

(SJ: 16 Jan 2021)

## Installing and Integrating PySpark with Jupyter Notebook on Linux

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

- Java SE Development Kit (Java jdk)
- Scala Build Tool
- Python3
- Jupyter Notebook (<https://jupyter.org/install>)
- Apache Spark 2.4 or higher (<https://spark.apache.org/downloads.html>). NOTE: select the “pre-built for Apache Hadoop 2.7 and later” installation)

Note: Both Python3 and Jupyter Notebook can be obtained by install Anaconda.

We can get all the required packages install by first update the Linux apt-get.

```
kali㉿kali:~$ sudo apt-get update
Ign:1 http://ppa.launchpad.net/deadsnakes/ppa/ubuntu hirsute InRelease
Err:2 http://ppa.launchpad.net/deadsnakes/ppa/ubuntu hirsute Release
      404  Not Found [IP: 91.189.95.85 80]
Hit:3 http://kali.cs.ntu.edu.tw/kali kali-rolling InRelease
Reading package lists... Done
E: The repository 'http://ppa.launchpad.net/deadsnakes/ppa/ubuntu hirsute Release' does not have a Release file.
N: Updating from such a repository can't be done securely, and is therefore disabled by default.
N: See apt-secure(8) manpage for repository creation and user configuration details.
kali㉿kali:~$
```

Next, we install the Java jdk and Scala. If your system has already installed with the respective packages, Linux will tell you so, otherwise, the required packages will be installed to your system.

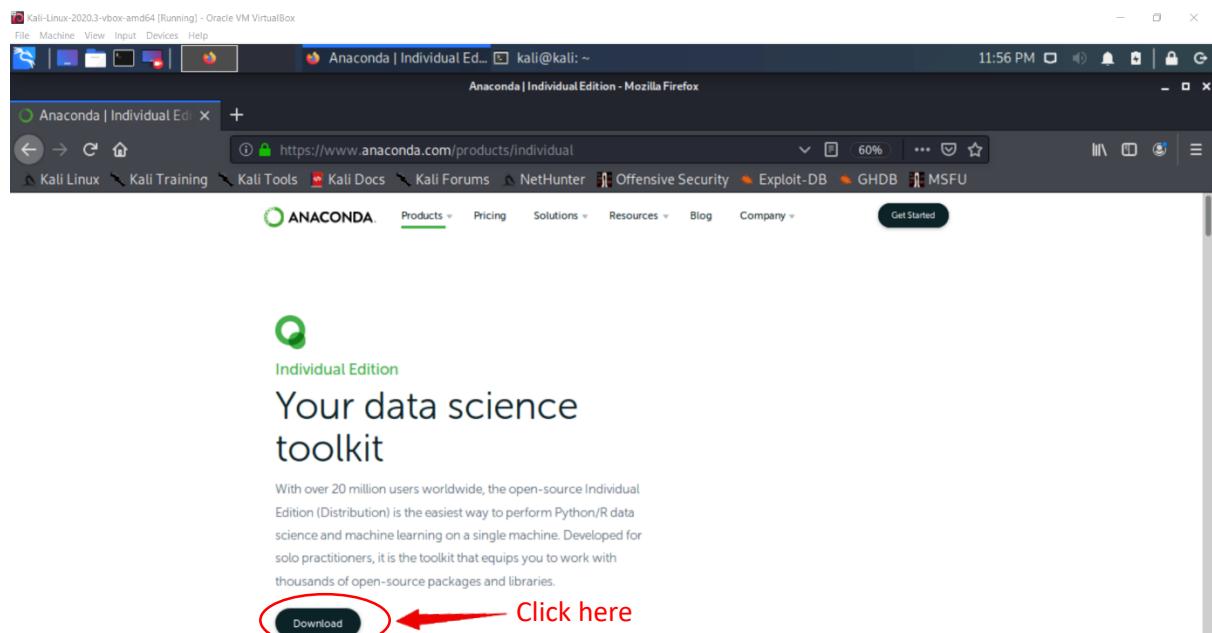
```
kali㉿kali:~$ sudo apt-get -y install default-jdk scala
[sudo] password for kali:
Reading package lists... Done
Building dependency tree
Reading state information... Done
default-jdk is already the newest version (2:1.11-72).
scala is already the newest version (2.11.12-4).
0 upgraded, 0 newly installed, 0 to remove and 1176 not upgraded.
kali㉿kali:~$
```

Upon completion of the installation, you can check the version of the respective packages.

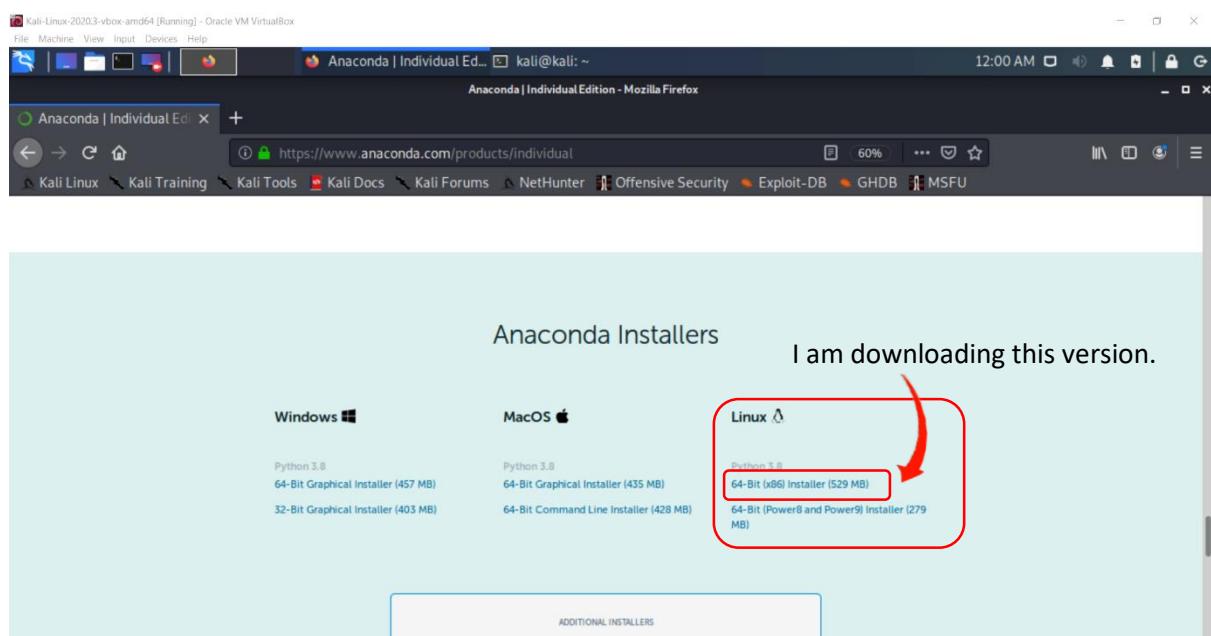
```
kali㉿kali:~$ java -version
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
openjdk version "11.0.8" 2020-07-14
OpenJDK Runtime Environment (build 11.0.8+10-post-Debian-1)
OpenJDK 64-Bit Server VM (build 11.0.8+10-post-Debian-1, mixed mode, sharing)
kali㉿kali:~$
```

```
kali㉿kali:~$ scala -version
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Scala code runner version 2.11.12 -- Copyright 2002-2017, LAMP/EPFL
kali㉿kali:~$
```

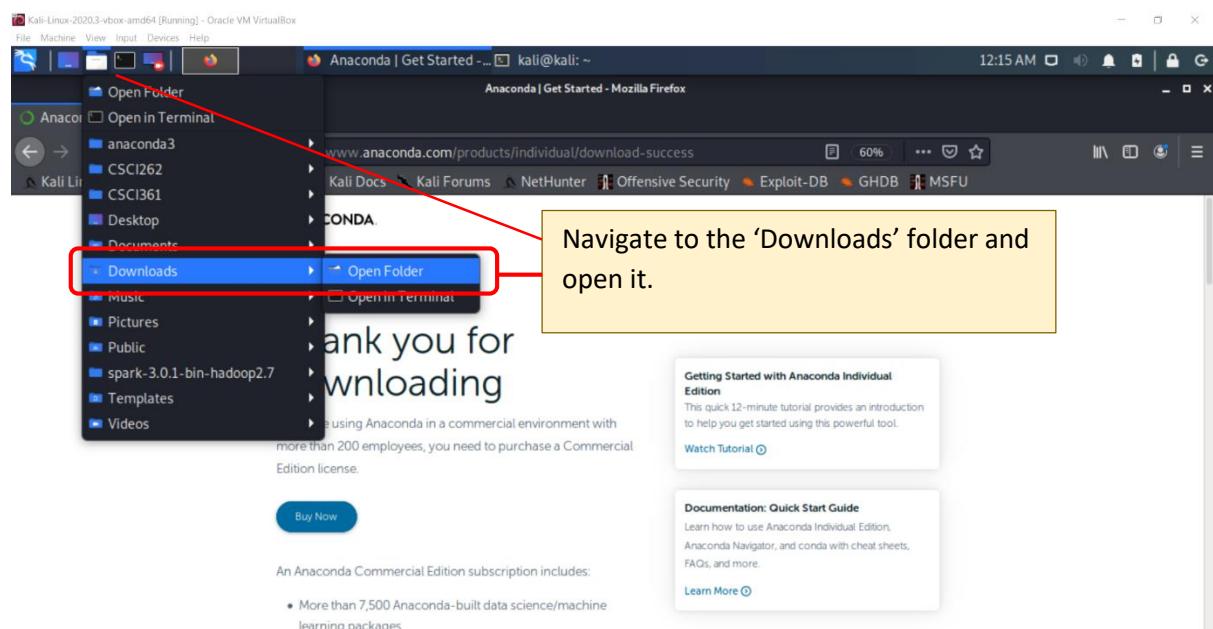
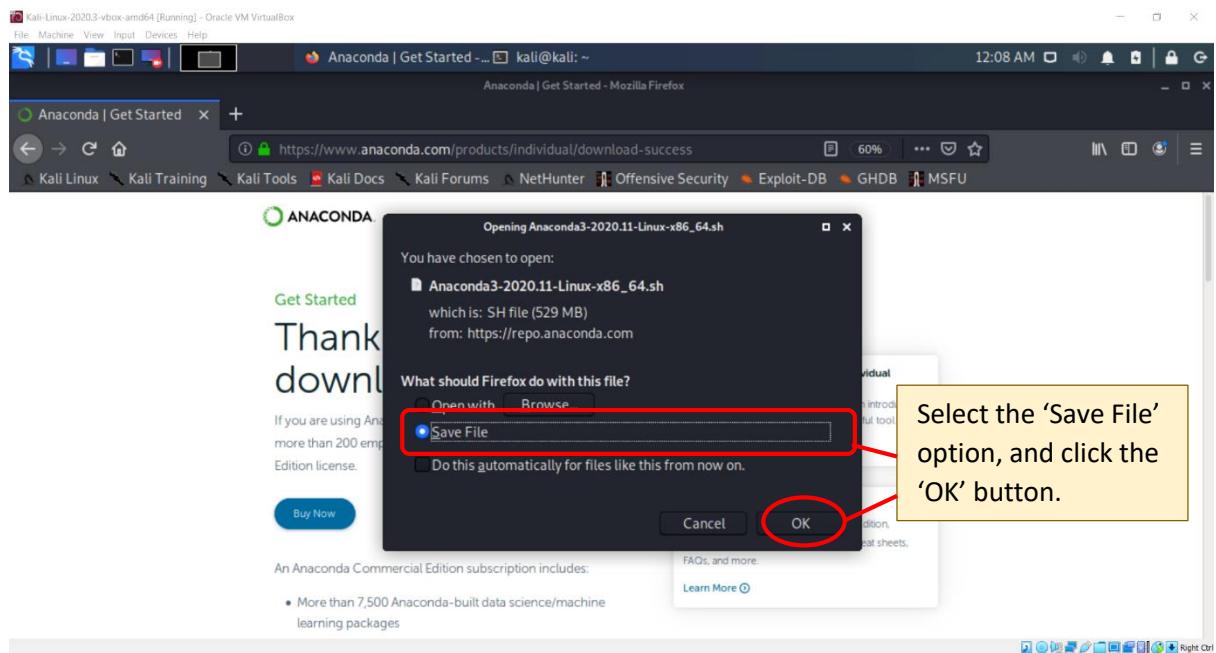
## Downloading and installing Anaconda



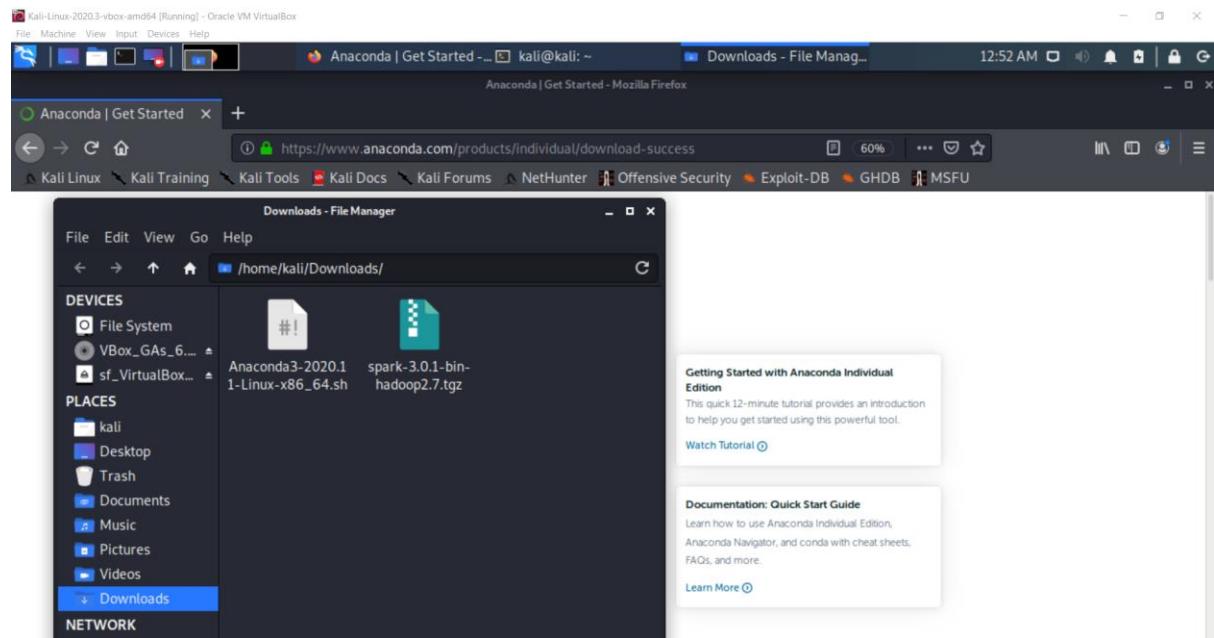
Choose the Linux version to download.



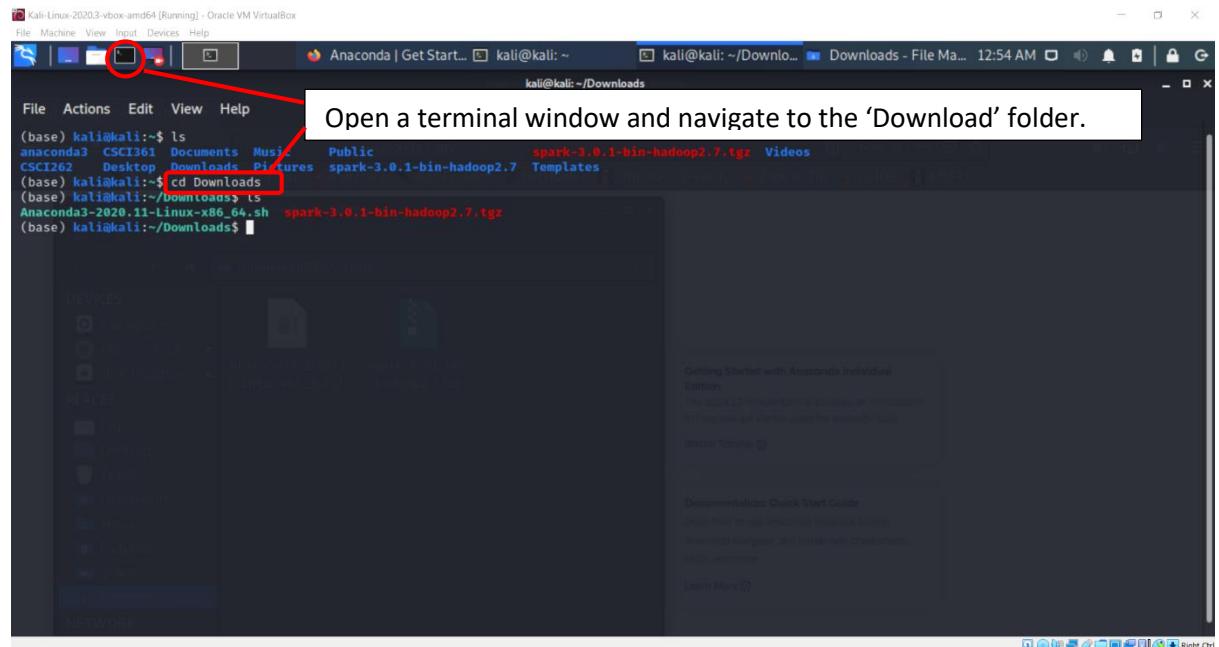
Save your download.



Check to ensure the download is successful.



Next, open a terminal window and navigate to the 'Download' folder.



Run the shell (.sh) script to install Anaconda3. Ensure you are in the directory where the installer script was downloaded. Issue the command at the terminal prompt \$ sh Anaconda3-2020.11-Linux-x86\_64.sh.

```
(base) kali㉿kali:~/Downloads$ sh Anaconda3-2020.11-Linux-x86_64.sh
[base] Kali-Linux-2020.3-vbox-amd64 [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
spark-3... Extract [Installati... kali@kali... [kali@kali... kali@kali... Download... 09:58 AM
kali@kali:~/anaconda3
File Actions Edit View Help
kali@kali:~/Downloads$ ls
Anaconda3-2020.11-Linux-x86_64.sh spark-3.0.3-bin-hadoop2.7.tgz
kali@kali:~/Downloads$ sh Anaconda3-2020.11-Linux-x86_64.sh

Welcome to Anaconda3 2020.11

In order to continue the installation process, please review the license
agreement.
Please, press ENTER to continue
>>>
End User License Agreement - Anaconda Individual Edition

Copyright 2015-2020, Anaconda, Inc.

All rights reserved under the 3-clause BSD License:

This End User License Agreement (the "Agreement") is a legal agreement between you and Anaconda, Inc. ("Anaconda") and governs your use of Anaconda Individual Edition (which was formerly known as Anaconda Distribution).

Subject to the terms of this Agreement, Anaconda hereby grants you a non-exclusive, non-transferable license to:


- Install and use the Anaconda Individual Edition (which was formerly known as Anaconda Distribution),
- Modify and create derivative works of sample source code delivered in Anaconda Individual Edition from Anaconda's repository; and
- Redistribute code files in source (if provided to you by Anaconda as source) and binary forms, with or without modification subject to the requirements set forth below.


Anaconda may, at its option, make available patches, workarounds or other updates to Anaconda Individual Edition. Unless the updates are provided with their separate governing terms, they are deemed part of Anaconda Individual Edition licensed to you as provided in this Agreement. This Agreement does not entitle you to any support for Anaconda Individual Edition.

Anaconda reserves all rights not expressly granted to you in this Agreement.
```

Accept the Licence Agreement and confirm the location where you want your Anaconda3 to be installed. For me, I just accept the proposed location.

```
Kali-Linux-2020.3-vbox-amd64 [Running] - Oracle VM VirtualBox  
File Machine View Input Devices Help  
spark-3.... Extract [Installat... kali@kali... kali@kal... kali@kal... Download... 10:01 AM  
kali@kali: ~/anaconda3  
File Actions Edit View Help  
Do you accept the license terms? [yes|no]  
[no] >>>  
Please answer 'yes' or 'no':  
>>> yes  
  
Anaconda3 will now be installed into this location:  
/home/kali/anaconda3  
  
- Press ENTER to confirm the location  
- Press CTRL-C to abort the installation  
- Or specify a different location below  
  
[~/home/kali/anaconda3] >>>  
PREFIX=/home/kali/anaconda3  
Unpacking payload ...  
Collecting package metadata (current_repodata.json): done  
Solving environment: done  
  
## Package Plan ##
```

```

Kali-Linux-2020.3-vbox.amd64 [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
spark-3... Extract [Installati... kali@kali... [kali@kali... kali@kali... Download... 10:03 AM
File Actions Edit View Help
zope.interface      pkgs/main/linux-64::zope.interface-5.1.2-py38h7b6447c_0
zstd                pkgs/main/linux-64::zstd-1.4.5-h9ceee32_0

Preparing transaction: done
Executing transaction: done
installation finished.
Do you wish the installer to initialize Anaconda3
by running conda init? [yes|no]
[no] >>> yes
Anaconda3-2020.11-Linux-x86_64.sh: 494: [: not found
no change  /home/kali/anaconda3/condabin/conda
no change  /home/kali/anaconda3/bin/conda
no change  /home/kali/anaconda3/bin/conda-env
no change  /home/kali/anaconda3/bin/activate
no change  /home/kali/anaconda3/bin/deactivate
no change  /home/kali/anaconda3/etc/profile.d/conda.sh
no change  /home/kali/anaconda3/etc/fish/conf.d/conda.fish
no change  /home/kali/anaconda3/shell/condabin/Conda.psm1
no change  /home/kali/anaconda3/shell/condabin/conda-hook.ps1
no change  /home/kali/anaconda3/lib/python3.8/site-packages/xontrib/conda.xsh
no change  /home/kali/anaconda3/etc/profile.d/conda.csh
modified   /home/kali/.bashrc

⇒ For changes to take effect, close and re-open your current shell. ←
If you'd prefer that conda's base environment not be activated on startup,
set the auto_activate_base parameter to false:
conda config --set auto_activate_base false

Thank you for installing Anaconda3!

```

Upon completion of the installation, you can verify that the installation is okay by opening the python3 at the terminal command line:

```

kali@kali:~/Downloads$ python3
Python 3.9.1 (default, Dec  8 2020, 07:51:42)
[GCC 10.2.0] on linux
Type "help", "copyright", "credits" or "license" for more information.
>>> exit()
kali@kali:~/Downloads$ 

```

You can type 'exit()' to quit/exit from pyton3.

Next, we can launch Jupyter Notebook through anaconda navigator, a desktop graphical user interface (GUI) that is included in the Anaconda distribution package.

Navigate to the bin folder of anaconda3.

```

kali@kali:~$ cd ~/anaconda3/bin
kali@kali:~/anaconda3/bin$ 

```

```

kali@kali:~$ ls
anaconda3  CSCI361  Documents  Music  Public          spark-3.0.1-bin-hadoop2.7.tgz  Videos
CSCI262  Desktop  Downloads  Pictures  spark-3.0.1-bin-hadoop2.7  Templates
kali@kali:~$ cd anaconda3
kali@kali:~/anaconda3$ ls
bin        condabin  doc  etc  info  libexec  man  phrasebooks  plugins  resources  share  ssl  var
compiler_compat  conda-meta  envs  include  lib  LICENSE.txt  mkspecs  pkgs  qml  sbin  shell  translations  x86_64-conda_cos6-linux-gnu
kali@kali:~/anaconda3$ cd bin
kali@kali:~/anaconda3/bin$ 

```

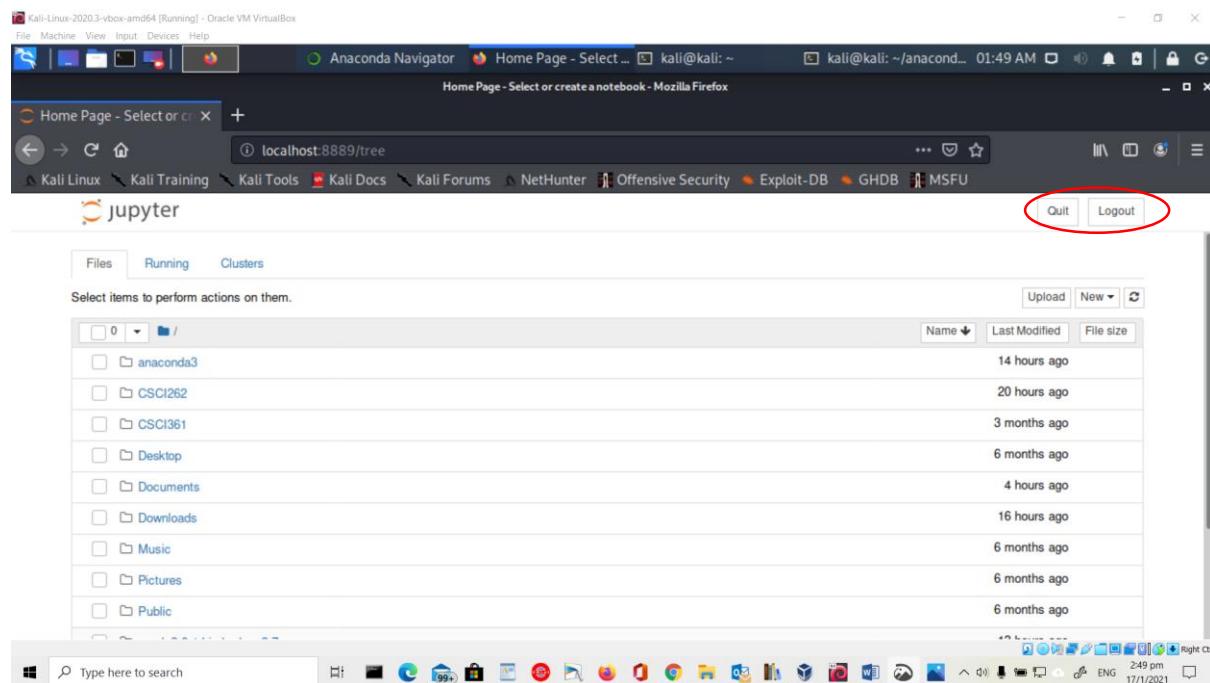
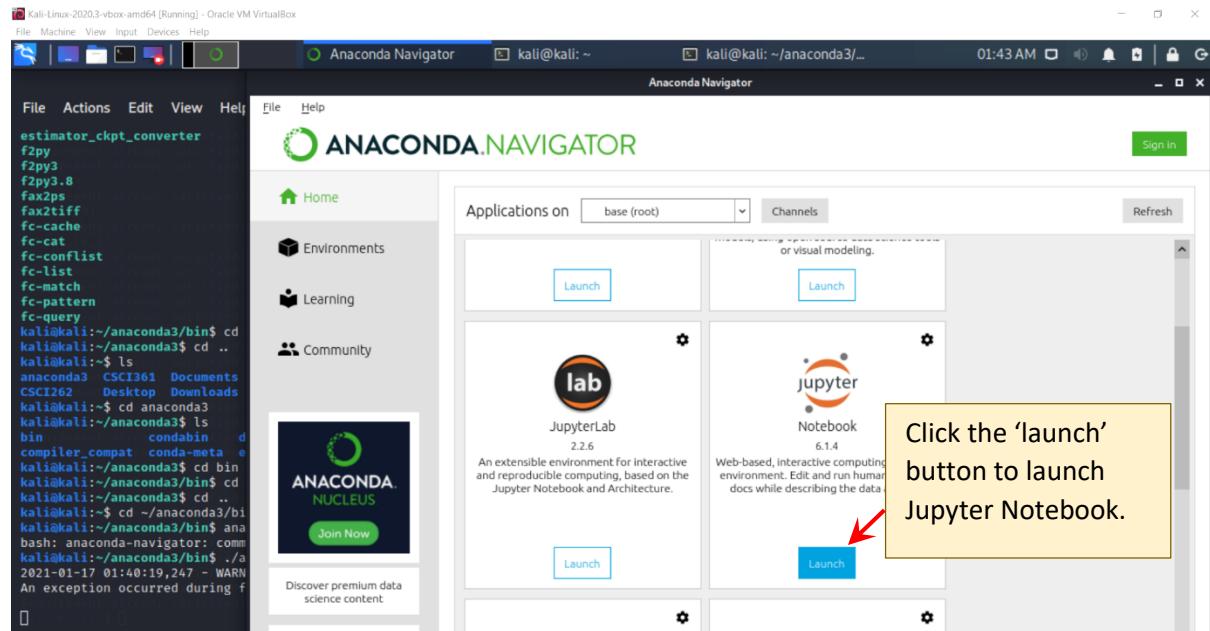
Next, execute the anaconda-navigator.

```

kali@kali:~/anaconda3/bin$ ./anaconda-navigator
2021-01-17 01:40:19,247 - WARNING linux_scaling.get_scaling_factor_using_dbus:27
An exception occurred during fetching list of system display settings.
Requirement already satisfied: science-content

```

An Anaconda Navigator GUI window is opened. Look for “Jupyter” and click on the ‘launch’ button to launch Jupyter Notebook.

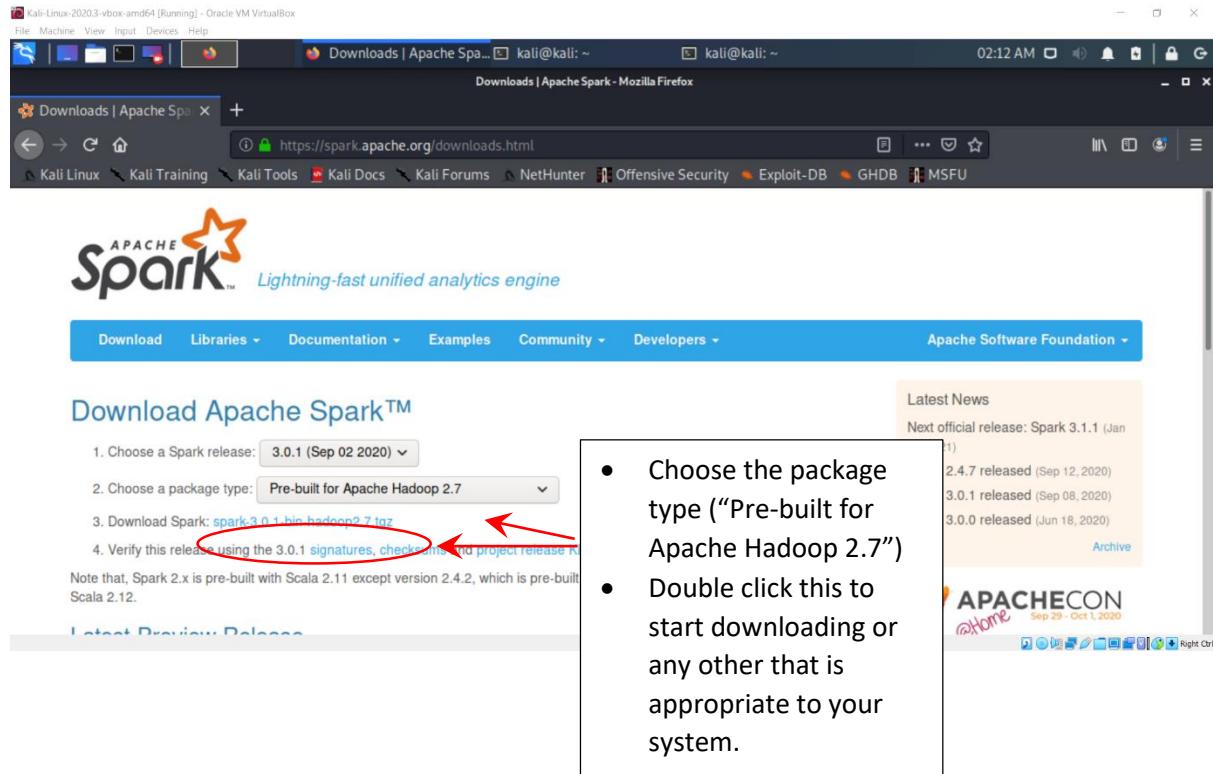


We will come back to Jupyter Notebook in a while. For now, you can quit and logout from Jupyter Notebook by clicking on the ‘Quit’ and ‘Logout’ button on the right-hand-top-corner of the screen.

# Downloading and setup PySpark

After installing all the required packages, we can now download and setup the PySpark.

Open your browser and enter the following URL: <https://spark.apache.org/downloads.html>



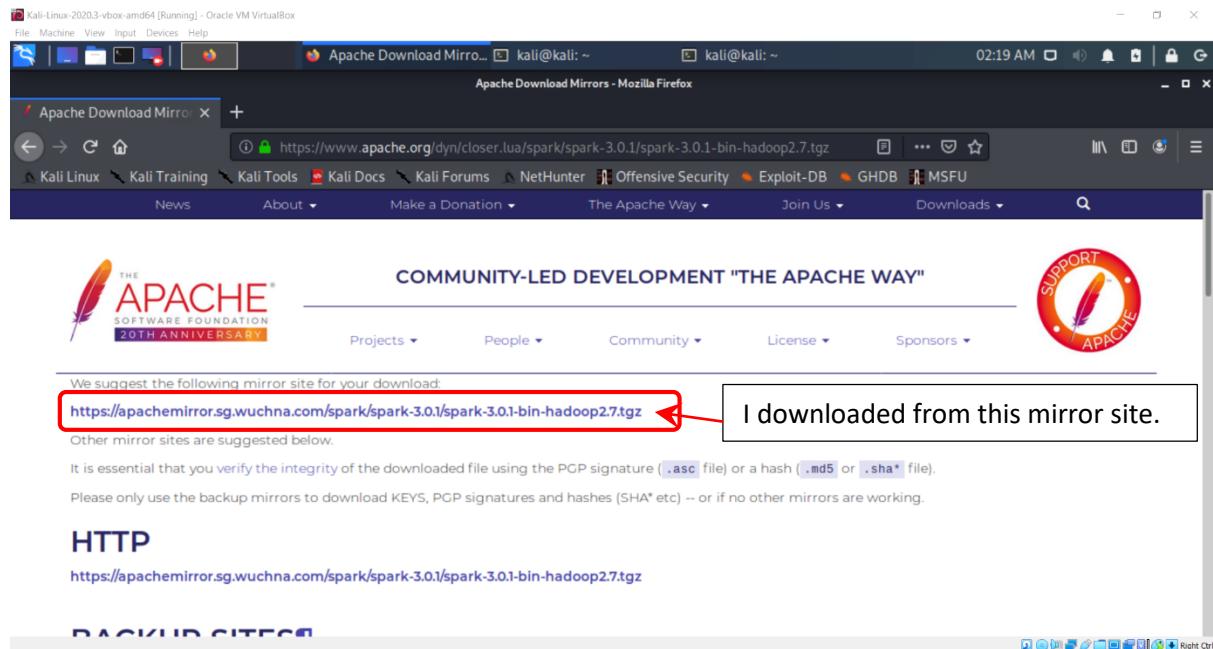
The screenshot shows the Apache Spark download page. The URL in the address bar is <https://spark.apache.org/downloads.html>. The page title is "Download Apache Spark™". The main content area has a list of steps:

1. Choose a Spark release: 3.0.1 (Sep 02 2020)
2. Choose a package type: Pre-built for Apache Hadoop 2.7
3. Download Spark: [spark-3.0.1-bin-hadoop2.7.tgz](https://spark.apache.org/spark-3.0.1-bin-hadoop2.7.tgz)
4. Verify this release using the 3.0.1 [signatures, checksums and project release notes](#)

A note below states: "Note that, Spark 2.x is pre-built with Scala 2.11 except version 2.4.2, which is pre-built Scala 2.12." To the right, there is a sidebar with "Latest News" and an "APACHECON @Home" section.

**• Choose the package type ("Pre-built for Apache Hadoop 2.7")**  
**• Double click this to start downloading or any other that is appropriate to your system.**

At the next screen, select the link from any one of the suggested mirror site.



The screenshot shows the Apache Download Mirrors page. The URL in the address bar is <https://www.apache.org/dyn/closer.lua/spark/spark-3.0.1/spark-3.0.1-bin-hadoop2.7.tgz>. The page title is "Apache Download Mirrors". The main content area says "COMMUNITY-LED DEVELOPMENT 'THE APACHE WAY'". It lists several mirror sites, with the first one highlighted:

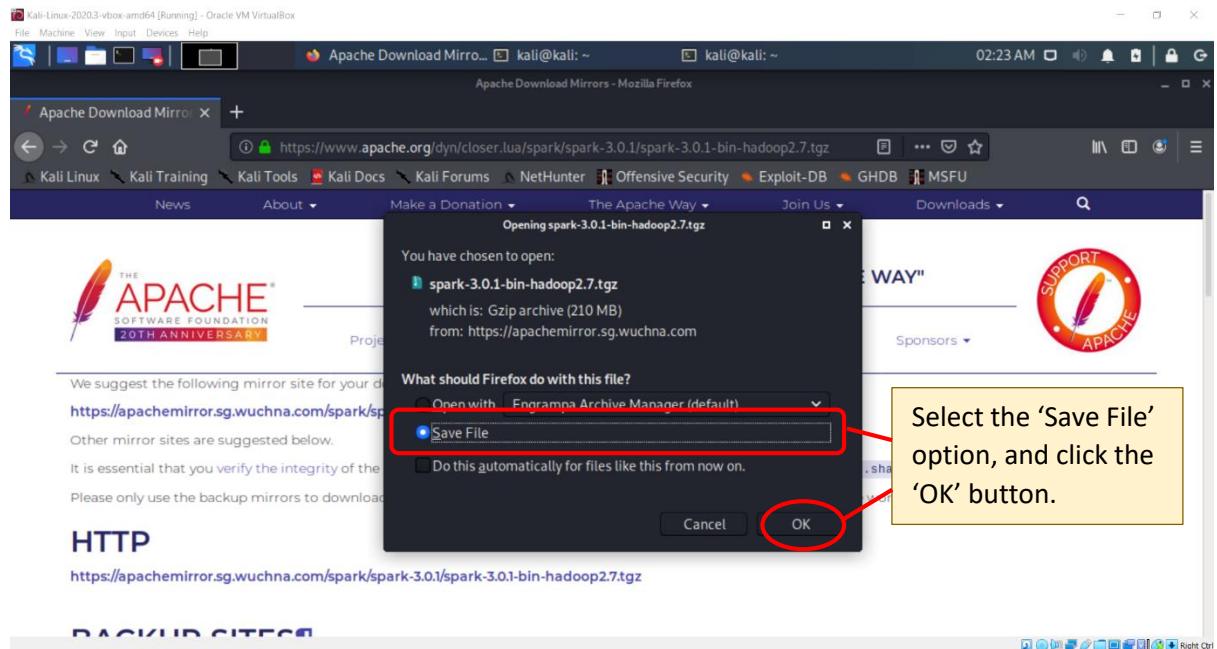
We suggest the following mirror site for your download:  
<https://apachemirror.sg.wuchna.com/spark/spark-3.0.1/spark-3.0.1-bin-hadoop2.7.tgz>

**I downloaded from this mirror site.**

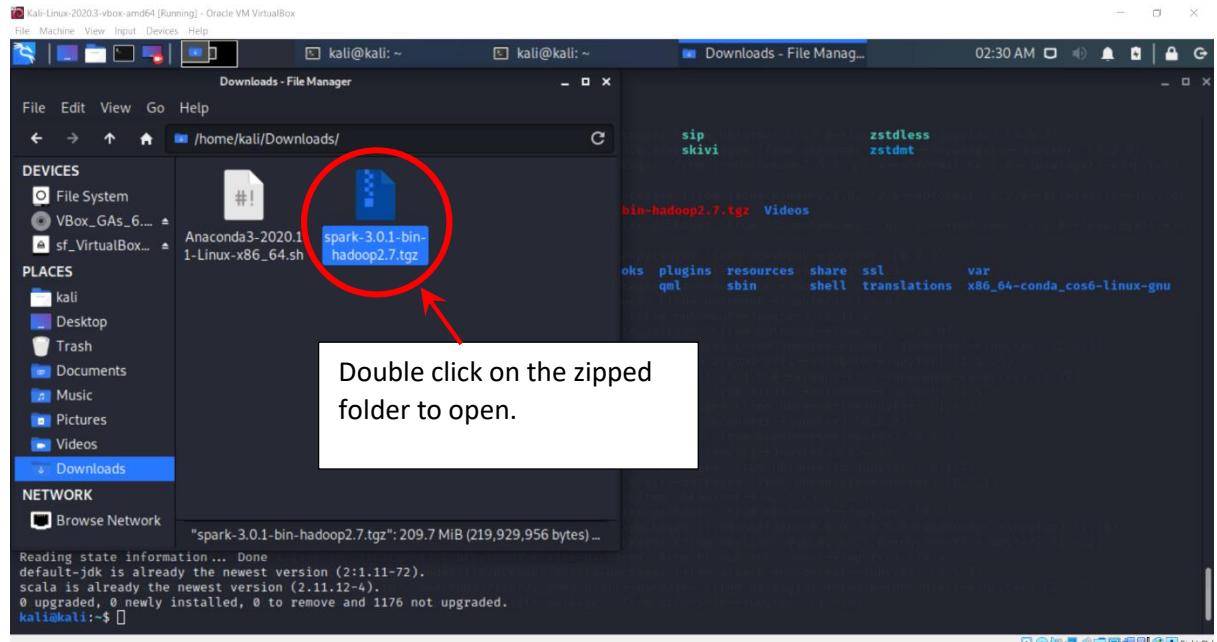
Other mirror sites are suggested below. It is essential that you verify the integrity of the downloaded file using the PGP signature ([.asc](#) file) or a hash ([.md5](#) or [.sha\\*](#) file). Please only use the backup mirrors to download KEYS, PGP signatures and hashes (SHA\* etc) -- or if no other mirrors are working.

**HTTP**  
<https://apachemirror.sg.wuchna.com/spark/spark-3.0.1/spark-3.0.1-bin-hadoop2.7.tgz>

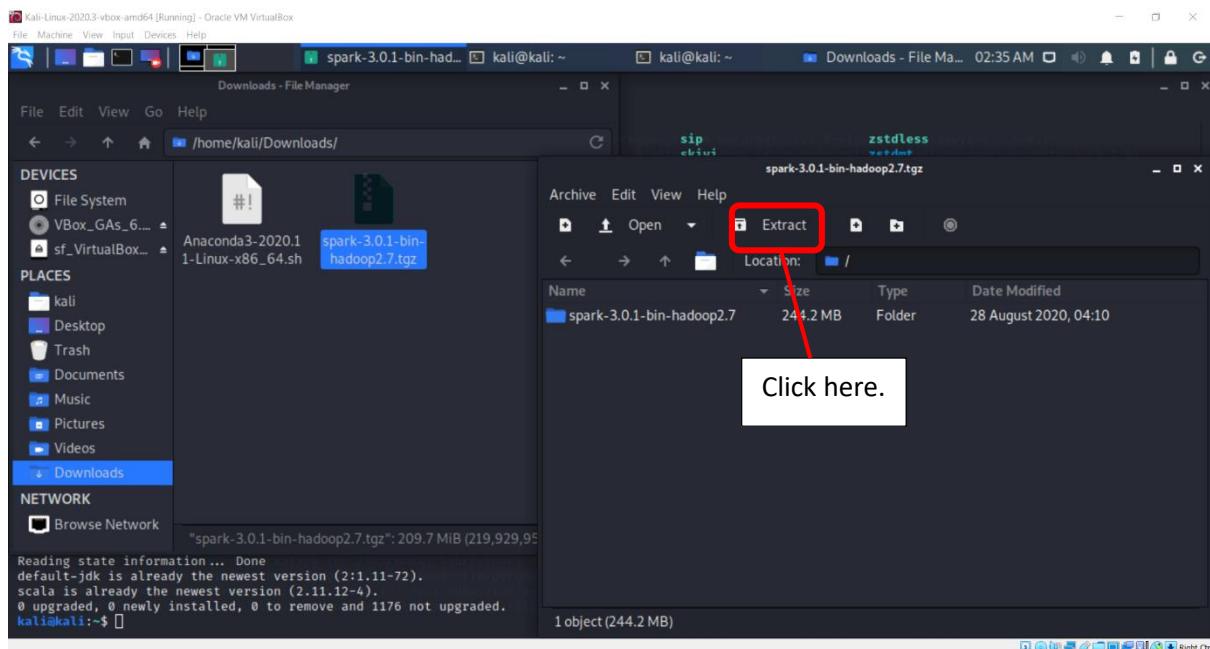
Next, select the ‘Save File’ option and click ‘OK’ button to save the downloaded package.



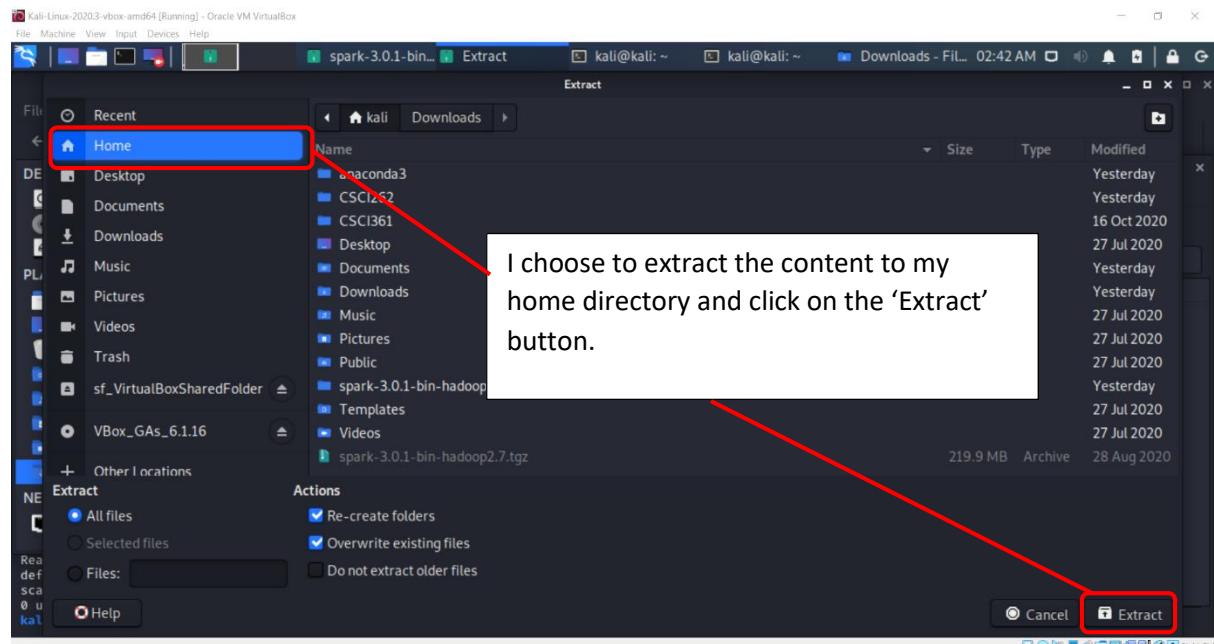
After the downloading is completed, navigate to the ‘Downloads’ folder and open the downloaded packages. (Note: The downloaded package is compressed using zip application. We need to extract (unzipped) the access to the information.)



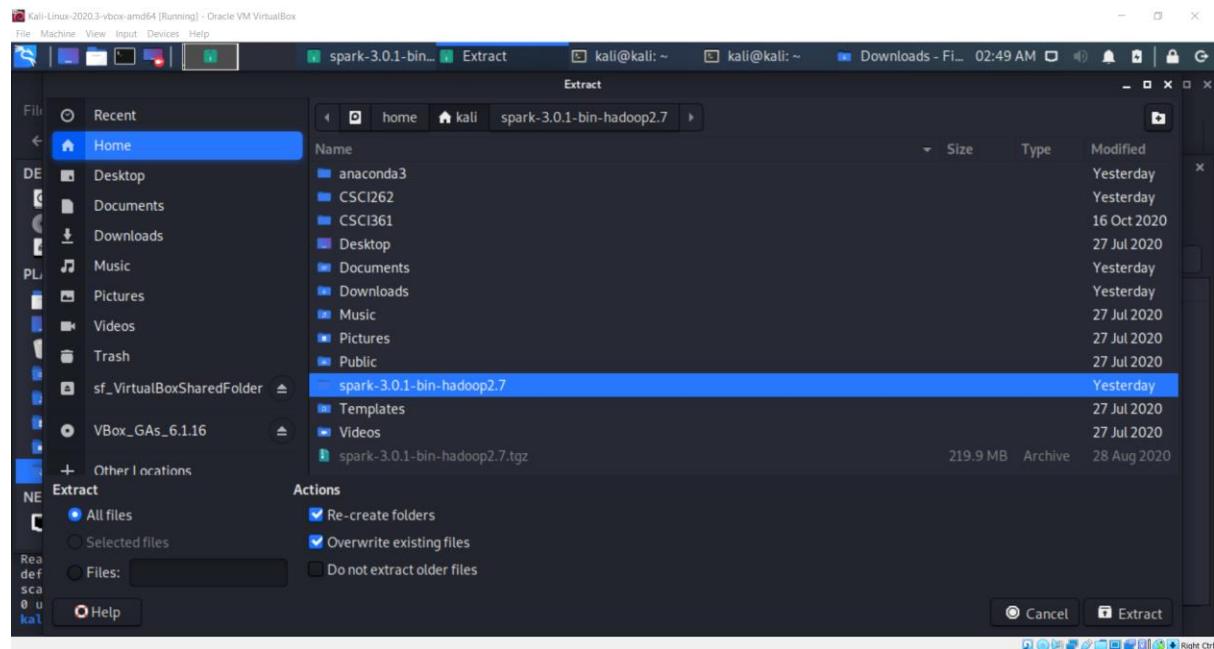
Click the ‘Extract’ button to start extracting the content of the package.



And select the folder where you want the content extracted to. I select to extract the content to my home directory.



Check that a folder spark-3.0.1-bin-hadoop2.7 is now appear in my home folder.



## Setting up the environment variables

Next, we need to set up a couple of environment variables to launch PySpark with Python3 and enable it to be called from Jupyter Notebook.

Before we edit the .bashrc, we need to get ready all the path (location) where the various required packages are located.

Finding Java jdk location:

```
kali㉿kali:~$ update-alternatives --list java
/usr/lib/jvm/java-11-openjdk-amd64/bin/java
/usr/lib/jvm/java-8-openjdk-amd64/jre/bin/java
kali㉿kali:~$
```

Finding Python3 location:

```
kali㉿kali:~$ which python3
/usr/bin/python3
```

Finding iPython3 location

```
kali㉿kali:~$ which ipython3  
/usr/bin/ipython3  
kali㉿kali:~$ █
```

Home location of PySpark:

/home/kali/spark-3.0.1-bin-hadoop2.7

Next, we need to export the paths to various packages to the respective environment variables. **Take a backup of .bashrc before proceeding**, just in case, you need to go back.

Open .baschrc using nano editor.

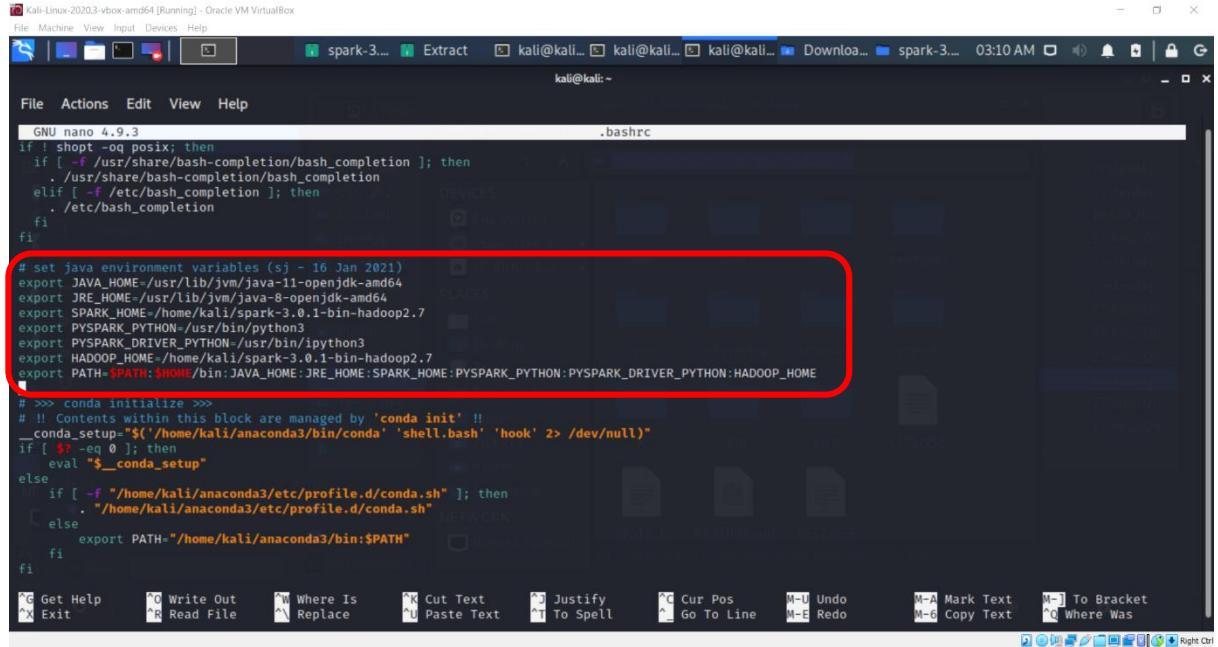
```
kali㉿kali:~$ nano .bashrc
```

Enter the following lines at the end of .bashrc:

```
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64    (Note: remove the /bin/java from the path.)  
export JRE_HOME=/usr/lib/jv/java-8-openjdk-amd64    (Note: remove the /jre/bin/java from the path.)  
export SPARK_HOME=/home/kali/spark-3.0.1-bin-hadoop2.7  
export PYSPARK_PYTHON=/usr/bin/python3  
export PYSPARK_DRIVER_PYTHON=/usr/bin/ipython3  
export HADOOP_HOME=/home/kali/spark-3.0.1-bin-hadoop2.7  
export  
PATH=$PATH:$HOME/bin:JAVA_HOME:JRE_HOME:SPARK_HOME:PYSPARK_PYTHON:PYSPARK_DRIVER_PYTHON:HADOOP_HOME
```

Check that your entries are correct. If any of the entry is incorrect, we may encounter quite a bit of problems later.

Press ^O (Control+O) to save (write out) the changes. To exit the nano editor, press ^X (Control+X).



```
GNU nano 4.9.3 .bashrc
if ! shopt -oq posix; then
  if [ -f /usr/share/bash-completion/bash_completion ]; then
    . /usr/share/bash-completion/bash_completion
  elif [ -f /etc/bash_completion ]; then
    . /etc/bash_completion
  fi

# set java environment variables (sj - 16 Jan 2021)
export JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
export JRE_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export SPARK_HOME=/home/kali/spark-3.0.1-bin-hadoop2.7
export PYSPARK_PYTHON=/usr/bin/python3
export PYSPARK_DRIVER_PYTHON=/usr/bin/ipython3
export HADOOP_HOME=/home/kali/spark-3.0.1-bin-hadoop2.7
export PATH=$PATH:$HOME/bin:JAVA_HOME:JRE_HOME:SPARK_HOME:PYSPARK_PYTHON:PYSPARK_DRIVER_PYTHON:HADOOP_HOME

# >>> conda initialize >>>
# !! Contents within this block are managed by 'conda init' !!
__conda_setup="$('/home/kali/anaconda3/bin/conda' 'shell.bash' '&gt;' '/dev/null')"
if [ $? -eq 0 ]; then
  eval "$__conda_setup"
else
  if [ -f "/home/kali/anaconda3/etc/profile.d/conda.sh" ]; then
    . "/home/kali/anaconda3/etc/profile.d/conda.sh"
  else
    export PATH="/home/kali/anaconda3/bin:$PATH"
  fi
fi

G Get Help   ^O Write Out   ^W Where Is   ^K Cut Text   ^J Justify   C Cur Pos   M-U Undo
X Exit      ^R Read File   ^B Replace   ^U Paste Text  ^T To Spell   M-E Redo   M-A Mark Text   M-[ To Bracket
                                         ^_ Go To Line   M-B Copy Text   ^Q Where Was
                                         Right Ctrl
```

## Installing important libraries and dependencies

Before we can start calling the PySpark using Jupyter Notebook, we need to import/install some Python libraries (all of which can be installed with pip/pip3). Note: Some libraries may have already installed together with Anaconda.

- i. IPython3 (already installed with Anaconda)
- ii. Jupyter Notebook (already installed with Anaconda)
- iii. Tensorflow2 (Need to install using PIP)
- iv. Scikit-Learn (Need to install using PIP)
- v. Scientific computing libraries: SciPy, NumPy, Pandas, Matplotlib, importlib, psutil (Some of these are already installed with Anaconda)
- vi. Some dependencies libraries: Py4J, pyarrow, psutil, cairocffi (Some of these are already install with Anaconda)
- vii. findspark (this library is important. Must import)
- viii. pyspark (this library is also important. Some important dependencies are installed here.)

Following are the steps to import using PIP. If those libraries are already installed together with Anaconda, the system will inform.)

### Installing findspark

```
(base) kali@kali:~$ python3 -m pip install findspark
Collecting findspark
  Downloading findspark-1.4.2-py2.py3-none-any.whl (4.2 kB)
Installing collected packages: findspark
Successfully installed findspark-1.4.2
```

### Installing pyspark

```
kali@kali:~$ pip3 install pyspark==3.0.1
Collecting pyspark==3.0.1
  Downloading pyspark-3.0.1.tar.gz (204.2 MB)
    |██████████| 204.2 MB 57 kB/s
Collecting py4j==0.10.9
  Downloading py4j-0.10.9-py2.py3-none-any.whl (198 kB)
    |██████████| 198 kB 72.1 MB/s
Building wheels for collected packages: pyspark
  Building wheel for pyspark: setup.py ... done
    Created wheel for pyspark: filename=pyspark-3.0.1-py2.py3-none-any.whl size=204612244 sha256=024e99e20bdd35357628b5f3718ec111683e0ffd7702f565999066be64c
97c7
  Stored in directory: /home/kali/.cache/pip/wheels/4d/60/70/4b354ff632e827ce13a755d886b704e306089e6c275be8aba4
Successfully built pyspark
Installing collected packages: py4j, pyspark
Successfully installed py4j-0.10.9 pyspark-3.0.1
kali@kali:~$
```

### Installing Scikit-Learn

```
kali@kali:~$ pip3 install Scikit-Learn
Collecting Scikit-Learn
  Downloading scikit_learn-0.24.0-cp39-cp39-manylinux2010_x86_64.whl (23.8 MB)
    |██████████| 23.8 MB 9.5 MB/s
Collecting threadpoolctl>=2.0.0
  Downloading threadpoolctl-2.1.0-py3-none-any.whl (12 kB)
Collecting joblib>=0.11
  Downloading joblib-1.0.0-py3-none-any.whl (302 kB)
    |██████████| 302 kB 49.5 MB/s
Requirement already satisfied: numpy>=1.13.3 in /usr/lib/python3/dist-packages (from Scikit-Learn) (1.19.4)
Requirement already satisfied: scipy>=0.19.1 in /usr/lib/python3/dist-packages (from Scikit-Learn) (1.5.4)
Installing collected packages: threadpoolctl, joblib, Scikit-Learn
Successfully installed Scikit-Learn-0.24.0 joblib-1.0.0 threadpoolctl-2.1.0
kali@kali:~$
```

### Installing scientific computing libraries

```
kali@kali:~$ pip3 install NumPy
Requirement already satisfied: NumPy in /usr/lib/python3/dist-packages (1.19.4)
kali@kali:~$ pip3 install SciPy
Requirement already satisfied: SciPy in /usr/lib/python3/dist-packages (1.5.4)
kali@kali:~$ pip3 install Pandas
Requirement already satisfied: Pandas in /usr/lib/python3/dist-packages (1.1.5)
kali@kali:~$ pip3 install Matplotlib
Requirement already satisfied: Matplotlib in /usr/lib/python3/dist-packages (3.3.2)
```

### Installing dependencies libraries

```
kali@kali:~$ pip3 install Py4J
Requirement already satisfied: Py4J in ./local/lib/python3.9/site-packages (0.10.9)

(base) kali@kali:~$ pip3 install pyarrow
Collecting pyarrow
  Downloading pyarrow-2.0.0-cp38-cp38-manylinux2014_x86_64.whl (17.8 MB)
    |██████████| 17.8 MB 10.6 MB/s
Requirement already satisfied: numpy>=1.14 in ./anaconda3/lib/python3.8/site-packages (from pyarrow) (1.19.2)
Installing collected packages: pyarrow
Successfully installed pyarrow-2.0.0
(base) kali@kali:~$
```

```
(base) kali㉿kali:~$ pip3 install psutil
Requirement already satisfied: psutil in /usr/lib/python3/dist-packages (5.7.3)
kali㉿kali:~$ 

(base) kali㉿kali:~$ pip3 install cairocffi
Collecting cairocffi
  Downloading cairocffi-1.2.0.tar.gz (70 kB)
    |██████████| 70 kB 3.4 MB/s
Requirement already satisfied: cffi>=1.1.0 in ./anaconda3/lib/python3.8/site-packages (from cairocffi) (1.14.3)
Requirement already satisfied: pycparser in ./anaconda3/lib/python3.8/site-packages (from cffi>=1.1.0→cairocffi) (2.20)
Building wheels for collected packages: cairocffi
  Building wheel for cairocffi (setup.py) ... done
    Created wheel for cairocffi: filename=cairocffi-1.2.0-py3-none-any.whl size=89546 sha256=65bcf974c0466c7e1dd919e0fe361d1156c5cce8609557d53faf307a4d12798
    Stored in directory: /home/kali/.cache/pip/wheels/e8/fa/11/ae7a999afdf58d7169974c79b765f4f03880c184d578b1ed445
Successfully built cairocffi
Installing collected packages: cairocffi
Successfully installed cairocffi-1.2.0
(base) kali㉿kali:~$ 
```

## Installling Tensorflow2

```
(base) kali㉿kali:~$ pip3 install --upgrade tensorflow
Collecting tensorflow
  Downloading tensorflow-2.4.0-cp38-cp38-manylinux2010_x86_64.whl (394.8 MB)
    |██████████| 394.8 MB 50 kB/s
Requirement already satisfied: typing-extensions~=3.7.4 in ./anaconda3/lib/python3.8/site-packages (from tensorflow) (3.7.4.3)
Requirement already satisfied: h5py<~2.10.0 in ./anaconda3/lib/python3.8/site-packages (from tensorflow) (2.10.0)
Requirement already satisfied: numpy~=1.19.2 in ./anaconda3/lib/python3.8/site-packages (from tensorflow) (1.19.2)
Requirement already satisfied: wheel<~0.35 in ./anaconda3/lib/python3.8/site-packages (from tensorflow) (0.35.1)
Requirement already satisfied: six~=1.15.0 in ./anaconda3/lib/python3.8/site-packages (from tensorflow) (1.15.0)
Collecting gast==0.3.3
  Downloading gast-0.3.3-py2.py3-none-any.whl (9.7 kB)
Collecting absl-py==0.10
  Downloading absl_py-0.11.0-py3-none-any.whl (127 kB)
    |██████████| 127 kB 59.2 MB/s
Collecting astunparse==1.6.3
  Downloading astunparse-1.6.3-py2.py3-none-any.whl (12 kB)
Collecting flatbuffers==1.12.0
  Downloading flatbuffers-1.12-py2.py3-none-any.whl (15 kB)
Collecting google-pasta==0.2
  Downloading google_pasta-0.2.0-py3-none-any.whl (57 kB)
    |██████████| 57 kB 8.8 MB/s
Collecting grpcio==1.32.0
  Downloading grpcio-1.32.0-cp38-cp38-manylinux2014_x86_64.whl (3.8 MB)
    |██████████| 3.8 MB 62.3 MB/s
Collecting keras-preprocessing==1.1.2
  Downloading Keras_Preprocessing-1.1.2-py2.py3-none-any.whl (42 kB)
    |██████████| 42 kB 1.8 MB/s
Collecting opt-einsum==3.3.0
  Downloading opt_einsum-3.3.0-py3-none-any.whl (65 kB)
    |██████████| 65 kB 3.5 MB/s
Collecting protobuf≥3.9.2
  Downloading protobuf-3.14.0-cp38-cp38-manylinux1_x86_64.whl (1.0 MB)
    |██████████| 1.0 MB 43.9 MB/s 
```

```
|██████████| 147 kB 35.6 MB/s
Collecting rsa<5, ≥ 3.1.4
  Downloading rsa-4.7-py3-none-any.whl (34 kB)
Collecting tensorboard-plugin-wit≥1.6.0
  Downloading tensorboard_plugin_wit-1.7.0-py3-none-any.whl (779 kB)
    |██████████| 779 kB 49.6 MB/s
Collecting tensorflow-estimator<2.5.0, ≥ 2.4.0rc0
  Downloading tensorflow_estimator-2.4.0-py2.py3-none-any.whl (462 kB)
    |██████████| 462 kB 59.4 MB/s
Collecting termcolor==1.1.0
  Downloading termcolor-1.1.0.tar.gz (3.9 kB)
Collecting wrapt==1.12.1
  Downloading wrapt-1.12.1.tar.gz (27 kB)
Building wheels for collected packages: termcolor, wrapt
  Building wheel for termcolor (setup.py) ... done
    Created wheel for termcolor: filename=termcolor-1.1.0-py3-none-any.whl size=4830 sha256=0edc79ddeeb63862b8171aa20898daff9a5006f6a7636d9c24d81e7e1d8930d9
    Stored in directory: /home/kali/.cache/pip/wheels/a0/16/9c/5473df82468f958445479c59e784896fa24f4a5fc024b0f501
  Building wheel for wrapt (setup.py) ... done
    Created wheel for wrapt: filename=wrapt-1.12.1-cp38-cp38-linux_x86_64.whl size=81380 sha256=4ce880d2eeadbc1dc0ef783035314b403155b6f78aeb3e6c5e8a0b7db3bde
cb9
    Stored in directory: /home/kali/.cache/pip/wheels/5f/fd/9e/b6cf5890494cb8ef0b5eaff72e5d5a70fb56316007d6dfe73
Successfully built termcolor wrapt
Installing collected packages: pyasn1, rsa, pyasn1-modules, oauthlib, cachetools, requests-oauthlib, google-auth, tensorflow-plugin-wit, protobuf, markdown, grpcio, google-auth-oauthlib, absl-py, wrapt, termcolor, tensorflow-estimator, tensorboard, opt-einsum, keras-preprocessing, google-pasta, gast, flatbuffers, astunparse, tensorflow
  Attempting uninstall: wrapt
    Found existing installation: wrapt 1.11.2
    Uninstalling wrapt-1.11.2:
      Successfully uninstalled wrapt-1.11.2
Successfully installed absl-py-0.11.0 astunparse-1.6.3 cachetools-4.2.0 flatbuffers-1.12 gast-0.3.3 google-auth-1.24.0 google-auth-oauthlib-0.4.2 google-pasta-0.2.0 grpcio-1.32.0 keras-preprocessing-1.1.2 markdown-3.3.3 oauthlib-3.1.0 opt-einsum-3.3.0 protobuf-3.14.0 pyasn1-0.4.8 pyasn1-modules-0.2.8 requests-oauthlib-1.3.0 rsa-4.7 tensorboard-2.4.1 tensorboard-plugin-wit-1.7.0 tensorflow-2.4.0 tensorflow-estimator-2.4.0 termcolor-1.1.0 wrapt-1.12.1
(base) kali㉿kali:~$ 
```

We should be able to activate Anaconda and call the Jupyter Notebook to test the PySpark now.

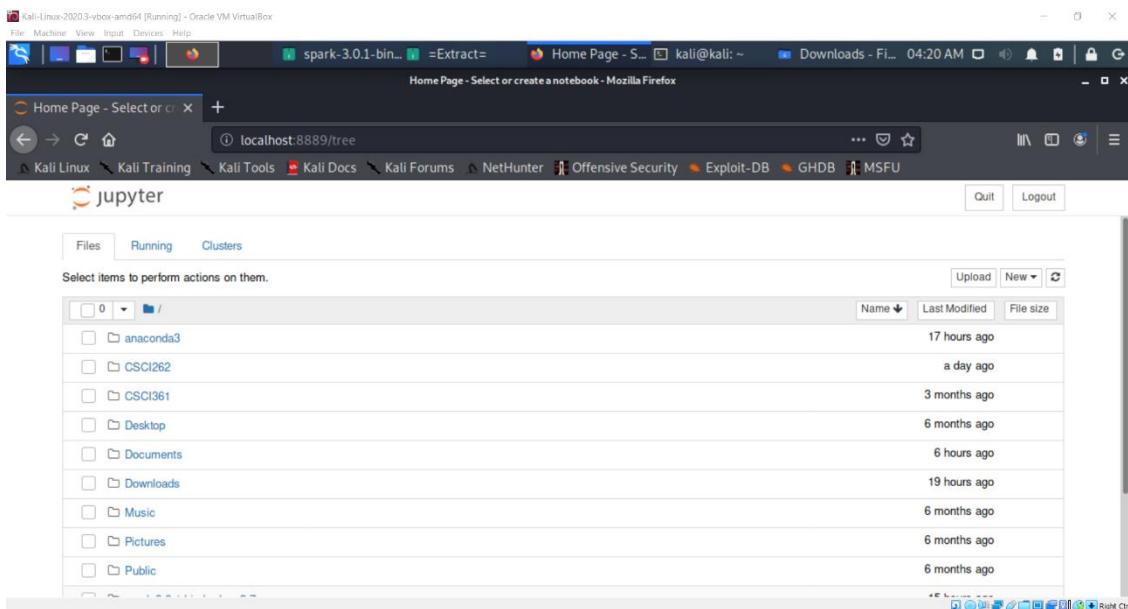
- Activate Anaconda

```
kali㉿kali:~$ conda activate  
(base) kali㉿kali:~$ █
```

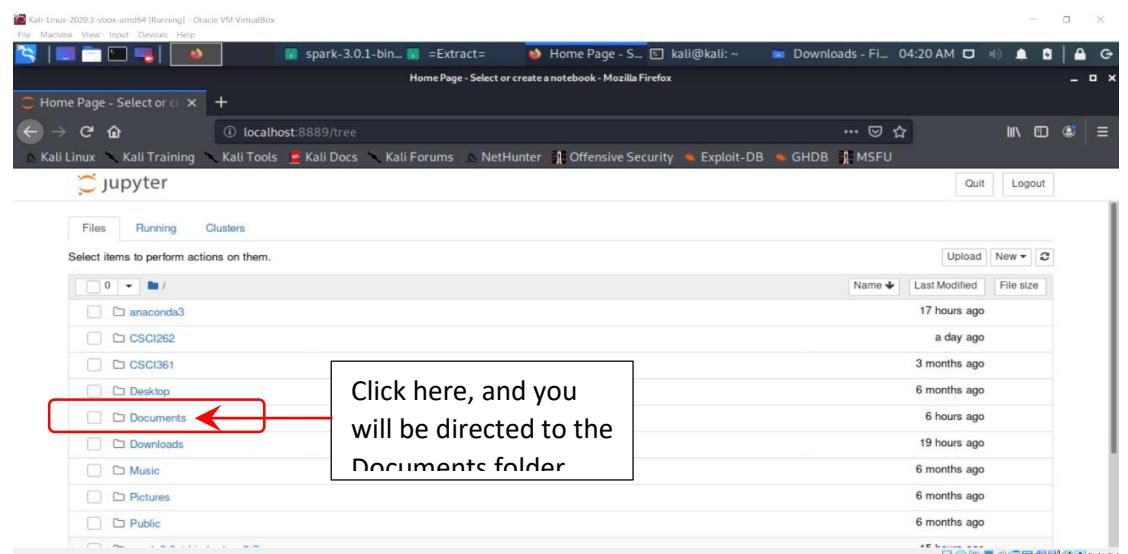
- Accessing Jupyter Notebook

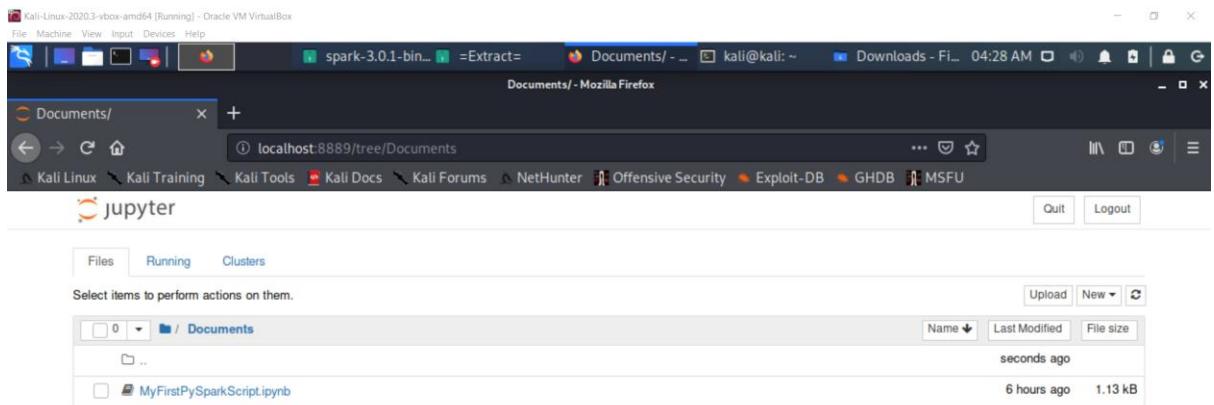
```
(base) kali㉿kali:~$ jupyter notebook █
```

Jupyter Notebook is invoked and we will see a locally hosted Jupyter Notebook as follow:

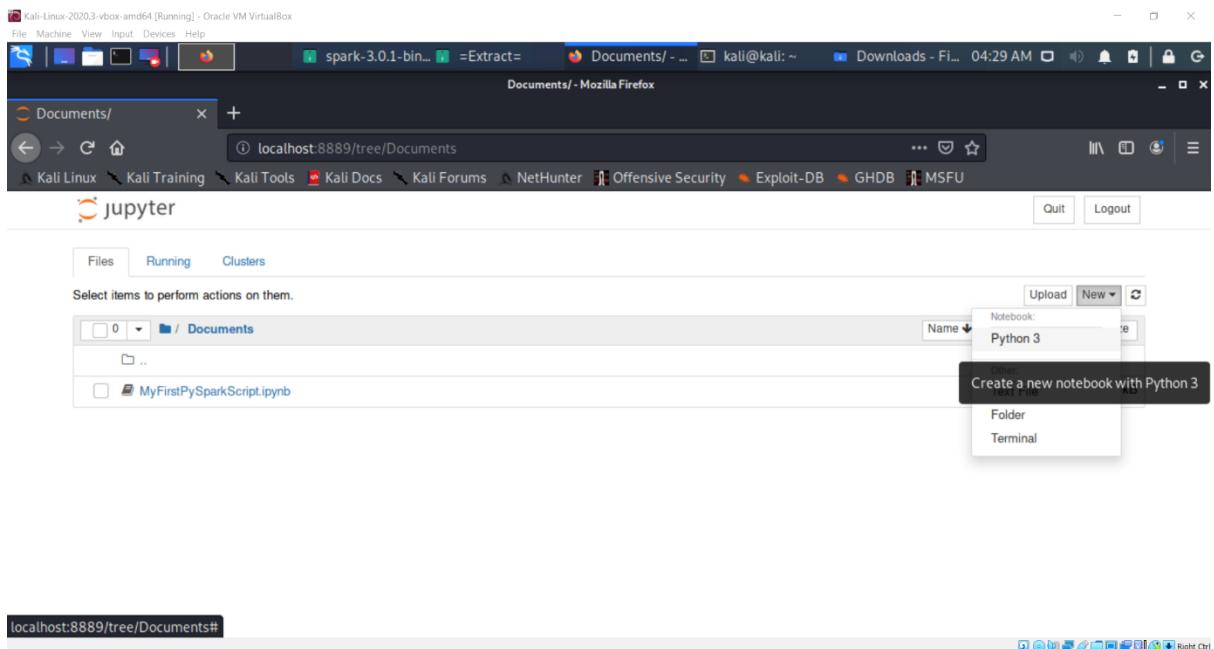


Choose a folder where you want your PySpark script to be stored. I choose my Documents folder. You can choose any folder you like.

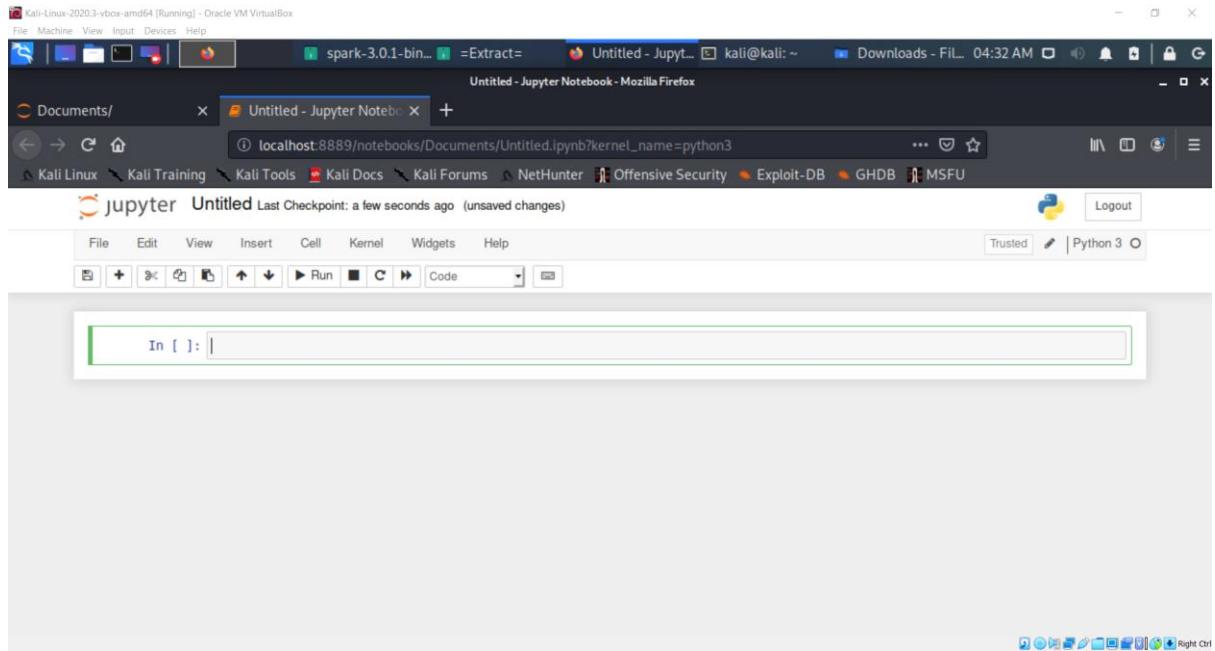




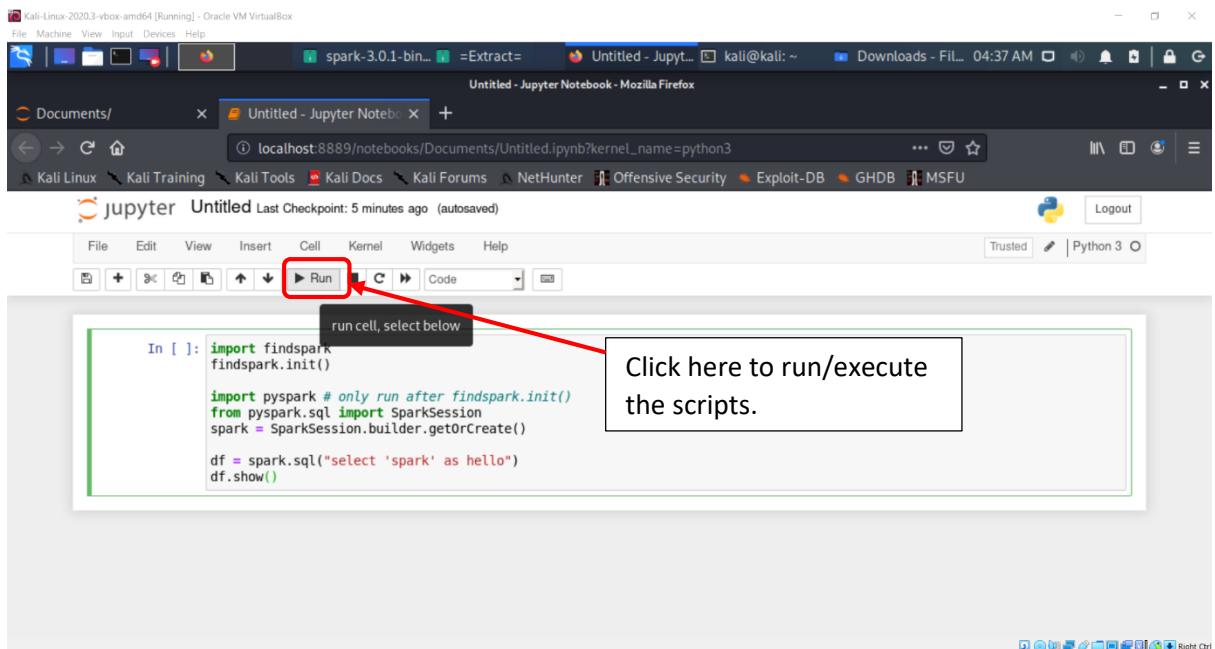
Open a new Python3 notebook. Click at the ‘New’ drop-down button, and a drop-down menu is displayed.



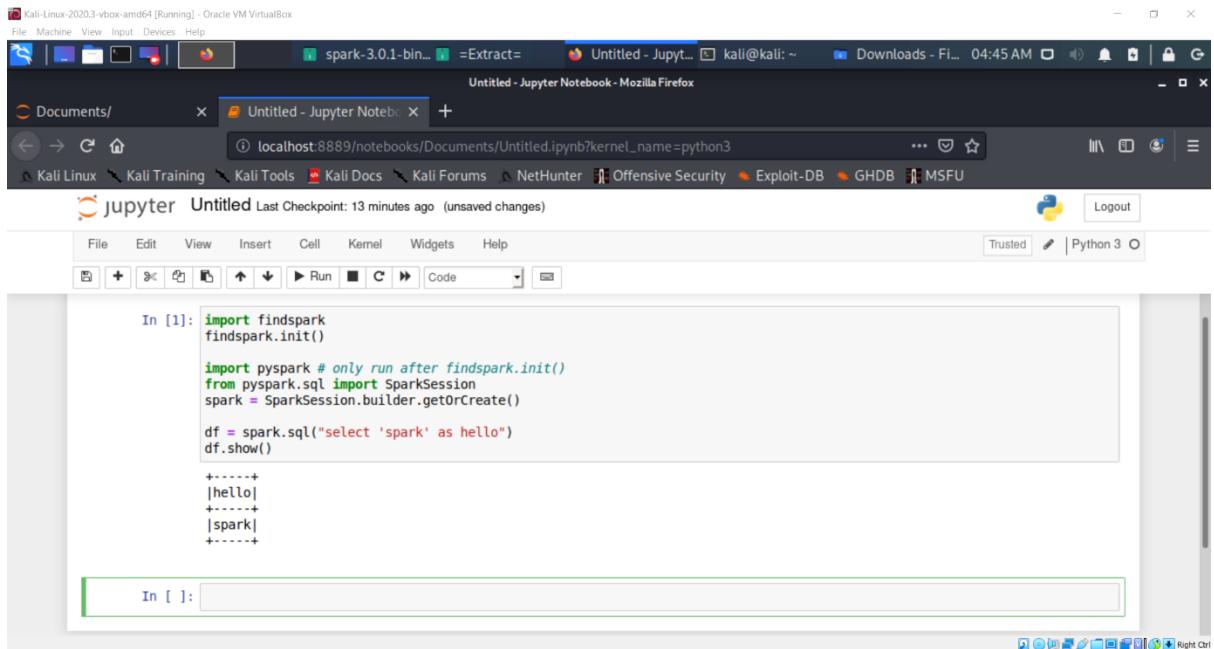
Select and click on the option ‘Python3’, and a new notebook is open as follow:



At the cell, type in the script (as provided in the notes “How to Install and Run PySpark in Jupyter Notebook on Windows, by Chang Hsin Lee”)



When you have finished typing the script, you can click on the “Run” button at the command ribbon. The script is run (executed) and you should see the output as shown below.



The screenshot shows a Jupyter Notebook interface running on a Kali Linux VM. The notebook title is "Untitled - Jupyter Notebook - Mozilla Firefox". The code cell (In [1]) contains the following PySpark script:

```
import findspark
findspark.init()

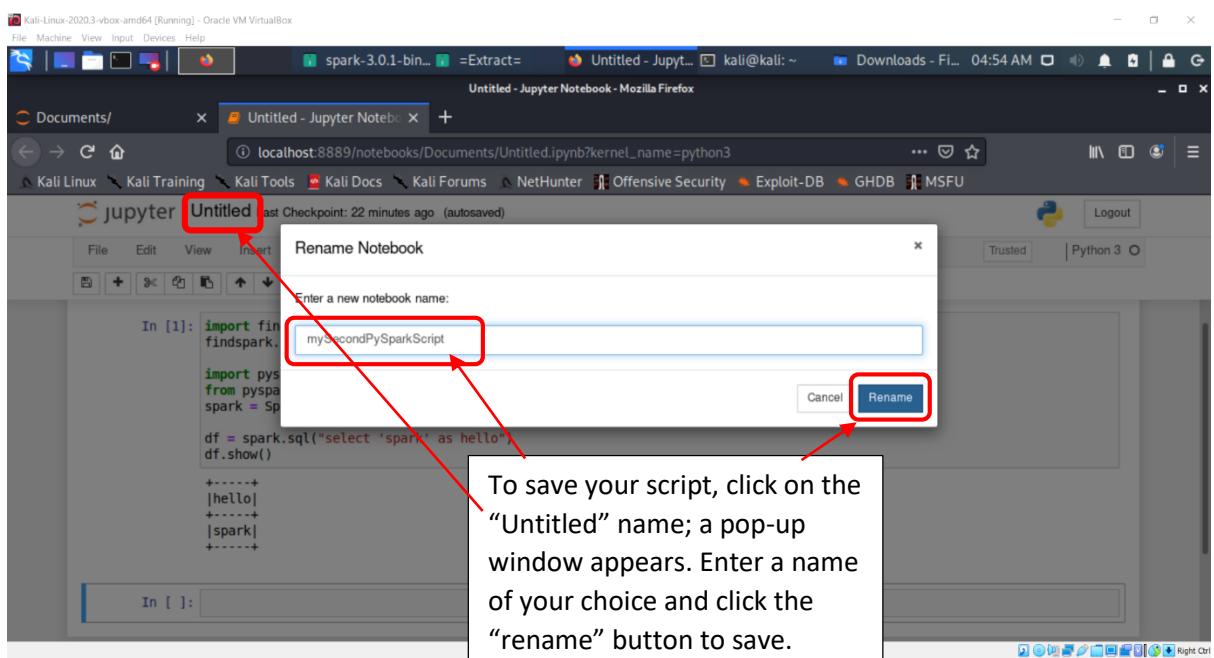
import pyspark # only run after findspark.init()
from pyspark.sql import SparkSession
spark = SparkSession.builder.getOrCreate()

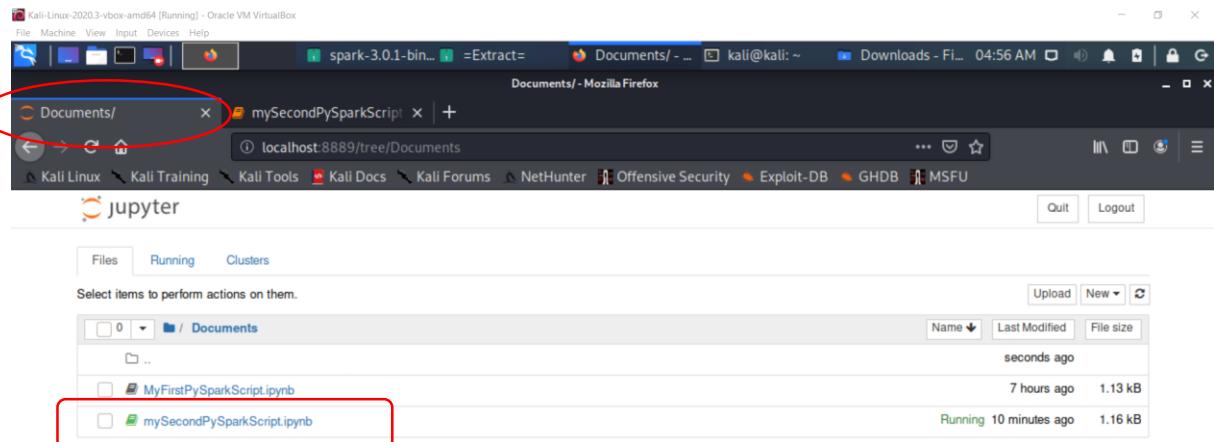
df = spark.sql("select 'spark' as hello")
df.show()
```

The output cell (In [1]) displays the result:

```
+----+
|hello|
+----+
|spark|
+----+
```

You can save your PySpark scripts to the Documents folder by clicking on the script filename, type in the new name and rename it.





The script can be found here.

Note: I do not have a Mac, hence I have not tried installing PySpark on Mac. If you need a guide on how to install PySpark on Mac, please refer to the video guide at

<https://medium.com/@GalarnykMichael/install-spark-on-ubuntu-pyspark-231c45677de0>

Or a step-by-step instructions on the installation by Kevin Vecmanis at

<https://kevinvecmanis.io/python/pyspark/install/2019/05/31/Installing-Apache-Spark.html>

## Installing ipython3 and jupyter separately:

You do the following steps only if you do not want to use Anaconda.

```
kali㉿kali:~$ sudo apt-get -y install ipython3
Reading package lists ... Done
Building dependency tree
Reading state information ... Done
The following additional packages will be installed:
  python3-ipython
Suggested packages:
  python-ipython-doc
The following packages will be upgraded:
  ipython3 python3-ipython
2 upgraded, 0 newly installed, 0 to remove and 1176 not upgraded.
Need to get 540 kB of archives.
After this operation, 8,192 B of additional disk space will be used.
Get:1 http://kali.cs.nctu.edu.tw/kali kali-rolling/main amd64 ipython3 all 7.19.0-3 [24.2 kB]
Get:2 http://kali.cs.nctu.edu.tw/kali kali-rolling/main amd64 python3-ipython all 7.19.0-3 [516 kB]
Fetched 540 kB in 2s (247 kB/s)
(Reading database ... 278086 files and directories currently installed.)
Preparing to unpack .../ipython3_7.19.0-3_all.deb ...
Unpacking ipython3 (7.19.0-3) over (7.16.1-1) ...
Preparing to unpack .../python3-ipython_7.19.0-3_all.deb ...
Unpacking python3-ipython (7.19.0-3) over (7.16.1-1) ...
Setting up python3-ipython (7.19.0-3) ...
Setting up ipython3 (7.19.0-3) ...
Processing triggers for man-db (2.9.3-2) ...
Processing triggers for kali-menu (2020.3.2) ...
kali㉿kali:~$ █
```

Next, we install the jupyter.

```

Requirement already satisfied: python-dateutil>=2.1 in /usr/lib/python3/dist-packages (from jupyter-client>=4.1→qtconsole→jupyter) (2.8.1)
Requirement already satisfied: pexpect>=4.3 in /usr/lib/python3/dist-packages (from ipython→jupyter-console→jupyter) (4.6.0)
Collecting nest-asyncio
  Downloading nest_asyncio-1.4.3-py3-none-any.whl (5.3 kB)
Collecting async-generator
  Downloading async_generator-1.10-py3-none-any.whl (18 kB)
Requirement already satisfied: six>1.9.0 in /usr/lib/python3/dist-packages (from bleach→nbconvert→jupyter) (1.15.0)
Requirement already satisfied: packaging in /usr/lib/python3/dist-packages (from bleach→nbconvert→jupyter) (20.3)
Requirement already satisfied: webencodings in /usr/lib/python3/dist-packages (from bleach→nbconvert→jupyter) (0.5.1)
Collecting pyprocess; os_name != nt
  Downloading pyprocess-0.7.0-py2.py3-none-any.whl (13 kB)
Collecting cffi>1.0.0
  Downloading cffi-1.14.4-cp39-cp39-manylinux1_x86_64.whl (405 kB)
|██████████| 405 kB 32.2 MB/s
Collecting pyparser
  Downloading pyparser-2.20-py2.py3-none-any.whl (112 kB)
|██████████| 112 kB 35.7 MB/s
Building wheels for collected packages: pandocfilters
  Building wheel for pandocfilters (setup.py) ... done
    Created wheel for pandocfilters: filename=pandocfilters-1.4.3-py3-none-any.whl size=7991 sha256=d44359e57213eb4aa1034740911030f85f658d18f46e7d52e03bb4504
3132aae
  Stored in directory: /root/.cache/pip/wheels/d7/2c/f8/55fc25b6130566bdddef382806ceaa33161b39a066c2de5912
Successfully built pandocfilters
ERROR: jupyterlab-pygments 0.1.2 has requirement pygments<3.2>4.1, but you'll have pygments 2.3.1 which is incompatible.
ERROR: nbconvert 6.0.7 has requirement pygments>3.4.1, but you'll have pygments 2.3.1 which is incompatible.
Installing collected packages: qtpy, pymq, jupyter-client, ipykernel, qtconsole, jupyter-console, jupyterlab-pygments, testpath, nest-asyncio, async-generator, nbclient, pandocfilters, entrypoints, bleach, defusedxml, nbconvert, Send2Trash, pyprocess, terminado, prometheus-client, pyparser, cffi, argon2-cffi, notebook, widgetsnbextension, jupyterlab-widgets, ipywidgets, jupyter
Successfully installed Send2Trash-1.5.0 argon2-cffi-20.1.0 async-generator-1.10 bleach-3.2.1 cffi-1.14.4 defusedxml-0.6.0 entrypoints-0.3 ipykernel-5.4.3 iipywidgets-7.6.3 jupyter-1.0.0 jupyter-client-6.1.11 jupyter-console-6.2.0 jupyterlab-pygments-0.1.2 jupyterlab-widgets-1.0.0 nbclient-0.5.1 nbconvert-6.0.7 nest-asyncio-1.4.3 notebook-6.2.0 pandocfilters-1.4.3 prometheus-client-0.9.0 pyprocess-0.7.0 pyparser-2.20 pyzmq-21.0.1 qtconsole-5.0.1 qtpy-1.9.0 termindo-0.9.2 testpath-0.4.4 widgetsnbextension-3.5.1
kali:kali:~$
```

Alternatively, you can also install jupyter using pip3, if you have already install python3 successfully.

```

kali:kali:~$ pip3 install jupyter
Requirement already satisfied: jupyter in ./anaconda3/lib/python3.8/site-packages (1.0.0)
Requirement already satisfied: jupyter-console in ./anaconda3/lib/python3.8/site-packages (from jupyter) (6.2.0)
Requirement already satisfied: qtconsole in ./anaconda3/lib/python3.8/site-packages (from jupyter) (4.7.7)
Requirement already satisfied: ipykernel in ./anaconda3/lib/python3.8/site-packages (from jupyter) (5.3.4)
Requirement already satisfied: notebook in ./anaconda3/lib/python3.8/site-packages (from jupyter) (6.1.4)
Requirement already satisfied: ipywidgets in ./anaconda3/lib/python3.8/site-packages (from jupyter) (7.5.1)
Requirement already satisfied: nbconvert in ./anaconda3/lib/python3.8/site-packages (from jupyter) (6.0.7)
Requirement already satisfied: tornado>4.2.0 in ./anaconda3/lib/python3.8/site-packages (from ipykernel→jupyter) (6.0.4)
Requirement already satisfied: traitlets>4.1.0 in ./anaconda3/lib/python3.8/site-packages (from ipykernel→jupyter) (5.0.5)
Requirement already satisfied: ipython>5.0.0 in ./anaconda3/lib/python3.8/site-packages (from ipykernel→jupyter) (7.19.0)
Requirement already satisfied: jupyter-client in ./anaconda3/lib/python3.8/site-packages (from ipykernel→jupyter) (6.1.7)
Requirement already satisfied: setuptools>18.5 in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (50.3.1.post20201107)
Requirement already satisfied: prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>2.0.0 in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (3.0.8)
Requirement already satisfied: pygments in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (2.7.2)
Requirement already satisfied: pickleshare in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (0.7.5)
Requirement already satisfied: decorator in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (4.4.2)
Requirement already satisfied: backcall in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (0.2.0)
Requirement already satisfied: jedi>0.10 in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (0.17.1)
Requirement already satisfied: parso<0.8.0,≥0.7.0 in ./anaconda3/lib/python3.8/site-packages (from ipython>5.0.0→ipykernel→jupyter) (4.8.0)
Requirement already satisfied: pyprocess≥0.5 in ./anaconda3/lib/python3.8/site-packages (from pexpect>4.3→ipython>5.0.0→ipykernel→jupyter) (0.6.0)
Requirement already satisfied: wcdwidth in ./anaconda3/lib/python3.8/site-packages (from prompt-toolkit!=3.0.0,!=3.0.1,<3.1.0,>2.0.0→ipython>5.0.0→ipykernel→jupyter) (0.2.5)
Requirement already satisfied: ipython-genutils in ./anaconda3/lib/python3.8/site-packages (from traitlets>4.1.0→ipykernel→jupyter) (0.2.0)
Requirement already satisfied: widgetsnbextension->3.5.0 in ./anaconda3/lib/python3.8/site-packages (from ipywidgets→jupyter) (3.5.1)
Requirement already satisfied: nbformat≥4.2.0 in ./anaconda3/lib/python3.8/site-packages (from ipywidgets→jupyter) (5.0.8)
Requirement already satisfied: jupyter-core in ./anaconda3/lib/python3.8/site-packages (from nbformat≥4.2.0→ipywidgets→jupyter) (4.6.3)
Requirement already satisfied: jsonschema≠2.5.0,≥2.4 in ./anaconda3/lib/python3.8/site-packages (from nbformat≥4.2.0→ipywidgets→jupyter) (3.2.0)
Requirement already satisfied: six≥1.11.0 in ./anaconda3/lib/python3.8/site-packages (from jsonschema≠2.5.0,≥2.4→nbformat≥4.2.0→ipywidgets→jupyter) (1.15.0)
Requirement already satisfied: attrs≥17.4.0 in ./anaconda3/lib/python3.8/site-packages (from jsonschema≠2.5.0,≥2.4→nbformat≥4.2.0→ipywidgets→jupyter)
```