

| Package or method | Description | Code example |
|------------------------------------|--|--|
| SparkSession.builder.getOrCreate() | Gets an existing SparkSession or, if there is no existing one, creates a new one based on the options set in this builder. | <pre>1 from pyspark.sql import SparkSession</pre> <p>Create a SparkSession or get an existing one:</p> <pre>1 spark = SparkSession.builder.appName("myApp").getOrCreate()</pre> <p>Now you can use the spark object to work with Spark</p> |
| cd | Used to move efficiently from the existing working directory to different directories on your system. | <p>Basic syntax of cd command:</p> <pre>1 cd [options]... [directory]</pre> <p>Example 1: Change directory location to folder1:</p> <pre>1 cd /usr/local/folder1</pre> <p>Example 2: Get back to the previous working directory</p> <pre>1 cd -</pre> <p>Example 3: Move up one level from the present working directory tree</p> <pre>1 cd ..</pre> |
| docker-compose | Tool for defining and running multicontainer Docker applications. It uses the YAML file to configure the services and enables us to create and start all the services from just one configuration file. | <p>Example (docker-compose.yml)</p> <pre>1 version: '3' 2 services: 3 web: 4 image: nginx:latest 5 ports: 6 - "80:80" 7 db: 8 image: postgres:latest</pre> |
| git clone | You can create a copy of a specific repository or branch within a repository. | <pre>1 git clone REPOSITORY_URL [DESTINATION_DIRECTORY]</pre> |
| import | Used to make code from one module accessible in another. Python imports are crucial for a successful code structure. You may reuse code and keep your projects manageable by using imports effectively to increase productivity. | <pre>1 from pyspark.sql import SparkSession</pre> |
| pip | To ensure that requests will function, the pip program searches for the package in the Python Package Index (PyPI), resolves any dependencies, and installs everything in your current Python environment. | <pre>1 pip list</pre> |
| pip install | The pip install <package> command looks for the latest version of the package and installs it. | <pre>1 pip install package_name</pre> |

| | | |
|------------------|---|--|
| pip3 install | <p>pip3 is the official package manager and pip command for Python 3. It enables installing and managing third-party software packages with features and functionality not found in the Python standard library. Pip3 installs packages from PyPI (Python Package Index), but it won't resolve dependencies or help you solve dependency conflicts.</p> | <pre>1 pip3 install package_name</pre> |
| print() | <p>Prints the specified message to the screen or other standard output device.</p> <p>The message can be a string or any other object; the object will be converted into a string before being written to the screen.</p> | <pre>1 print("Hello, World!")</pre> |
| python3 | <p>Python 3 is a widely used programming language known for its readability and versatility.</p> | <pre>1 sudo apt-get update 2 sudo apt-get install python3 3 python3 --version #to verify the python version</pre> |
| sc.setLogLevel() | <p>Using this method, you can change the log level to the desired level. Valid log levels include ALL, DEBUG, ERROR, FATAL, INFO, OFF, TRACE, and WARN.</p> | <p>Import necessary modules:</p> <pre>1 from pyspark import SparkContext</pre> <p>Create a SparkContext:</p> <pre>1 sc = SparkContext("local", "LogLevelExample")</pre> <p>Set the log level to a desired level (e.g., INFO, ERROR):</p> <pre>1 sc.setLogLevel("INFO")</pre> <p>Now, any logging messages with a severity level equal to or higher than INFO will be displayed</p> |
| setMaster() | <p>Denotes where to run your Spark application, local or cluster. When you run on a cluster, you need to specify the address of the Spark master or Driver URL for a distributed cluster. We usually assign a local[*] value to setMaster() in Spark while doing internal testing.</p> | <pre>1 <pre>from pyspark import SparkContext </pre> 2 <p>Create a SparkContext with a local master:</p> 3 <pre>sc = SparkContext("local[*]", "MyApp")</pre></td></pre> |
| show() | <p>Spark DataFrame show() is used to display the contents of the DataFrame in a table row and column format. By default, it shows only twenty rows, and the column values are truncated at twenty characters.</p> | <pre>1 df.show()</pre> |
| source | <p>Used to execute a script file in the current shell environment, allowing you to modify the current shell environment in the same way you would if you had typed commands manually.</p> | <p>Assuming a Bash shell:</p> <pre>1 source myscript.sh</pre> |
| virtualenv | <p>Primarily a command-line application. It modifies the environment variables in a shell to create an isolated Python environment, so you'll need to have a shell to run it. You can type in virtualenv (name of the application), followed by flags that control its behavior.</p> | <p>Creating a virtual environment named "myenv":</p> <pre>1 virtualenv myenv</pre> |