

Module 4 Cheat Sheet: DataFrames and Spark SQL

Package/Method	Description	Code Example
appName()	A name for your job to display on the cluster web UI.	<pre>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("MyApp").getOrCreate()</pre>
createDataFrame()	Used to load the data into a Spark DataFrame.	<pre>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("MyApp").getOrCreate() data = [("Jhon", 30), ("Peter", 25), ("Bob", 35)] columns = ["name", "age"] Creating a DataFrame df = spark.createDataFrame(data, columns)</pre>
createTempView()	Create a temporary view that can later be used to query the data. The only required parameter is the name of the view.	<pre>df.createOrReplaceTempView("cust_tbl1")</pre>
fillna()	Used to replace NULL/None values on all or selected multiple DataFrame columns with either zero (0), empty string, space, or any constant literal values.	<pre>Replaced NULL/None values in a DataFrame filled_df = df.fillna(0) Replace with zero</pre>
filter()	Returns an iterator where the items are filtered through a function to test if the item is accepted or not.	<pre>filtered_df = df.filter(df['age'] > 30)</pre>
getOrCreate()	Get or instantiate a SparkContext and register it as a singleton object.	<pre>spark = SparkSession.builder.getOrCreate()</pre>
groupby()	Used to collect the identical data into groups on DataFrame and perform count, sum, avg, min, max functions on the grouped data.	<pre>Grouping data and performing aggregation grouped_df = df.groupBy("age").agg({"age": "count"})</pre>
head()	Returns the first n rows for the object based on position.	<pre>Returning the first 5 rows first_5_rows = df.head(5)</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, which can increase your productivity.	<pre>from pyspark.sql import SparkSession</pre>
pd.read_csv()	Required to access data from the CSV file from Pandas that retrieves data in the form of the data frame.	<pre>import pandas as pd Reading data from a CSV file into a DataFrame df_from_csv = pd.read_csv("data.csv")</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>pip list</pre>
pip install	The pip install <package> command looks for the latest version of the package and installs it.	<pre>pip install pyspark</pre>
printSchema()	Used to print or display the 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.	<pre>df.printSchema()</pre>
rename()	Used to change the row indexes and the column labels.	<pre>import pandas as pd Create a sample DataFrame data = {'A': [1, 2, 3], 'B': [4, 5, 6]} df = pd.DataFrame(data) Rename columns df = df.rename(columns={'A': 'X', 'B': 'Y'}) The columns 'A' and 'B' are now renamed to 'X' and 'Y'</pre>

Package/Method	Description	Code Example
		<pre>print(df)</pre>
<code>select()</code>	Used to select one or multiple columns, nested columns, column by index, all columns from the list, by regular expression from a DataFrame. <code>select()</code> is a transformation function in Spark and returns a new DataFrame with the selected columns.	<pre>selected_df = df.select('name', 'age')</pre>
<code>show()</code>	Spark DataFrame <code>show()</code> 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.	<pre>df.show()</pre>
<code>sort()</code>	Used to sort DataFrame by ascending or descending order based on single or multiple columns.	<p>Sorting DataFrame by a column in ascending order</p> <pre>sorted_df = df.sort("age")</pre> <p>Sorting DataFrame by multiple columns in descending order</p> <pre>sorted_df_desc = df.sort(["age", "name"], ascending=[False, True])</pre>
<code>SparkContext()</code>	It is an entry point to Spark and is defined in <code>org.apache.spark</code> package since version 1.x and used to programmatically create Spark RDD, accumulators, and broadcast variables on the cluster.	<pre>from pyspark import SparkContext</pre> <p>Creating a SparkContext</p> <pre>sc = SparkContext("local", "MyApp")</pre>
<code>SparkSession</code>	It is an entry point to Spark, and creating a <code>SparkSession</code> instance would be the first statement you would write to the program with RDD, DataFrame, and dataset	<pre>from pyspark.sql import SparkSession</pre> <p>Creating a SparkSession</p> <pre>spark = SparkSession.builder.appName("MyApp").getOrCreate()</pre>
<code>spark.read.json()</code>	Spark SQL can automatically infer the schema of a JSON data set and load it as a DataFrame. The <code>read.json()</code> function loads 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.	<pre>json_df = spark.read.json("customer.json")</pre>
<code>spark.sql()</code>	To issue any SQL query, use the <code>sql()</code> method on the <code>SparkSession</code> instance. All <code>spark.sql</code> queries executed in this manner return a DataFrame on which you may perform further Spark operations if required.	<pre>result = spark.sql("SELECT name, age FROM cust_tbl WHERE age > 30") result.show()</pre>
<code>spark.udf.register()</code>	In PySpark DataFrame, it is used to register a user-defined function (UDF) with Spark, making it accessible for use in Spark SQL queries. This allows you to apply custom logic or operations to DataFrame columns using SQL expressions.	<p>Registering a UDF (User-defined Function)</p> <pre>from pyspark.sql.functions import udf from pyspark.sql.types import StringType def my_udf(value): return value.upper() spark.udf.register("my_udf", my_udf, StringType())</pre>
<code>where()</code>	Used to filter the rows from DataFrame based on the given condition. Both <code>filter()</code> and <code>where()</code> functions are used for the same purpose.	<p>Filtering rows based on a condition</p> <pre>filtered_df = df.where(df['age'] > 30)</pre>
<code>withColumn()</code>	Transformation function of DataFrame used to change the value, convert the data type of an existing column, create a new column, and many more.	<p>Adding a new column and performing transformations</p> <pre>from pyspark.sql.functions import col new_df = df.withColumn("age_squared", col("age") ** 2)</pre>
<code>withColumnRenamed()</code>	Returns a new DataFrame by renaming an existing column.	<p>Renaming an existing column</p> <pre>renamed_df = df.withColumnRenamed("age", "years_old")</pre>



Skills Network