

Package/Method	Description	Code Example
agg()	Used to get the aggregate values like count, sum, avg, min, and max for each group.	<pre>1 agg_df = df.groupBy("column_name").agg({"column_to_aggregate": "sum"})</pre>
cache()	Apache Spark transformation that is often used on a DataFrame, data set, or RDD when you want to perform more than one action. cache() caches the specified DataFrame, data set, or RDD in the memory of your cluster's workers. Since cache() is a transformation, the caching operation takes place only when a Spark action (for example, count(), show(), take(), or write()) is also used on the same DataFrame, Dataset, or RDD in a single action.	<pre>1 df = spark.read.csv("customer.csv") 2 df.cache()</pre>
cd	Used to move efficiently from the existing working directory to different directories on your system.	<p>Basic syntax of the 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>
def	Used to define a function. It is placed before a function name that is provided by the user to create a user-defined function.	<pre>1 def greet(name):</pre> <p>This function takes a name as a parameter and prints a greeting.</p> <pre>1 print(f"Hello, {name}!")</pre> <p>Calling the function:</p> <pre>1 greet("John")</pre>
docker exec	Runs a new command in a running container. Only runs while the container's primary process is running, and it is not restarted if the container is restarted.	<pre>1 docker exec -it container_name command_to_run 2 docker exec -it my_container /bin/bash</pre>
docker rm	Used to remove one or more containers.	<p>To remove a single container by name or ID:</p> <pre>1 docker rm container_name_or_id</pre> <p>To remove multiple containers by specifying their names or IDs:</p> <pre>1 docker rm container1_name_or_id container2_name_or_id</pre> <p>To remove all stopped containers:</p> <pre>1 docker rm \$(docker ps -aq)</pre>
docker run	It runs a command in a new container, getting the image and starting the container if needed.	<pre>1 docker run [OPTIONS] IMAGE [COMMAND] [ARG...]</pre>

for	The <i>for</i> loop operates on lists of items. It repeats a set of commands for every item in a list.	<div><div>1fruits = ["apple", "banana", "cherry"]</div><div>Iterating through the list using a <i>for</i> loop for fruit in fruits:</div><div><div>1print(f"I like {fruit}s")</div></div></div>
groupby()	Used to collect the identical data into groups on DataFrame and perform count, sum, avg, min, max functions on the grouped data.	<div><div>1import pandas as pd</div><div>Sample DataFrame:</div><div><div>1data = {'Category': ['A', 'B', 'A', 'B', 'A', 'B'], 2Value': [10, 20, 15, 25, 30, 35]} 3df = pd.DataFrame(data)</div></div><div>Grouping by "Category" and performing aggregation operations:</div><div><div>1grouped = df.groupby('Category').agg({'Value': ['count', 'sum', 'mean', 'min', 'max']}) 2print(grouped)</div></div></div>
repartition()	Used to increase or decrease the RDD or DataFrame partitions by number of partitions or by a single column name or multiple column names.	<div><div>Create a sample DataFrame:</div><div><div>1data = [("John", 25), ("Peter", 30), ("Julie", 35), ("David", 40), ("Eva", 45)] 2columns = ["Name", "Age"] 3df = spark.createDataFrame(data, columns)</div></div><div>Show the current number of partitions.</div><div><div>1print("Number of partitions before repartitioning: ", df.rdd.getNumPartitions())</div></div><div>Repartition the DataFrame to 2 partitions.</div><div><div>1df_repartitioned = df.repartition(2)</div></div><div>Show the number of partitions after repartitioning.</div><div><div>1print("Number of partitions after repartitioning: ", df_repartitioned.rdd.getNumPartitions())</div></div><div>Stop the SparkSession.</div><div><div>1spark.stop()</div></div></div>

return	Used to end the execution of the function call and returns the result (value of the expression following the return keyword) to the caller.	<pre> 1 def add_numbers(a, b): 2 result = a + b 3 return result </pre> <p>Calling the function and capturing the returned value:</p> <pre> 1 sum_result = add_numbers(5, 6) </pre> <p>Printing the result.</p> <pre> 1 print("The sum is:", sum_result) </pre> <p>Output.</p> <pre> 1 The sum is: 11 </pre>
show()	Spark DataFrame show() is used to display the contents of the DataFrame in a table row and column format. By default, it shows only 20 rows, and the column values are truncated at 20 characters.	<pre> 1 df.show() </pre>

spark.read.csv("path")	Using this, you can read a CSV file with fields delimited by pipe, comma, tab (and many more) into a Spark DataFrame.	<pre> 1 from pyspark.sql import SparkSession </pre> <p>Create a SparkSession.</p> <pre> 1 spark = SparkSession.builder.appName("CSVReadExample").getOrCreate() </pre> <p>Read a CSV file into a Spark DataFrame.</p> <pre> 1 df = spark.read.csv("path_to_csv_file.csv", header=True, inferSchema=True) </pre> <p>Show the first few rows of the DataFrame.</p> <pre> 1 df.show() </pre> <p>Stop the SparkSession.</p> <pre> 1 spark.stop() </pre>
wget	Stands for web get. The wget is a free noninteractive file downloader command. Noninteractive means that it can work in the background when the user is not logged in.	<p>Basic syntax of the wget command; commonly used options are <code>[-V], [-h], [-b], [-e], [-o], [-a], [-q]</code></p> <pre> 1 wget [options]... [URL]... </pre> <p>Example 1: Specifies to download file.txt over HTTP website URL into the working directory.</p> <pre> 1 wget http://example.com/file.txt </pre> <p>Example 2: Specifies to download the archive.zip over HTTP website URL in the background and returns you to the command prompt in the interim.</p> <pre> 1 wget -b http://www.example.org/files/archive.zip </pre>

withColumn()

Transformation function of DataFrame which is used to change the value, convert the datatype of an existing column, create a new column, and many more.

Sample DataFrame:

```
1 data = [("John", 25), ("Peter", 30), ("David", 35)]
2 columns = ["Name", "Age"]
3 df = spark.createDataFrame(data, columns)
```

Using withColumn to create a new column and change values

```
1 updated_df = df \
2     .withColumn("DoubleAge", col("Age") * 2) # Create a new column "DoubleAge" by doubling the "Age" column
3 updated_df = updated_df \
4     .withColumn("AgeGroup", when(col("Age") <= 30, "Young")
5         .when((col("Age") > 30) & (col("Age") <= 40), "Middle-aged")
6         .otherwise("Old")) # Create a new column "AgeGroup" based on conditions
7 updated_df.show()
```

Stop the SparkSession.

```
1 spark.stop()
```