**Name: Pradip Bochare** 

♣ Installation of Apache Spark and Setup

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

## Install Java

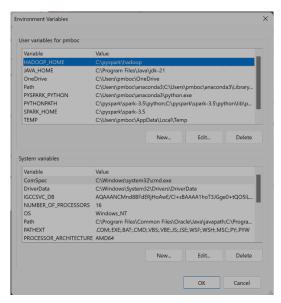
- Firstly, we have to install java jdk version which is compatible to your system. It's good to download and install latest standard version of java.
- After that we have to check in command prompt by typing java –version which displays the version of java which you downloaded.
- In means that you have successfully downloaded and installed in your system.

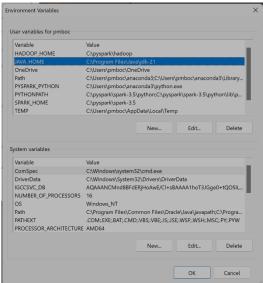
- Install Python
- Along with java, we have to install python environment into our system. It's good to download and install latest and standard version of python.
- After that we have to check in command prompt by typing python –version which displays the version of java which you downloaded.
- In means that you have successfully downloaded and installed in your system.

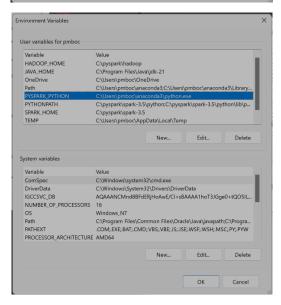


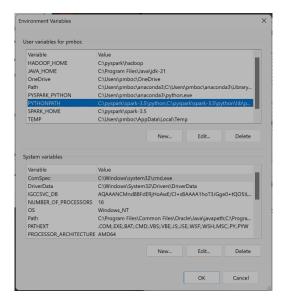
- Install Apache Spark
- Visit the Apache Spark download page.
- Choose the latest version of Spark and download the pre-built package for Hadoop. It will be a tarball (.tgz) file.
- Extract the downloaded tarball to a location on your machine.

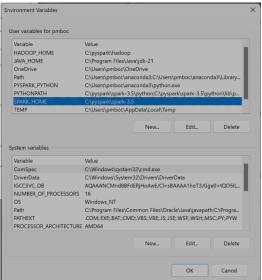
## Set Environment Variables

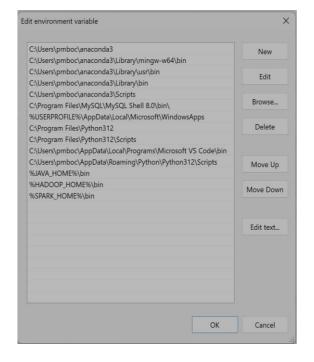




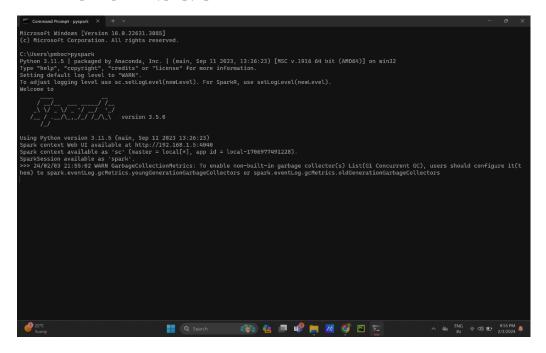




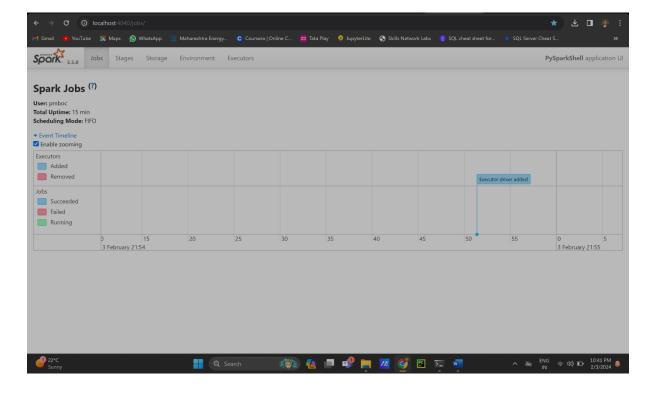




After setting up the environment variables, we need to save all of them and have to go to the command prompt and type **pyspark** as below.



- If it shows like these, you have successfully set up the Apache Spark setup.
- If you want to see pyspark web UI, you have to type in chrome as like
- http://localhost:4040/jobs/ you will get the web page as follows which indicates your pyspark is installed as follows



Aprohe Sporth.  - personal purpose cluster computing system - personal purpose cluster computing system - provide aptinized engine that supposts general - provide aptinized engine that supposts general - provide aptinized engine that supposts general - provide aptinized engine - sports (components:  Sports (core - sports (core - sports (core - standalone)  A pache Sports (core - sports core - standalone)  A pache Sports (core - sports core - sports core - delivers speed by providing in memory computation capability:  Ley features - in charge of essential I/O functionalities - significant in programming to observing role of - sports cluster  - task dispatching fault recovery - Overcornes sway of MapReduce by using in-memory - computation		Data
Spark Sql Spark stokening MLib Greathx  Spark Core  Spark Core  Spark Hadoup Mesos  Standalone MARH  A packe Spark Core  - All functionalities built on top of spark core  - delivers speed by providing in-memory computation capability.  Ley features  - in charge of essential I/O functionalities.  Significant in programming to observing role of spark cluster.  Task dispatching.  - fault recovery  - Overcomes snag of ManRodines I	*	- proteine aprilar seet
Sparty Core  Sparty Hadoup Mesos  Standalme MARH  A packe Spark Core  - All functionalities built as top of spark core  - delivers speed by providing in-memory computation capability.  Ley features  - in charge of essential I/O functionalities.  - significant in programming to observing role of spark cluster.  Task dispatching.  - fault recovery  - Overcomes snag of ManRadues 1	100	AND THE RESIDENCE OF THE PARTY
Standalone PARH  A packe Spack Cose  - All functionalities built as top of spack core  - delivers speed by providing in-memory computation capability.  Key features  - in charge of essential I/O functionalities.  - significant in programming & observing role of spark cluster.  - Task dispatching.  - fault recovery  - Overcomes snag of ManReduce 1	12,	Shaed 2dr Shaed oran 101
- All functionalities built as top of spark care  - delivers speed by providing in-memory computation capability:  Very features  - in charge of essential I/O functionalities.  - significant in programming & observing role of spark clyster:  - Task dispatching.  - fault recovery  - Overcomes snag of ManReduce I		
		- All functionalities built on top of spork corse  - delivers speed by providing in-memory computation capability.  Very features  - in charge of essential I/O functionalities.  - significant in programming & observing role of sporks cluster.  - Task dispatching.  - fault recovery  - Overcomes snag of ManReduce I

	Page No.
	1 1 1 1 1 A Johnson .
	RDD -+ Resilient distributed dataset.
	- sports and handles postitioning data across all the nodes in a clyster.
	the nodes in a city.
	two operations performed on 200
	Frans formation Action.
	Fransformation Action.  Produces new RDD from Work with actual dataset.
	\$100 EX 2000 EX
	* Apache spark SQL
	- distributed framework for structured data  processing.
	- It does not depend on API tangunge !
	- Works to access structured & semi-structured
Circ	information.
	features - cost based optimizer.
	- mid overy fault tolerance
70	- full compatibility outil entitle common every to greets
	consert of data sources.
	sparks treatenes doing
	Data France API

Date - Allows scalable, high-throughput, fault-tolerance stream processing of live data streams 1 Gathering - Basic Sources: file systems, socket connecting - Advanced sources: Kafka, Flume, Kinesis available through enter utility dosses. (b) Processing - gethered data is processed using complex algorithms. @ Data Storage - processed data is pushed out to file systems, databases and live dashboards. of Apache Spark Milib (machine Learning library) - MLIB in sparty is a scalable machine learning library that discusses both high-quality algorithm & high speed. - clustering, regression, classification collaborative filtering.

