



Module 3 Cheat Sheet: Apache Spark

Package/Method	Description	Code Example
appName()	<p>A name for your job to display on the cluster web UI.</p> <p>An Apache Spark transformation often used on a DataFrame, data set, or RDD when you want to perform multiple actions. 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, data set, or RDD in a single action.</p>	<pre>1. 1 2. 2 1. from pyspark.sql import SparkSession 2. spark = SparkSession.builder.appName("MyApp").getOrCreate()</pre>
cache()	<p>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, data set, or RDD in a single action.</p>	<pre>1. 1 2. 2 1. df = spark.read.csv("customer.csv") 2. df.cache()</pre>
count()	<p>Returns the number of elements with</p>	<pre>1. 1 2. 2 1. count = df.count()</pre>

Copied!

Copied!

	the specified value.	2. print(count)	
			Copied!
createTempView()	Creates a temporary view that can later be used to query the data. The only required parameter is the name of the view.	1. 1 1. df.createOrReplaceTempView("cust_tbl")	Copied!
filter()	Returns an iterator where the items are filtered through a function to test if the item is accepted or not.	1. 1 1. filtered_df = df.filter(df['age'] > 30)	Copied!
getOrCreate()	Get or instantiate a SparkContext and register it as a singleton object.	1. 1 1. spark = SparkSession.builder.getOrCreate()	Copied!
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, which can increase your productivity.	1. 1 1. from pyspark.sql import SparkSession	Copied!
len()	Returns the number of items in an object. When the object is a string, the len() function	1. 1 2. 2 1. row_count = len(df.collect()) 2. print(row_count)	Copied!

	returns the number of characters in the string.	
	Returns a map object (an iterator) of the results after applying the given function to each item of a given iterable (list, tuple, etc.)	<pre>1. 1 2. 2</pre>
map()		<pre>1. rdd = df.rdd.map(lambda row: (row['name'], 2. row['age']))</pre> <div>Copied!</div>
	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. 1 1. pip list</pre> <div>Copied!</div>
pip		
	The pip install <package> command looks for the latest version of the package and installs it.	<pre>1. 1 1. pip install pyspark</pre> <div>Copied!</div>
pip install		
	Prints the specified message to the screen or other standard output device. The message can be a string or any other object; the object will be converted into a string before being written to the screen.	<pre>1. 1 1. print("Hello, PySpark!")</pre> <div>Copied!</div>
print()		
	Used to print or display the	<pre>1. 1 1. df.printSchema()</pre>
printSchema()		

schema of the DataFrame or data set in tree format along with the column name and data type. If you have a DataFrame or data set with a nested structure, it displays the schema in a nested tree format.

Copied!

sc.parallelize()

Creates a parallelized collection. Distributes a local Python collection to form an RDD. Using range is recommended if the input represents a range for performance.

1. 1

1. rdd = sc.parallelize([1, 2, 3, 4, 5])

Copied!

Used to select one or multiple columns, nested columns, column by index, all columns from the list, by regular expression from a DataFrame.

select()

select() is a transformation function in Spark and returns a new DataFrame with the selected columns.

1. 1

1. selected_df = df.select('name', 'age')

Copied!

show()

Spark DataFrame show() is used to display the

1. 1

1. df.show()

	<p>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> <p>Spark SQL can automatically infer the schema of a JSON data set and load it as a DataFrame. The read.json() function loads</p>	<div>Copied!</div>
spark.read.json	<p>data from a directory of JSON files where each line of the files is a JSON object. Note that the file offered as a JSON file is not a typical JSON file.</p> <p>To issue any SQL query, use the sql() method on the SparkSession instance . All spark.sql queries executed in this manner return a DataFrame on which you may perform further Spark operations if required.</p>	<pre>1. 1 1. json_df = spark.read.json("customer.json")</pre> <div>Copied!</div>
spark.sql()	<p>spark.sql queries executed in this manner return a DataFrame on which you may perform further Spark operations if required.</p>	<pre>1. 1 2. 2 1. result = spark.sql("SELECT name, age FROM cust_tbl WHERE age > 30") 2. result.show()</pre> <div>Copied!</div>
time()	<p>Returns the current time in the number of seconds since</p>	<pre>1. 1 2. 2 3. 3 1. from pyspark.sql.functions import current_timestamp 2. current_time = df.select(current_timestamp().alias("current_time"))</pre>

the Unix
Epoch.

3. current_time.show()

Copied!

Changelog

Date	Version	Changed by	Change Description
2023-09-06	1.0	Sameeksha Saxena	Initial version created

IBM Corporation 2023. All rights reserved.