# Module 4 Cheat Sheet: DataFrames and Spark SQL

| Package/Method | Description | Code Example |
|---|---|---|
| | | 1. 1<br>2. 2 |
| appName() | A name for your job to display on the cluster web UI. | 1. `from pyspark.sql import SparkSession`<br>2. `spark = SparkSession.builder.appName("MyApp").getOrCreate()`<br><br>`Copied!`<br>1. 1<br>2. 2<br>3. 3<br>4. 4 |
| createDataFrame() | Used to load the data into a Spark DataFrame. | 1. `from pyspark.sql import SparkSession`<br>2. `spark = SparkSession.builder.appName("MyApp").getOrCreate()`<br>3. `data = [("Jhon", 30), ("Peter", 25), ("Bob", 35)]`<br>4. `columns = ["name", "age"]`<br><br>`Copied!`<br><br>Creating a DataFrame<br><br>1. 1<br><br>1. `df = spark.createDataFrame(data, columns)`<br><br>`Copied!` |
| createTempView() | Create a temporary view that can later be used to query the data. The only required parameter is the name of the view. | 1. 1<br><br>1. `df.createOrReplaceTempView("cust_tbl")`<br><br>`Copied!`<br><br>Replace NULL/None values in a DataFrame |
| 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. | 1. 1<br><br>1. `filled_df = df.fillna(0)`<br><br>`Copied!`<br><br>Replace with zero<br>1. 1 |
| filter() | Returns an iterator where the items are filtered through a function to test if the item is accepted or not. | 1. `filtered_df = df.filter(df['age'] > 30)`<br><br>`Copied!`<br>1. 1 |
| getOrCreate() | Get or instantiate a SparkContext and register it as a singleton object. | 1. `spark = SparkSession.builder.getOrCreate()`<br><br>`Copied!`<br><br>Grouping data and performing aggregation<br>1. 1 |
| groupby() | Used to collect the identical data into groups on DataFrame and perform count, sum, avg, min, max functions on the grouped data. | 1. `grouped_df = df.groupBy("age").agg({"age": "count"})`<br><br>`Copied!`<br><br>Returning the first 5 rows |
| head() | Returns the first *n* rows for the object based on position. | 1. 1<br><br>1. `first_5_rows = df.head(5)`<br><br>`Copied!`<br>1. 1 |
| 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. `from pyspark.sql import SparkSession`<br><br>`Copied!`<br><br>1. 1<br><br>1. `import pandas as pd`<br><br>`Copied!` |
| pd.read_csv() | Required to access data from the CSV file from Pandas that retrieves data in the form of the data frame. | Reading data from a CSV file into a DataFrame<br><br>1. 1<br><br>1. `df_from_csv = pd.read_csv("data.csv")`<br><br>`Copied!` |
| 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. | 1. 1<br><br>1. `pip list`<br><br>`Copied!` |
| pip install | The pip install <package> command looks for the latest version of the package and | 1. 1<br><br>1. `pip install pyspark` |

| Package/Method | Description | Code Example |
|---|---|---|
| | installs it. | `Copied!` |
| 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. | 1. 1<br><br>1. `df.printSchema()`<br>`Copied!`<br><br>1. 1<br><br>1. `import pandas as pd`<br>`Copied!`<br><br>Create a sample DataFrame<br><br>1. 1<br>2. 2<br><br>1. `data = {'A': [1, 2, 3], 'B': [4, 5, 6]}`<br>2. `df = pd.DataFrame(data)`<br>`Copied!` |
| rename() | Used to change the row indexes and the column labels. | Rename columns<br><br>1. 1<br><br>1. `df = df.rename(columns={'A': 'X', 'B': 'Y'})`<br>`Copied!`<br><br>The columns 'A' and 'B' are now renamed to 'X' and 'Y'<br><br>1. 1<br><br>1. `print(df)`<br>`Copied!` |
| select() | Used to select one or multiple columns, nested columns, column by index, all columns from the list, by regular expression from a DataFrame. select() is a transformation function in Spark and returns a new DataFrame with the selected columns. | 1. 1<br><br>1. `selected_df = df.select('name', 'age')`<br>`Copied!` |
| 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 twenty rows, and the column values are truncated at twenty characters. | 1. 1<br><br>1. `df.show()`<br>`Copied!` |
| sort() | Used to sort DataFrame by ascending or descending order based on single or multiple columns. | Sorting DataFrame by a column in ascending order<br><br>1. 1<br><br>1. `sorted_df = df.sort("age")`<br>`Copied!`<br><br>Sorting DataFrame by multiple columns in descending order<br><br>1. 1<br><br>1. `sorted_df_desc = df.sort(["age", "name"], ascending=[False, True])`<br>`Copied!` |
| SparkContext() | It is an entry point to Spark and is defined in org.apache.spark package since version 1.x and used to programmatically create Spark RDD, accumulators, and broadcast variables on the cluster. | 1. 1<br><br>1. `from pyspark import SparkContext`<br>`Copied!`<br><br>Creating a SparkContext<br><br>1. 1<br><br>1. `sc = SparkContext("local", "MyApp")`<br>`Copied!` |
| SparkSession | It is an entry point to Spark, and creating a SparkSession instance would be the first statement you would write to the program with RDD, DataFrame, and dataset | 1. 1<br><br>1. `from pyspark.sql import SparkSession`<br>`Copied!`<br><br>Creating a SparkSession<br><br>1. 1<br><br>1. `spark = SparkSession.builder.appName("MyApp").getOrCreate()`<br>`Copied!` |

| Package/Method | Description | Code Example |
|---|---|---|
| spark.read.json() | Spark SQL can automatically infer the schema of a JSON data set and load it as a DataFrame. The read.json() 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. | ```1. 1```<br><br>```1. json_df = spark.read.json("customer.json")```<br><br>Copied! |
| spark.sql() | 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. | ```1. 1```<br>```2. 2```<br><br>```1. result = spark.sql("SELECT name, age FROM cust_tbl WHERE age > 30")```<br>```2. result.show()```<br><br>Copied! |
| spark.udf.register() | Registering a UDF (User-defined Function)<br><br>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. | ```1. 1```<br>```2. 2```<br>```3. 3```<br>```4. 4```<br>```5. 5```<br><br>```1. from pyspark.sql.functions import udf```<br>```2. from pyspark.sql.types import StringType```<br>```3. def my_udf(value):```<br>```4.   return value.upper()```<br>```5. spark.udf.register("my_udf", my_udf, StringType())```<br><br>Copied! |
| where() | Filtering rows based on a condition<br><br>Used to filter the rows from DataFrame based on the given condition. Both filter() and where() functions are used for the same purpose. | ```1. 1```<br><br>```1. filtered_df = df.where(df['age'] > 30)```<br><br>Copied! |
| withColumn() | Adding a new column and performing transformations<br><br>Transformation function of DataFrame used to change the value, convert the data type of an existing column, create a new column, and many more. | ```1. 1```<br>```2. 2```<br><br>```1. from pyspark.sql.functions import col```<br>```2. new_df = df.withColumn("age_squared", col("age") ** 2)```<br><br>Copied! |
| withColumnRenamed() | Renaming an existing column<br><br>Returns a new DataFrame by renaming an existing column. | ```1. 1```<br><br>```1. renamed_df = df.withColumnRenamed("age", "years_old")```<br><br>Copied! |


Skills Network