# Pyspark installation in windows

In this video I am explaining how to install Pyspark in windows. How to create environment to practice pyspark. Explaining in this document.

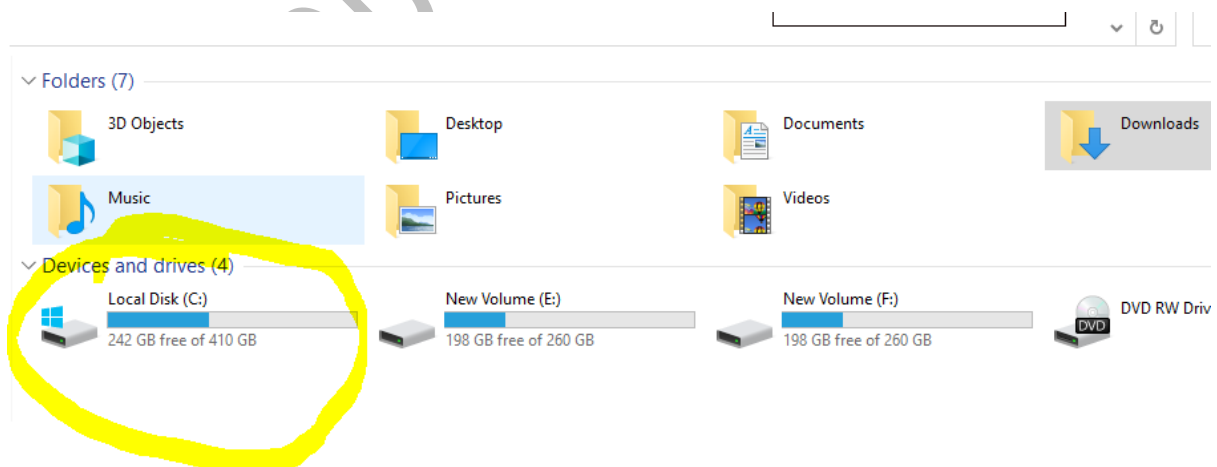If you want to install Pyspark, first you must download these three softwares from internet

https://dlcdn.apache.org/spark/spark-3.1.2/spark-3.1.2-bin-hadoop3.2.tgz

https://archive.apache.org/dist/hadoop/core/hadoop-3.2.2/hadoop-3.2.2.tar.gz

https://github.com/aitraining/datasets/blob/main/hadoop-dependencies-3.2.2.zip

https://repo.anaconda.com/archive/Anaconda3-2021.05-Windows-x86_64.exe

https://www.jetbrains.com/pycharm/download/

https://www.oracle.com/in/java/technologies/javase/jdk11-archive-downloads.html
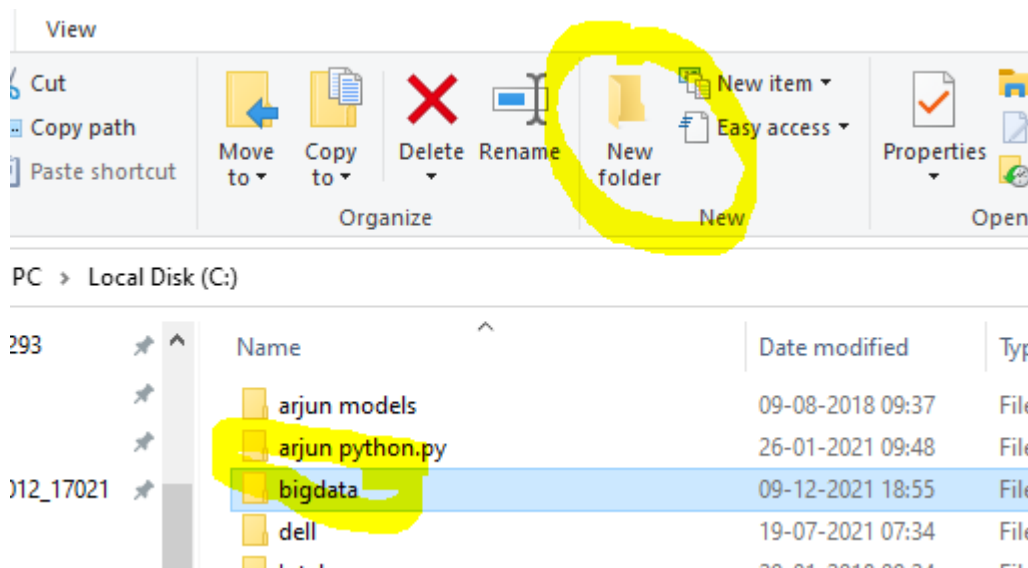
In your computer you have multiple drives like C, D , E where you have maximum space in that drive create a folder called bigdata.
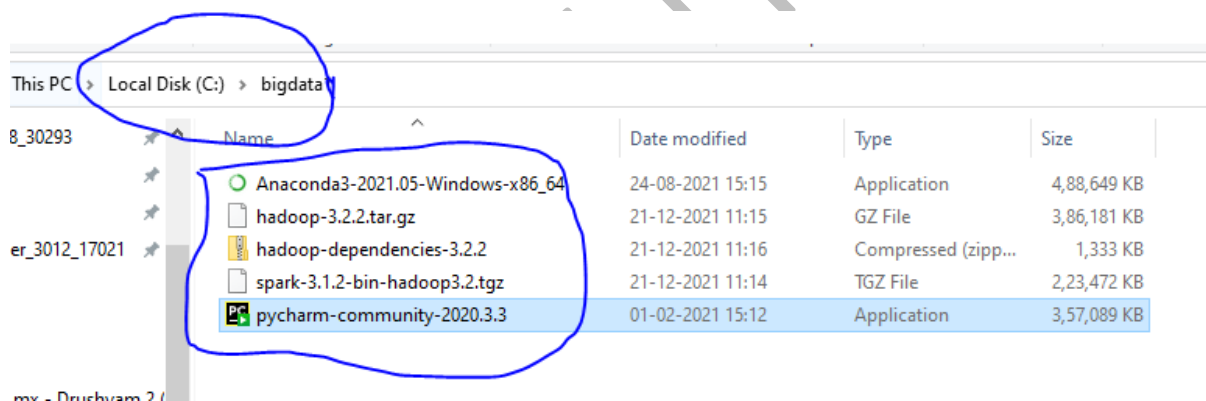
Important note: Don't create "big data" must create "bigdata" without any space. Otherwise its create new problems.
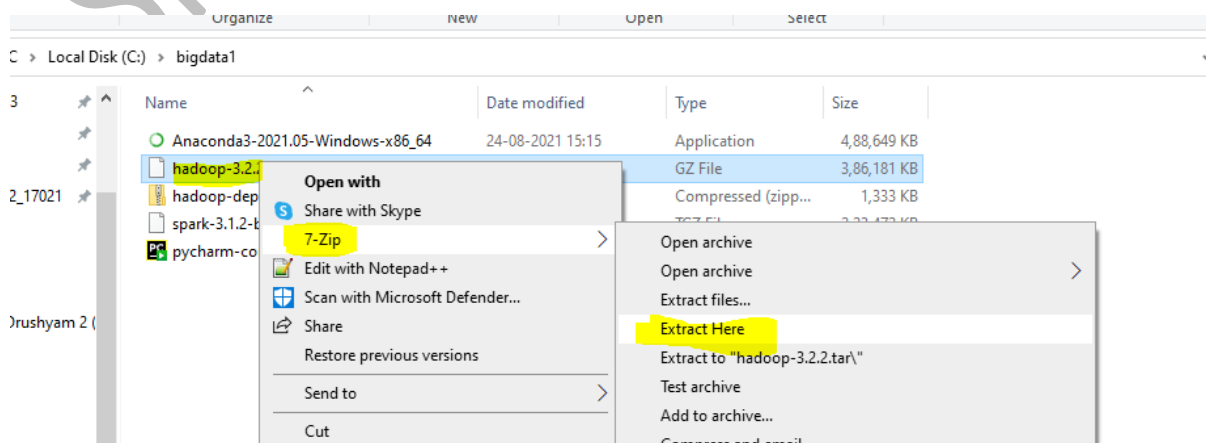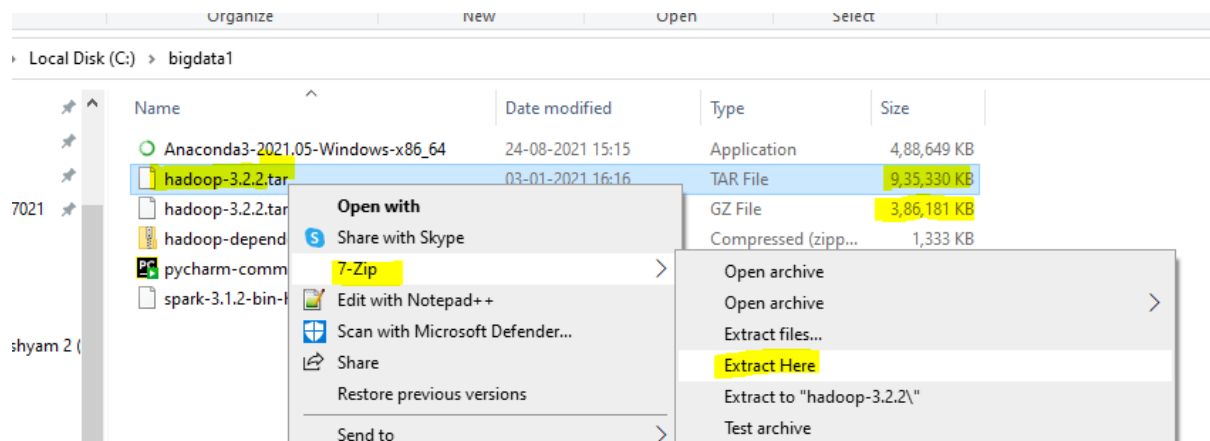
Now.

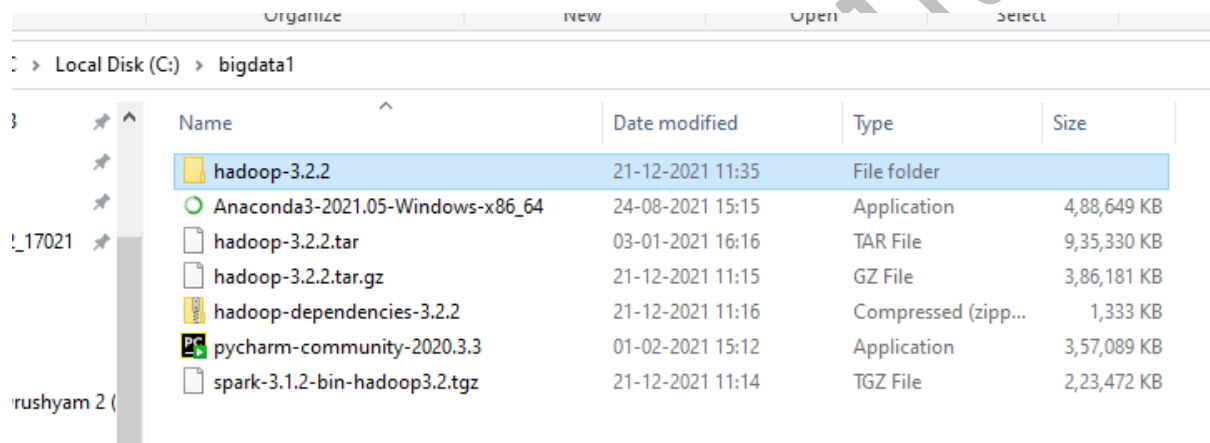After download above softwares place these folders in bigdata folders like this



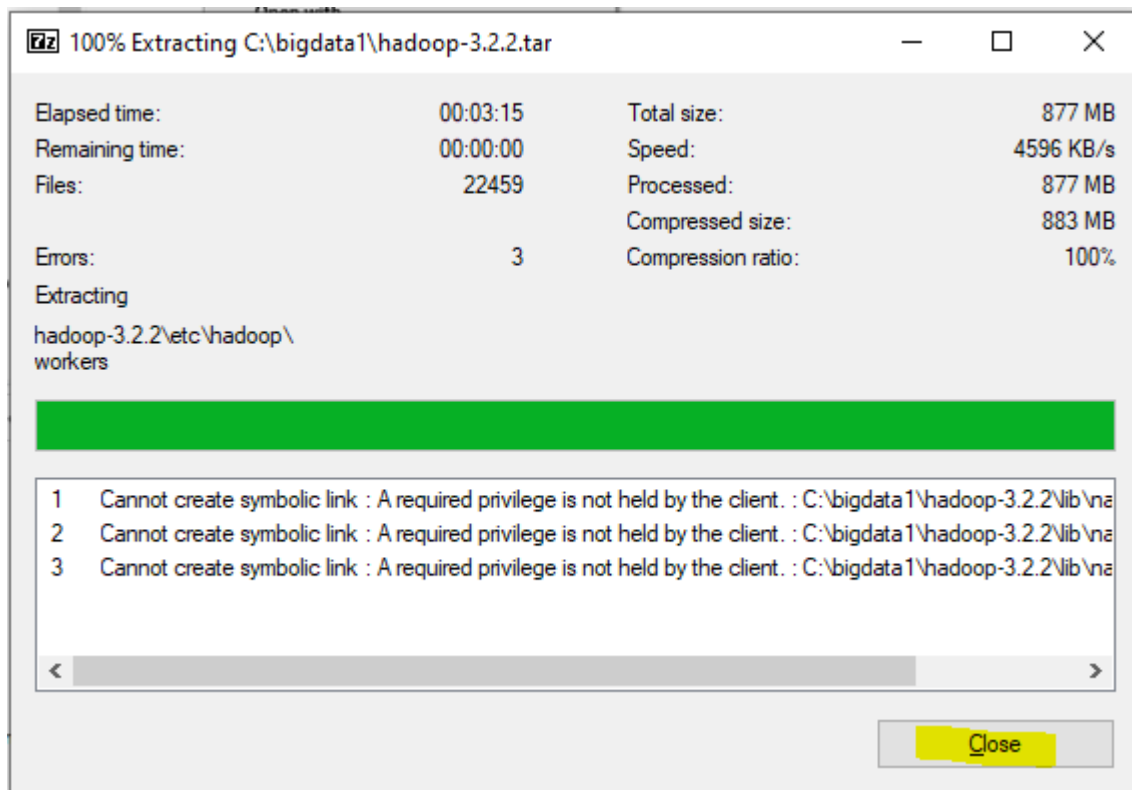After download please extract all folders. Right click and .. 7zip> extract here.

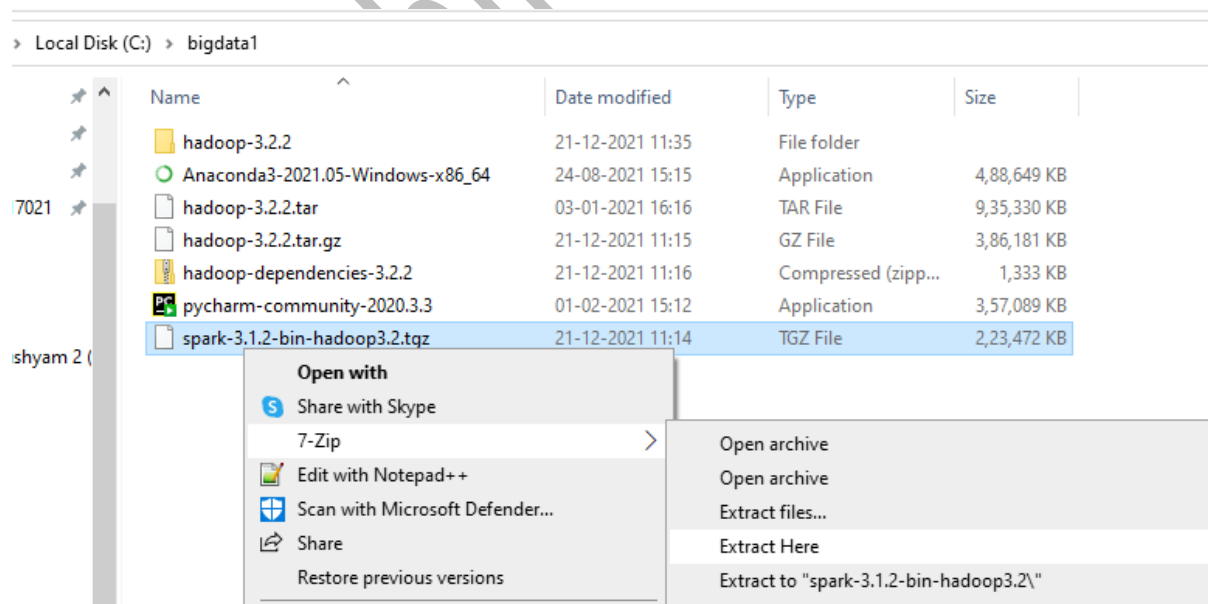After extract you will get another zip folder. Please extract it again. Please check size its changed



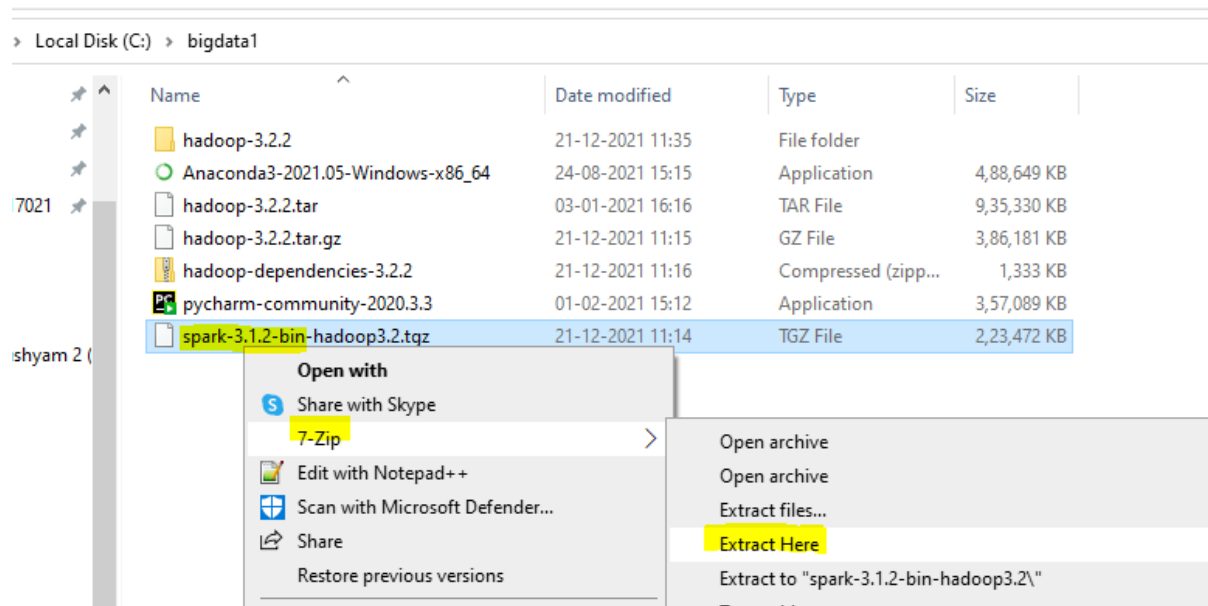Now you will get folder like this.



After 5 minutes you will get a small warning like this. Please ignore, simply click close.
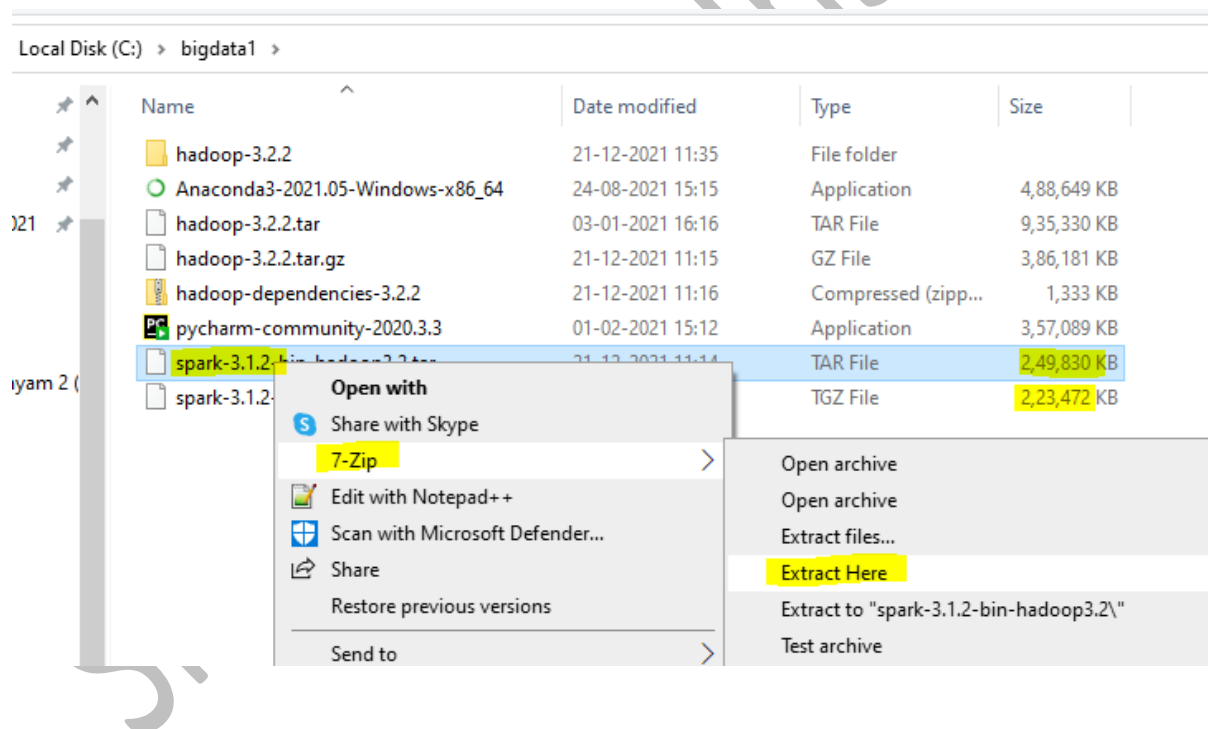
## 100% Extracting C:\bigdata1\hadoop-3.2.2.tar

| | | | |
|---|---|---|---|
| Elapsed time: | 00:03:15 | Total size: | 877 MB |
| Remaining time: | 00:00:00 | Speed: | 4596 KB/s |
| Files: | 22459 | Processed: | 877 MB |
| | | Compressed size: | 883 MB |
| Errors: | 3 | Compression ratio: | 100% |

Extracting

hadoop-3.2.2\etc\hadoop\
workers

| 1 | Cannot create symbolic link : A required privilege is not held by the client. : C:\bigdata1\hadoop-3.2.2\lib\na |
|---|---|
| 2 | Cannot create symbolic link : A required privilege is not held by the client. : C:\bigdata1\hadoop-3.2.2\lib\na |
| 3 | Cannot create symbolic link : A required privilege is not held by the client. : C:\bigdata1\hadoop-3.2.2\lib\na |

Close

Same way extract Spark folder as well.



Local Disk (C:) > bigdata1

| Name | Date modified | Type | Size |
|---|---|---|---|
| hadoop-3.2.2 | 21-12-2021 11:35 | File folder | |
| Anaconda3-2021.05-Windows-x86_64 | 24-08-2021 15:15 | Application | 4,88,649 KB |
| hadoop-3.2.2.tar | 03-01-2021 16:16 | TAR File | 9,35,330 KB |
| hadoop-3.2.2.tar.gz | 21-12-2021 11:15 | GZ File | 3,86,181 KB |
| hadoop-dependencies-3.2.2 | 21-12-2021 11:16 | Compressed (zipp... | 1,333 KB |
| pycharm-community-2020.3.3 | 01-02-2021 15:12 | Application | 3,57,089 KB |
| spark-3.1.2-bin-hadoop3.tqz | 21-12-2021 11:14 | TGZ File | 2,23,472 KB |

**Open with**
- Share with Skype
- 7-Zip  >
  - Open archive
  - Open archive
  - Extract files...
  - Extract Here
  - Extract to "spark-3.1.2-bin-hadoop3.2\"
- Edit with Notepad++
- Scan with Microsoft Defender...
- Share
- Restore previous versions

| Name | Date modified | Type | Size |
|---|---|---|---|
| hadoop-3.2.2 | 21-12-2021 11:35 | File folder | |
| Anaconda3-2021.05-Windows-x86_64 | 24-08-2021 15:15 | Application | 4,88,649 KB |
| hadoop-3.2.2.tar | 03-01-2021 16:16 | TAR File | 9,35,330 KB |
| hadoop-3.2.2.tar.gz | 21-12-2021 11:15 | GZ File | 3,86,181 KB |
| hadoop-dependencies-3.2.2 | 21-12-2021 11:16 | Compressed (zipp... | 1,333 KB |
| pycharm-community-2020.3.3 | 01-02-2021 15:12 | Application | 3,57,089 KB |
| spark-3.1.2-bin-hadoop3.2.tgz | 21-12-2021 11:14 | TGZ File | 2,23,472 KB |

Open with
Share with Skype
7-Zip  >
Edit with Notepad++
Scan with Microsoft Defender...
Share
Restore previous versions

Open archive
Open archive
Extract files...
Extract Here
Extract to "spark-3.1.2-bin-hadoop3.2\"

You will get another file. Please extract again

| Name | Date modified | Type | Size |
|---|---|---|---|
| hadoop-3.2.2 | 21-12-2021 11:35 | File folder | |
| Anaconda3-2021.05-Windows-x86_64 | 24-08-2021 15:15 | Application | 4,88,649 KB |
| hadoop-3.2.2.tar | 03-01-2021 16:16 | TAR File | 9,35,330 KB |
| hadoop-3.2.2.tar.gz | 21-12-2021 11:15 | GZ File | 3,86,181 KB |
| hadoop-dependencies-3.2.2 | 21-12-2021 11:16 | Compressed (zipp... | 1,333 KB |
| pycharm-community-2020.3.3 | 01-02-2021 15:12 | Application | 3,57,089 KB |
| spark-3.1.2-bin-hadoop3.2.tar | 21-12-2021 11:14 | TAR File | 2,49,830 KB |
| spark-3.1.2- | | TGZ File | 2,23,472 KB |

Open with
Share with Skype
7-Zip  >
Edit with Notepad++
Scan with Microsoft Defender...
Share
Restore previous versions
Send to  >

Open archive
Open archive
Extract files...
Extract Here
Extract to "spark-3.1.2-bin-hadoop3.2\"
Test archive

**Sreyobhilashi Institute Hyderabad ... +91-8500002025**

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| hadoop-3.2.2 | 03-01-2021 15:41 | File folder | |
| spark-3.1.2-bin-hadoop3.2 | 21-12-2021 11:41 | File folder | |
| Anaconda3-2021.05-Windows-x86_64 | 24-08-2021 15:15 | Application | 4,88,649 KB |
| hadoop-3.2.2.tar | 03-01-2021 16:16 | TAR File | 9,35,330 KB |
| hadoop-3.2.2.tar.gz | 21-12-2021 11:15 | GZ File | 3,86,181 KB |
| hadoop-dependencies-3.2.2 | 21-12-2021 11:16 | Compressed (zipp... | 1,333 KB |
| pycharm-community-2020.3.3 | 01-02-2021 15:12 | Application | 3,57,089 KB |
| spark-3.1.2-bin-hadoop3.2.tar | 21-12-2021 11:14 | TAR File | 2,49,830 KB |
| spark-3.1.2-bin-hadoop3.2.tgz | 21-12-2021 11:14 | TGZ File | 2,23,472 KB |

Now configure in environment variable, that's y copy this path

This PC > Local Disk (C:) > bigdata1 > spark-3.1.2-bin-hadoop3.2

| Name | Date modified | Type | Siz |
|------|---------------|------|-----|
| bin | 24-05-2021 10:15 | File folder | |
| conf | 24-05-2021 10:15 | File folder | |
| data | 24-05-2021 10:15 | File folder | |
| examples | 24-05-2021 10:15 | File folder | |
| jars | 24-05-2021 10:15 | File folder | |
| kubernetes | 24-05-2021 10:15 | File folder | |
| licenses | 24-05-2021 10:15 | File folder | |

Now Start button > env > Environment variable >

After this step



Finally click ok

Similarly Hadoop too



Copy this path and paste in environment variable

Next Same way path also just paste in path

Edit environment variable

| | |
|---|---|
| C:\Program Files\MySQL\MySQL Shell 8.0\bin\ | |
| %USERPROFILE%\AppData\Local\Microsoft\WindowsApps | |
| C:\Users\dell\AppData\Local\Programs\Microsoft VS Code\bin | |
| %IntelliJ IDEA Community Edition% | |
| C:\bigdata\Java\jdk1.8.0_291\bin | |
| C:\bigdata\spark-2.4.7-bin-hadoop2.7\bin | |
| C:\bigdata\spark-2.4.7-bin-hadoop2.7\sbin | |
| %PyCharm Community Edition% | 2 |
| C:\bigdata\hadoop-3.2.2\bin | 3 |
| C:\bigdata\hadoop-3.2.2\sbin | |

1 New

Edit

Browse...

Delete

Move Up

Move Down

Edit text...

4 OK    Cancel

5 OK    Cancel

New

C:\bigdata1\hadoop-3.2.2 …

Paste this Hadoop-dependencies-322 zip file in Hadoop 3.2.2 folder … now right click and extract it.



You will get replace all , Plese click replace all

Name

- bin
- etc
- include
- lib
- libexec
- sbin
- share
- hadoop-dependen
- LICENSE
- NOTICE
- README

**Confirm File Replace** ✕

Destination folder already contains processed file.

Would you like to replace the existing file

C:\bigdata\hadoop-3.2.2\etc\hadoop\
capacity-scheduler.xml
9213 bytes
Modified: 2021-01-03 15:24:47

with this one?

etc\hadoop\
capacity-scheduler.xml
9213 bytes
Modified: 2021-01-03 15:24:47

| Yes | Yes to All | Auto Rename |
|-----|-----------|-------------|
| No | No to All | Cancel |

Elapsed ti
Remaining
Files:

Extracting
etc\hado

5095 KB

0
0

Now somethime if you want to change namenode and datanode place please modify these

> Local Disk (C:) > bigdata > hadoop-3.2.2

| Name | Date modified | Type | Size |
|------|---------------|------|------|
| bin | 14-10-2021 15:29 | File folder | |
| dfs | 24-11-2021 16:24 | File folder | |
| etc | 14-10-2021 15:26 | File folder | |
| include | 21-12-2021 12:00 | File folder | |
| lib | 21-12-2021 12:00 | File folder | |
| libexec | 21-12-2021 12:00 | File folder | |
| sbin | 21-12-2021 12:00 | File folder | |
| share | 21-12-2021 12:23 | File folder | |
| hadoop-dependencies-3.2.2 | 21-12-2021 11:16 | Compressed (zipp... | 1,333 KB |
| LICENSE | 05-12-2020 20:39 | Text Document | 148 KB |
| NOTICE | 05-12-2020 20:39 | Text Document | 22 KB |
| README | 05-12-2020 20:39 | Text Document | 2 KB |

7021

shyam 2 (

Click Hadoop

PC > Local Disk (C:) > bigdata > hadoop-3.2.2 > etc

| Name | Date modified |
|------|---------------|
| hadoop | 14-10-2021 15:26 |

93

Now right click core.site.xml edit with notepad++

: > Local Disk (C:) > bigdata > hadoop-3.2.2 > etc > hadoop >

| Name | Date modified | Type |
|------|---------------|------|
| shellprofile.d | 14-10-2021 15:26 | File folder |
| capacity-scheduler | 03-01-2021 15:24 | XML Docume |
| configuration | 03-01-2021 15:27 | XSL Styleshee |
| container-executor.cfg | 03-01-2021 15:24 | CFG File |
| core-site | 22-11-2021 20:16 | XML Docume |
| hadoop-env | | ws Cor |
| hadoop-env | | |
| hadoop-metrics2.propertie | | ERTIES |
| hadoop-policy | | ocume |
| hadoop-user-functions.sh. | | PLE File |
| hdfs-site | | ocume |
| httpfs-env | | |
| httpfs-log4j.properties | | ERTIES |

Open

Edit

S  Share with Skype

7-Zip  >

Edit with Notepad++

Scan with Microsoft Defender...

Share

Open with  >

If you want to change namenodde or datanode other information pls change these

```xml
<configuration>
    <property>
        <name>fs.defaultFS</name>
        <value>hdfs://localhost:9000</value>
    </property>
    <property>
        <name>hadoop.tmp.dir</name>
        <value>/C:/bigdata/hadoop-3.2.2/dfs/tempdir</value>
    </property>
        <property>
        <name>fs.trash.interval</name>
        <value>1440</value>
    </property>



</configuration>
```

In hdfs-site.xml also pls change



If you want to change pls change

```
<configuration>
 <property>
        <name>dfs.replication</name>
        <value>1</value>
    </property>
    <property>
        <name>dfs.namenode.name.dir</name>
        <value>/C:/bigdata/hadoop-3.2.2/dfs/namenode</value>
    </property>
    <property>
        <name>dfs.datanode.data.dir</name>
        <value>/C:/bigdata/hadoop-3.2.2/dfs/datanode</value>
    </property>


</configuration>
```

Now save core-site.xml and hdfs-site.xml

That's it 99.9% Hadoop installation completed. Verification do it later

Please note if you want to do anything Java 11 is mandatory so install **Java 11 and Anaconda**

Download java from this link

https://www.oracle.com/in/java/technologies/javase/jdk11-archive-downloads.html

Now you will get popup like this enter these credentials

[venumssi@gmail.om](venumssi@gmail.om)

Opassword.1

## Oracle account sign in

Username

**venumssi@gmail.com** ℹ

Password

**••••••••••** ℹ

**Sign in**

Need help?

## Don't have an Oracle Account?

Create Account

Now after download click to install java 11

You will get like this

Install anyway next



Click yes

Next



Click Change

Next

Now rename this as bigdata instead of Program Files



Now click Ok

Now click Next

Now java 11 installing in ur system wait sometime



Now click Close

Now java installation completed, I recommend configure in environment variable

Means



Copy this path and paste inenvironment variable

JAVA_HOME must be capital (case sensitive)

Environment Variables                                                            ✕

User variables for dell

| Variable | Value |
|---|---|
| FREI0R_PATH | C:\PROGRA~2\ACETHI~1\ACETHI~1\frei0r;C:\Program Files (x86)\... |
| HADOOP_HOME | C:\bigdata\hadoop-3.2.2 |
| IntelliJ IDEA Community Edit... | C:\bigdata\IntelliJ IDEA Community Edition 2021.1\bin; |
| OneDrive | C:\Users\dell\OneDrive |
| Path | C:\Program Files\MySQL\MySQL Shell 8.0\bin\;C:\Users\dell\AppD... |
| PyCharm Community Edition | C:\bigdata\PyCharm Community Edition 2021.2.3\bin; |

ew User Variable

ariable name:     JAVA_HOME

ariable value:    C:\bigdata\Java\jdk-11.0.12

| Browse Directory... | Browse File... | | OK | Cancel |
|---|---|---|---|---|

| | |
|---|---|
| DriverData | C:\Windows\System32\Drivers\DriverData |
| NUMBER_OF_PROCESSORS | 4 |
| OS | Windows_NT |
| Path | C:\Program Files\Common Files\Oracle\Java\javapath;C:\Program ... |
| PATHEXT | .COM;.EXE;.BAT;.CMD;.VBS;.VBE;.JS;.JSE;.WSF;.WSH;.MSC;.PY;.PYW |
| PROCESSOR_ARCHITECTURE | AMD64 |

| New... | Edit... | Delete |
|---|---|---|

|  | OK | Cancel |
|---|---|---|

Configure in path as well

Testing purpose java installed or not click on start … cmd …

Java -version



Please install anaconda as well

Local Disk (C:) > bigdata1

| Name |
|------|
| hadoop-3.2.2 |
| spark-3.1.2-bin-hadoop3 |
| Anaconda3-2021.05-Win     **1** |
| hadoop-3.2.2.tar |
| hadoop-3.2.2.tar.gz |
| hadoop-dependencies-3 |
| pycharm-community-20 |
| spark-3.1.2-bin-hadoop3 |
| spark-3.1.2-bin-hadoop3 |

7021

shyam 2 (

**Anaconda3 2021.05 (64-bit) Setup**

**Welcome to Anaconda3 2021.05 (64-bit) Setup**

Setup will guide you through the installation of Anaconda3 2021.05 (64-bit).

It is recommended that you close all other applications before starting Setup. This will make it possible to update relevant system files without having to reboot your computer.

Click Next to continue.

**2**

Next >    Cancel

---

**Anaconda3 2021.05 (64-bit) Setup**

**License Agreement**

Please review the license terms before installing Anaconda3 2021.05 (64-bit).

Press Page Down to see the rest of the agreement.

```
=====================================
End User License Agreement - Anaconda Individual Edition
=====================================

Copyright 2015-2021, 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).
```

If you accept the terms of the agreement, click I Agree to continue. You must accept the agreement to install Anaconda3 2021.05 (64-bit).

Anaconda, Inc.

< Back    I Agree    Cancel

Anaconda3 2021.05 (64-bit) Setup

**Select Installation Type**
Please select the type of installation you would like to perform for
Anaconda3 2021.05 (64-bit).

Install for:

◉ Just Me (recommended)

○ All Users (requires admin privileges)

Anaconda, Inc.

< Back    Next >    Cancel

Anaconda3 2021.05 (64-bit) Setup

**Choose Install Location**
Choose the folder in which to install Anaconda3 2021.05 (64-bit).

Setup will install Anaconda3 2021.05 (64-bit) in the following folder. To install in a different
folder, click Browse and select another folder. Click Next to continue.

Destination Folder

C:\Users\dell\anaconda3     Browse...

Space required: 2.9GB
Space available: 242.1GB

Anaconda, Inc.

< Back    Next >    Cancel

Its take a lot of time

After 30 minutes anaconda installed

Finally you will get like this



Click Next

Finally uncheck and click finish



Testing purpose anaconda installed or not verify like this

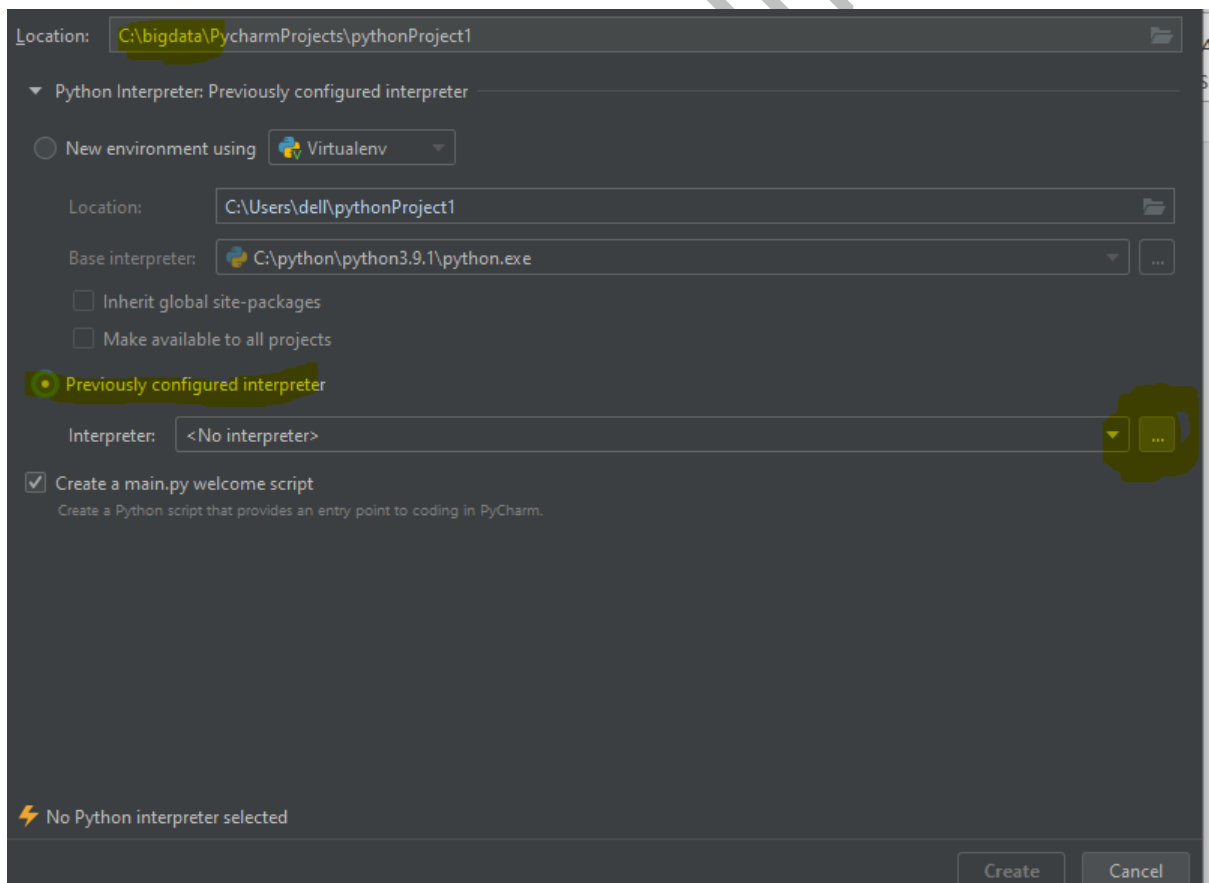If you get like this Anaconda Prompt (anaconda3) .. congrats anaconda installed.

Now python 3.6 recommended that's y execute this command.

conda create -n py36 python=3.6



Sometime you will get like this enter Y

```
The following NEW packages will be INSTALLED:

  certifi          pkgs/main/noarch::certifi-2020.6.20-pyhd3eb1b0_3
  pip              pkgs/main/win-64::pip-21.2.2-py36haa95532_0
  python           pkgs/main/win-64::python-3.6.13-h3758d61_0
  setuptools       pkgs/main/win-64::setuptools-58.0.4-py36haa95532_0
  sqlite           pkgs/main/win-64::sqlite-3.36.0-h2bbff1b_0
  vc               pkgs/main/win-64::vc-14.2-h21ff451_1
  vs2015_runtime   pkgs/main/win-64::vs2015_runtime-14.27.29016-h5e58377_2
  wheel            pkgs/main/noarch::wheel-0.37.0-pyhd3eb1b0_1
  wincertstore     pkgs/main/win-64::wincertstore-0.2-py36h7fe50ca_0


Proceed ([y]/n)? y
```

After 10 minutes py36 installed you can practice python 3.6 version in py36 environment.

By default you will get like this

By default python (base) you will get 3.8 or 3.9 but I want to practice python 3.6 (pyspark support by default 3.6 that's y) activate py36 using this command.

```
#
# To activate this environment, use
#
#     $ conda activate py36
#
# To deactivate an active environment, use
#
#     $ conda deactivate


(base) C:\Users\dell>python
Python 3.8.8 (default, Apr 13 2021, 15:08:03) [MSC v.1916 64 bit (AMD64)] :: Anaconda, Inc. on win32
Type "help", "copyright", "credits" or "license" for more information.
>>>

(base) C:\Users\dell>conda activate py36

(py36) C:\Users\dell>python
Python 3.6.13 |Anaconda, Inc.| (default, Mar 16 2021, 11:37:27) [MSC v.1916 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>>
```

Now Its recommended not mandatory

%SPARK_HOME%\python

%SPARK_HOME%\python\lib\py4j-0.10.7-src.zip

```
(py36) C:\Users\dell>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
```

Within this time install pycharm

https://www.jetbrains.com/pycharm/download/

in this link download community version only

Paste that Pycharm in bigdata folder



Double click

Click Next

My recommendation this is also install in bigdata folder



Finally go to C drive bigdata folder click OK

Now you will get like this

It means pycharm u r installing in bigdata folder

All these are optional pls ignore



Click install

It will take 5 minutes maximum

Finally you will get like this



Finally check Run Pycharm Community edition and Click Finish

Now 99% pycham also installation completed.

How to practice pyspark ill guide now.

Now you will get home screen like this

File > setting



If python version wrong here change



Next

Its important step. Click py4j and pyspark.zip file



Now you will get like this click ok

If you don't prefer black colour screen(theme) change like this (its optional)



Click Appearance & behaviour

Now click Appearance ..
Under Theme ... instead of Darcula choose IntelliJ Light ... Click Ok

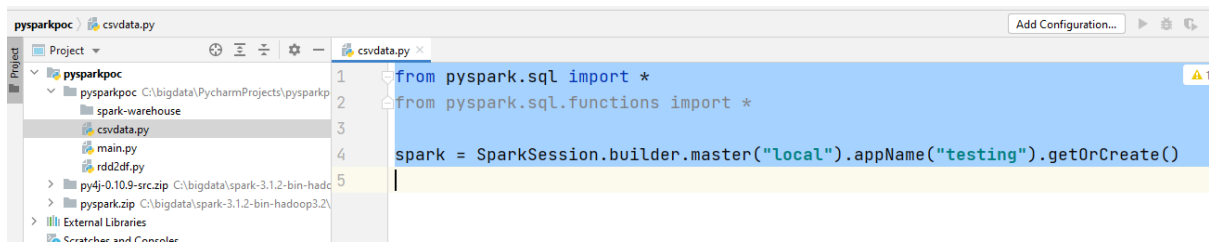You will get like this like light colour



Right click on Pysparkpoc … Click New …Choose Python File

Write something python file name like csvdata
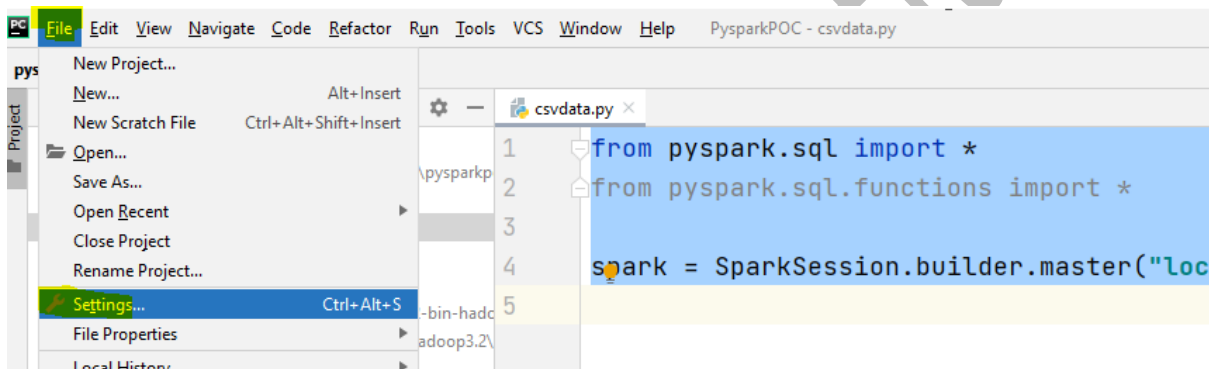


Now write something program like this

```python
from pyspark.sql import *
from pyspark.sql.functions import *

spark =
SparkSession.builder.master("local").appName("testing")
.getOrCreate()
```

These three lines ur using most frequently in any program so its recommended like this

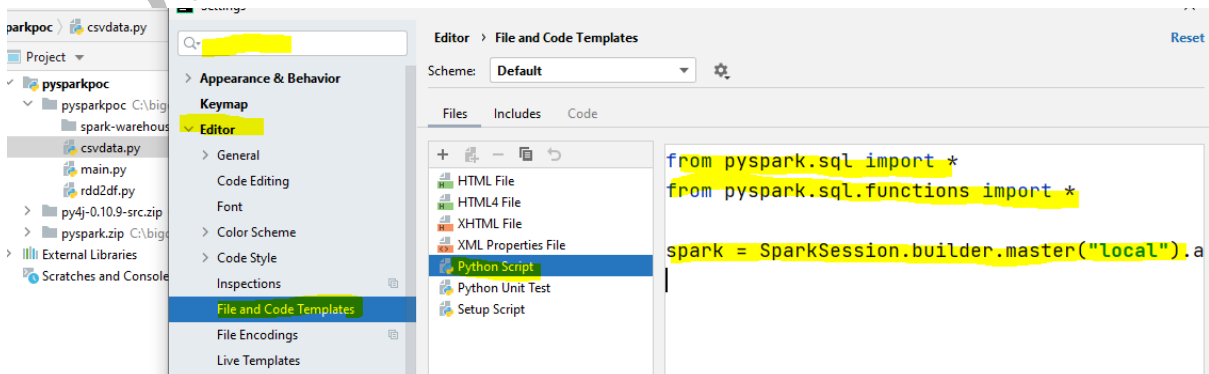(its optional)

Goto file .. setting



Now click on Editor .. Click on File and Code Templates ... Click PythonScript...

I recommend paste this code here

from pyspark.sql import *

from pyspark.sql.functions import *


spark = SparkSession.builder.master("local").appName("testing").getOrCreate()

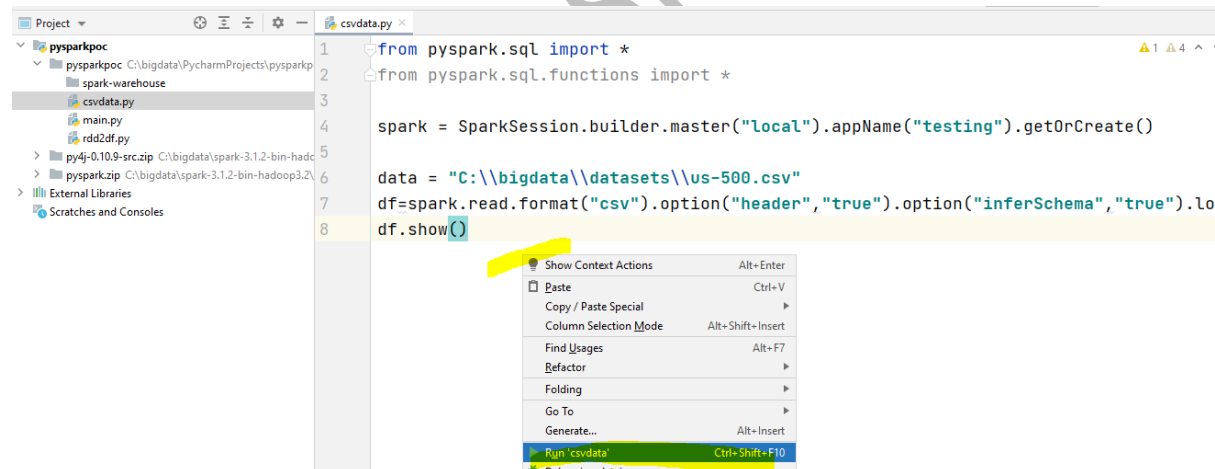If you paste here what happen, if you create any python object automatically you will get this pyspark code

Testing purpose paste this code and run this program

```python
from pyspark.sql import *
from pyspark.sql.functions import *

spark = SparkSession.builder.master("local").appName("testing").getOrCreate()

data = "C:\\bigdata\\datasets\\us-500.csv"
df=spark.read.format("csv").option("header","true").option("inferSchema","true").load(data)
df.show()
```
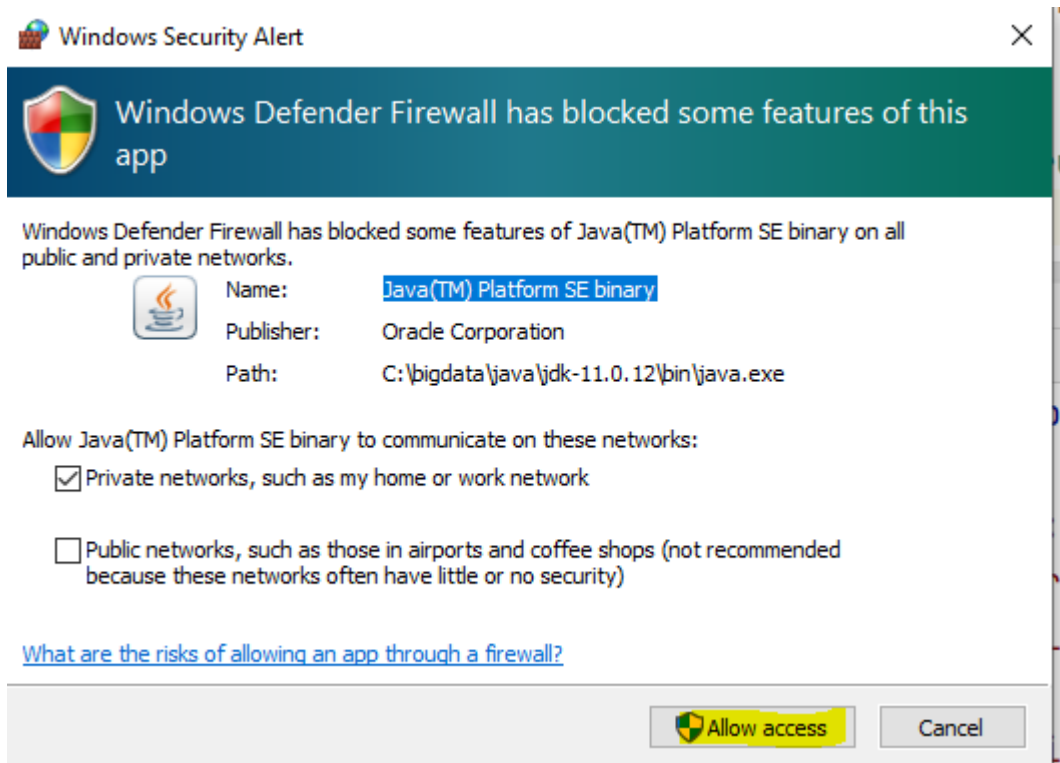
Here you must replace data location data = **"C:\\bigdata\\datasets\\us-500.csv"**
Similarly must use \\ otherwise python consider as escape character … now right click run this program



If you running first time you will get a popup like this

It means spark want to access java, you are granting permission .. Click Allow access.

Now results you will get like this



It means Spark 100% working fine.

Similarly Hadoop working or not. Open cmd anaconda or normal cmd anything ok

Pls execute this command



If you get like this Hadoop 100% formatted now Hadoop 100% working

```
921-12-21 15:57:14,473 INFO namenode.FSImage: Allocated new BlockPoolId: BP-1177568279-192.168.11.1-1640082454443
921-12-21 15:57:14,773 INFO common.Storage: Storage directory C:\bigdata\hadoop-3.2.2\dfs\namenode has been success-
 formatted.
921-12-21 15:57:14,913 INFO namenode.FSImageFormatProtobuf: Saving image file C:\bigdata\hadoop-3.2.2\dfs\namenode\
nt\fsimage.ckpt_0000000000000000000 using no compression
921-12-21 15:57:15,352 INFO namenode.FSImageFormatProtobuf: Image file C:\bigdata\hadoop-3.2.2\dfs\namenode\current
age.ckpt_0000000000000000000 of size 396 bytes saved in 0 seconds .
921-12-21 15:57:15,452 INFO namenode.NNStorageRetentionManager: Going to retain 1 images with txid >= 0
921-12-21 15:57:15,462 INFO namenode.FSImage: FSImageSaver clean checkpoint: txid=0 when meet shutdown.
921-12-21 15:57:15,462 INFO namenode.NameNode: SHUTDOWN_MSG:
************************************************************
HUTDOWN_MSG: Shutting down NameNode at DESKTOP-2IVLKVD/192.168.11.1
************************************************************/

:\Users\dell>start-dfs && start-yarn
tarting yarn daemons

:\Users\dell>
```

Now try any Hadoop commands its working

Now run start-dfs && start-yarn

You will get 4 popup boxes don't worry just minimise it.(don't close)

If you get like this



```
C:\Users\dell>start-dfs && start-yarn
starting yarn daemons

C:\Users\dell>jps
496 NameNode
13476 DataNode
12552
3816 Jps
7212 NodeManager
8124 ResourceManager
```

100% Hadoop working fine.

You can start any Hadoop command now.