this window stays open.

A Jupyter Notebook client is open. Navigate to the desired working directory. In my setup, I have a directory named **'Documents'** created in my user's name in Windows. I will use this directory as my working directory for CSCI316.
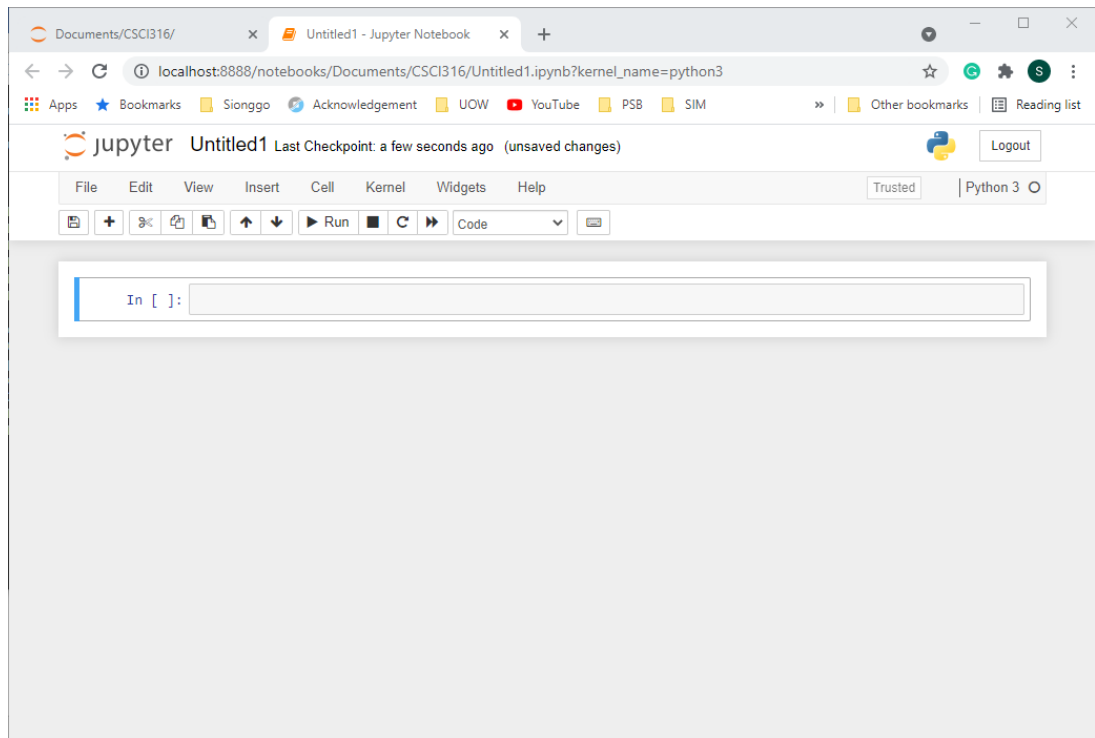


If you do not have a suitable working directory and want to create a new directory, you can click on the New ▾ icon and use the option to create a new folder. Once the new folder is created, you can rename it to your choice.

I am now in my working directory. I create a new Python Jupyter Notebook by selecting the '**New**' icon and from the drop-down option, choose '**Python 3**' option.
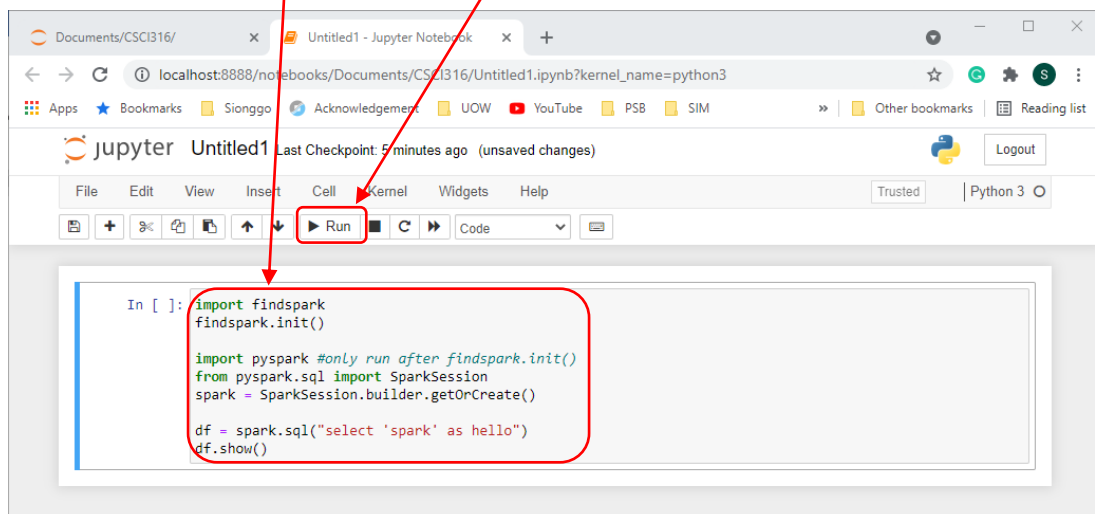
A new Jupyter Notebook node is created.



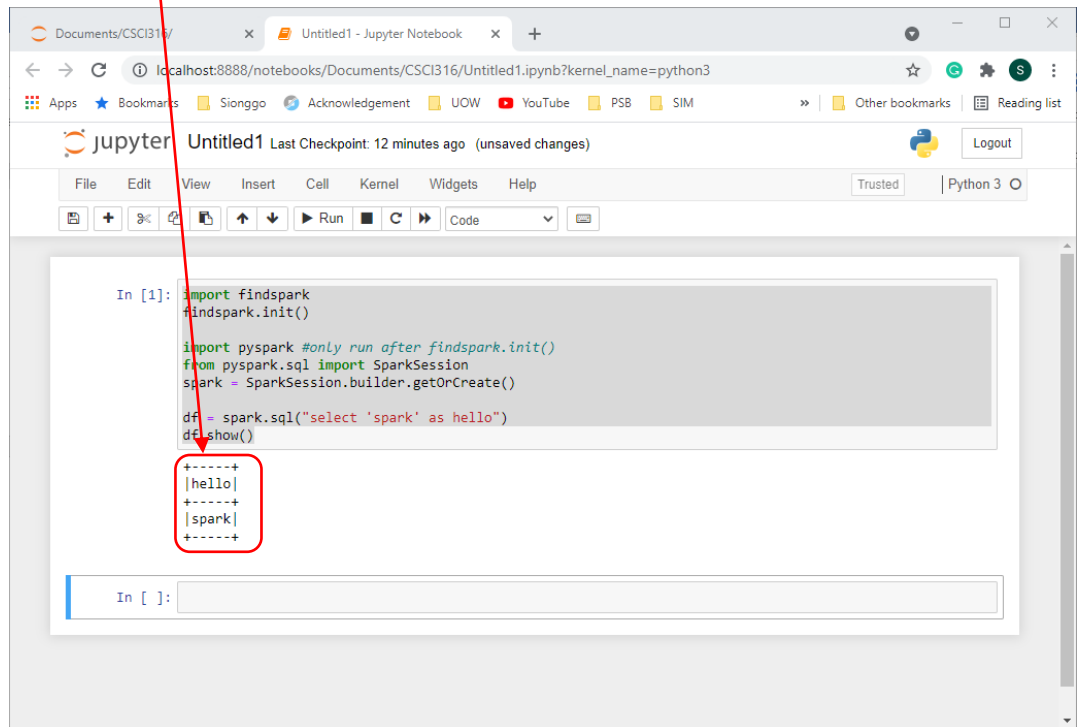Type the following **codes** in the Jupyter notebook and **run**/execute the code.

import findspark
findspark.init()

import pyspark #only run after findspark.init()
from pyspark.sql import SparkSession
spark = SparkSession.builder.getOrCreate()

df = spark.sql("select 'spark' as hello")
df.show()

If you see this, 'Congratulation' you have successfully install PySpark in your system.