

Hoja de trucos del Módulo 4: DataFrames y Spark SQL

Paquete/Método	Descripción	Ejemplo de Código
appName()	Un nombre para tu trabajo que se mostrará en la interfaz web del clúster.	<pre>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("MyApp").getOrCreate()</pre>
createDataFrame()	Se utiliza para cargar los datos en un DataFrame de Spark.	<pre>from pyspark.sql import SparkSession spark = SparkSession.builder.appName("MyApp").getOrCreate() data = [("Jhon", 30), ("Peter", 25), ("Bob", 35)] columns = ["name", "age"]</pre> <p>Creando un DataFrame</p> <pre>df = spark.createDataFrame(data, columns)</pre>
createTempView()	Crea una vista temporal que se puede usar más tarde para consultar los datos. El único parámetro requerido es el nombre de la vista.	<pre>df.createOrReplaceTempView("cust_tbl1")</pre>
fillna()	Se utiliza para reemplazar valores NULL/None en todas o en múltiples columnas seleccionadas del DataFrame con cero (0), cadena vacía, espacio, o cualquier valor literal constante.	<p>Reemplazar valores NULL/None en un DataFrame</p> <pre>filled_df = df.fillna(0)</pre> <p>Reemplazar con cero</p>
filter()	Devuelve un iterador donde los elementos se filtran a través de una función para probar si el elemento es aceptado o no.	<pre>filtered_df = df.filter(df['age'] > 30)</pre>

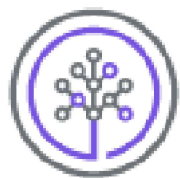
Paquete/Método	Descripción	Ejemplo de Código
getOrCreate()	Obtiene o instancia un SparkContext y lo registra como un objeto singleton.	<pre>spark = SparkSession.builder.getOrCreate()</pre>
groupby()	Se utiliza para agrupar los datos idénticos en grupos en el DataFrame y realizar funciones de conteo, suma, promedio, mínimo, máximo sobre los datos agrupados.	<p>Agrupando datos y realizando agregaciones</p> <pre>grouped_df = df.groupby("age").agg({"age": "count"})</pre>
head()	Devuelve las primeras <i>n</i> filas del objeto según la posición.	<p>Devolviendo las primeras 5 filas</p> <pre>first_5_rows = df.head(5)</pre>
import	Se utiliza para hacer que el código de un módulo sea accesible en otro. Las importaciones en Python son cruciales para una estructura de código exitosa. Puedes reutilizar código y mantener tus proyectos manejables utilizando importaciones de manera efectiva, lo que puede aumentar tu productividad.	<pre>from pyspark.sql import SparkSession</pre>
pd.read_csv()	Requerido para acceder a datos desde un archivo CSV de Pandas que recupera datos en forma de DataFrame.	<pre>import pandas as pd</pre> <p>Leyendo datos de un archivo CSV en un DataFrame</p> <pre>df_from_csv = pd.read_csv("data.csv")</pre>
pip	Para asegurar que las solicitudes funcionen, el programa pip busca el paquete en el Índice de Paquetes de Python (PyPI), resuelve cualquier dependencia e instala todo en tu entorno actual de Python.	<pre>pip list</pre>

Paquete/Método	Descripción	Ejemplo de Código
pip install	El comando pip install <package> busca la versión más reciente del paquete e instala.	<pre>pip install pyspark</pre>
printSchema()	Se utiliza para imprimir o mostrar el esquema del DataFrame o conjunto de datos en formato de árbol junto con el nombre de la columna y el tipo de dato. Si tienes un DataFrame o conjunto de datos con una estructura anidada, muestra el esquema en un formato de árbol anidado.	<pre>df.printSchema()</pre>
rename()	Se utiliza para cambiar los índices de fila y las etiquetas de columna.	<pre>import pandas as pd</pre> <p>Crear un DataFrame de ejemplo</p> <pre>data = {'A': [1, 2, 3], 'B': [4, 5, 6]} df = pd.DataFrame(data)</pre> <p>Renombrar columnas</p> <pre>df = df.rename(columns={'A': 'X', 'B': 'Y'})</pre> <p>Las columnas 'A' y 'B' ahora se han renombrado a 'X' y 'Y'</p> <pre>print(df)</pre>
select()	Se utiliza para seleccionar una o múltiples columnas, columnas anidadas, columna por índice, todas las columnas de la lista, por expresión regular de un DataFrame. select() es una función de transformación en Spark y	<pre>selected_df = df.select('name', 'age')</pre>

Paquete/Método	Descripción	Ejemplo de Código
	devuelve un nuevo DataFrame con las columnas seleccionadas.	
show()	El método show() de Spark DataFrame se utiliza para mostrar el contenido del DataFrame en formato de tabla con filas y columnas. Por defecto, muestra solo veinte filas, y los valores de las columnas se truncarán en veinte caracteres.	<code>df.show()</code>
sort()	Se utiliza para ordenar el DataFrame en orden ascendente o descendente según una o múltiples columnas.	<p>Ordenando el DataFrame por una columna en orden ascendente</p> <pre>sorted_df = df.sort("age")</pre> <p>Ordenando el DataFrame por múltiples columnas en orden descendente</p> <pre>sorted_df_desc = df.sort(["age", "name"], ascending=[False, True])</pre>
SparkContext()	Es un punto de entrada a Spark y se define en el paquete org.apache.spark desde la versión 1.x y se utiliza para crear programáticamente RDD de Spark, acumuladores y variables de difusión en el clúster.	<pre>from pyspark import SparkContext</pre> <p>Creamos un SparkContext</p> <pre>sc = SparkContext("local", "MyApp")</pre>
SparkSession	Es un punto de entrada a Spark, y crear una instancia de SparkSession sería la primera instrucción que escribirías en el programa con RDD, DataFrame y conjunto de datos.	<pre>from pyspark.sql import SparkSession</pre> <p>Creamos una SparkSession</p> <pre>spark = SparkSession.builder.appName("MyApp").getOrCreate()</pre>

Paquete/Método	Descripción	Ejemplo de Código
spark.read.json()	Spark SQL puede inferir automáticamente el esquema de un conjunto de datos JSON y cargarlo como un DataFrame. La función read.json() carga datos de un directorio de archivos JSON donde cada línea de los archivos es un objeto JSON. Ten en cuenta que el archivo ofrecido como archivo JSON no es un archivo JSON típico.	<pre>json_df = spark.read.json("customer.json")</pre>
spark.sql()	Para emitir cualquier consulta SQL, utiliza el método sql() en la instancia de SparkSession. Todas las consultas spark.sql ejecutadas de esta manera devuelven un DataFrame sobre el cual puedes realizar más operaciones de Spark si es necesario.	<pre>result = spark.sql("SELECT name, age FROM cust_tbl WHERE age > 30") result.show()</pre>
spark.udf.register()	En el DataFrame de PySpark, se utiliza para registrar una función definida por el usuario (UDF) con Spark, haciéndola accesible para su uso en consultas SQL de Spark. Esto te permite aplicar lógica o operaciones personalizadas a las columnas del DataFrame utilizando expresiones SQL.	<p>Registrando una UDF (Función definida por el usuario)</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>
where()	Se utiliza para filtrar las filas del DataFrame en función de la condición dada. Tanto las funciones filter() como where() se utilizan para el mismo propósito.	<p>Filtrando filas basadas en una condición</p> <pre>filtered_df = df.where(df['age'] > 30)</pre>
withColumn()	Función de transformación del DataFrame utilizada para cambiar el valor, convertir el tipo de dato de una columna existente, crear una nueva columna, y mucho más.	<p>Agregando una nueva columna y realizando transformaciones</p> <pre>from pyspark.sql.functions import col new_df = df.withColumn("age_squared", col("age") ** 2)</pre>
withColumnRenamed()	Devuelve un nuevo DataFrame renombrando una columna existente.	<p>Renombrando una columna existente</p> <pre>renamed_df = df.withColumnRenamed("age", "years_old")</pre>

Paquete/Método	Descripción	Ejemplo de Código



Skills Network