# Assignment

**Installation and Setup of Apache Spark:**

**Apache PySpark:**

- Apache Spark is an open-source distributed computing system that provides a fast and general-purpose cluster-computing framework for big data processing.

- It is a data processing framework that can quickly perform processing tasks on very large data sets.

- It can also distribute data processing tasks across multiple computers, either on its own or in tandem with other distributed computing tools.

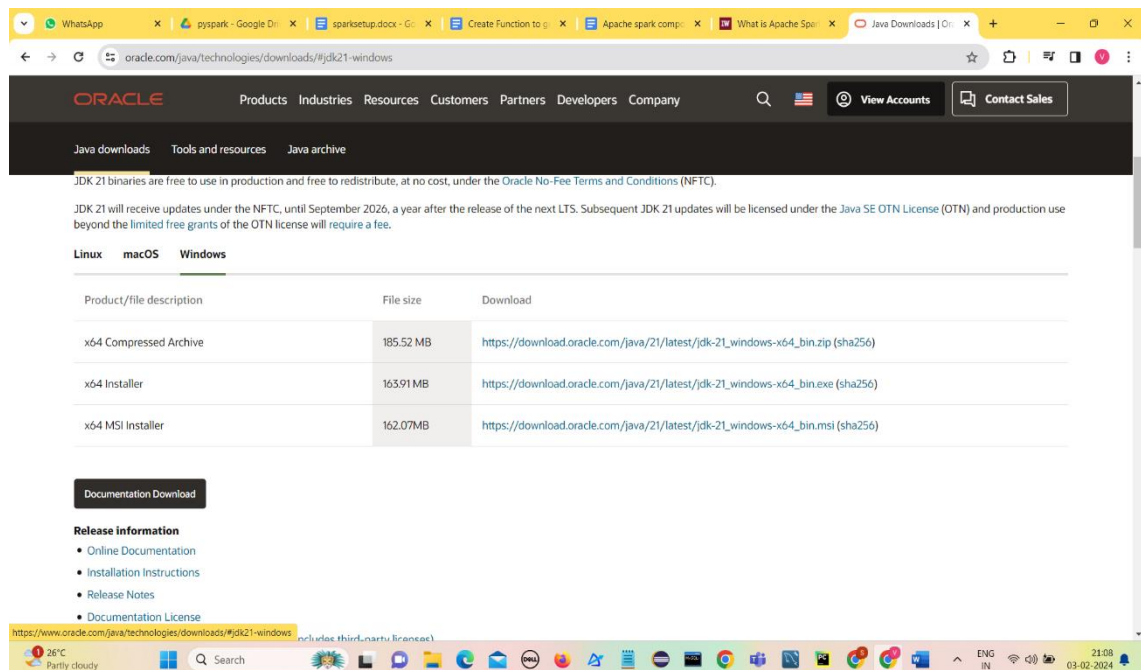    **PySpark** : It is the Python API for Apache Spark, allowing you to write Spark applications using Python.

    **Python+Spark=PySpark**

**Steps to be followed for Installation and Setup**

**1) Install Java**

- Firstly, we have to install java jdk version which is compatible to our system. It is advisable to download and install latest standard version of java**. Download Jdk from the official website oracle.com**

**Check the version of Java installed**

- After that we have to check in command prompt by typing **java –version** which displays the version of java which you downloaded.
- In means that you have successfully downloaded and installed in your system.



## 2) Install Python

- Along with java, we have to install python environment into our system. Its good to download and install latest and standard version of python.

  **Download python from official website in the following way**



- After that we have to check in command prompt by typing **python –version** which displays the version of python which you downloaded.
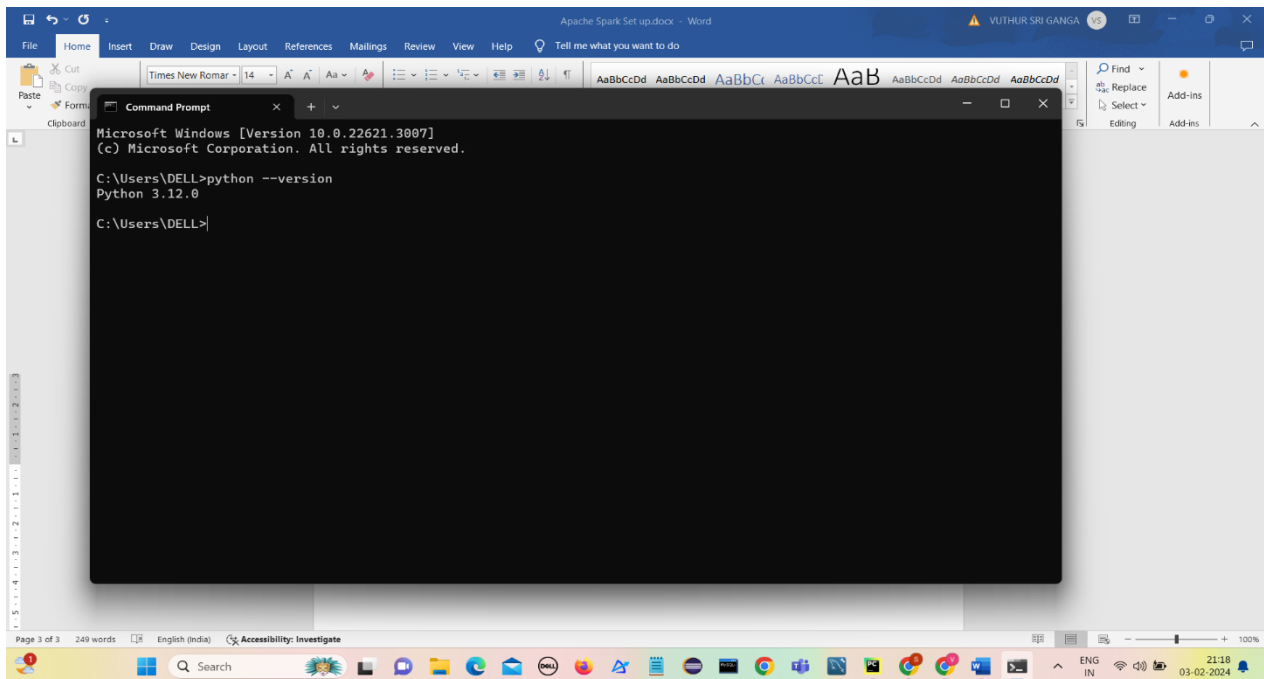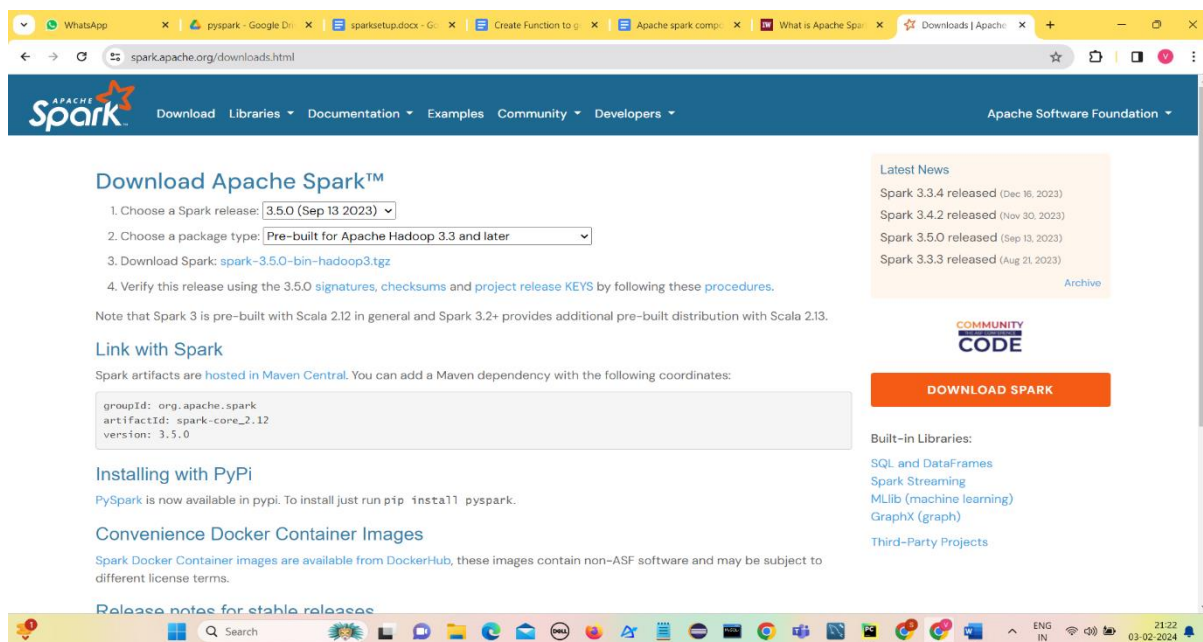- In means that you have successfully downloaded and installed in your system.

# Check the version of Python in Command Prompt



## Install Apache Spark

- Visit the Apache Spark download page.
- Choose the latest version of Spark and download the pre-built package for Hadoop. It will be a tarball (.tgz) file.
- Extract the downloaded tarball to a location on your machine.

## Install winutils from the following website

https://github.com/steveloughran/winutils
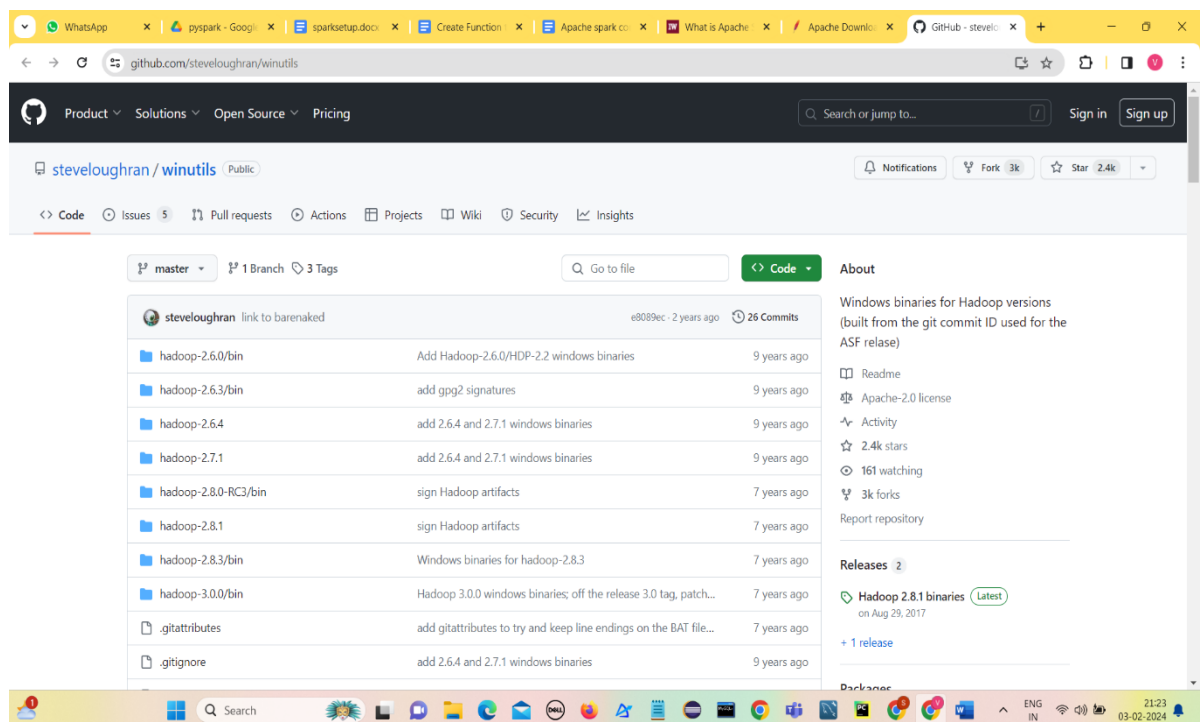
Ensure that the version of winutils should be same as that of Apache Spark.

# Set Environment Variables

**Click OK and exit.**

**To check whether spark is installed or not.**

After setting up the environment variables,we need to save all of them and have to go to the command prompt and type **spark-shell** as below.

**To check whether pyspark is installed or not.**

After setting up the environment variables,we need to save all of them and have to go to the command prompt and type **pyspark** as below.