

⇒ SQL UDF ⇒

Cmd 5

Markdown - x

User-Defined Functions

User Defined Functions (UDFs) in Spark SQL allow you to register custom SQL logic as functions in a database, making these methods reusable anywhere SQL can be run on Databricks. These functions are registered natively in SQL and maintain all of the optimizations of Spark when applying custom logic to large datasets.

At minimum, creating a SQL UDF requires a function name, optional parameters, the type to be returned, and some custom logic.

Below, a simple function named `sale_announcement` takes an `item_name` and `item_price` as parameters. It returns a string that announces a sale for an item at 80% of its original price.

Cmd 6

SQL - x

```
1 CREATE OR REPLACE FUNCTION sale_announcement(item_name STRING, item_price INT)
2 RETURNS STRING
3 RETURN concat("The ", item_name, " is on sale for $", round(item_price * 0.8, 0));
4
5 SELECT *, sale_announcement(name, price) AS message FROM item_lookup
```

⇒ To see details of a UDF :-
[DESCRIBE FUNCTION EXTENDED <udf-name>]

⇒ UDF persists b/w notebooks, jobs (environments).

⇒ UDFs exists as objects in the metastore, and is governed by same table, view, ACLs.

Simple Control Flow Functions

Combining SQL UDFs with control flow in the form of `CASE` / `WHEN` clauses provides optimized execution for control flows within SQL workloads. The standard SQL syntactic construct `CASE` / `WHEN` allows the evaluation of multiple conditional statements with alternative outcomes based on table contents.

Here, we demonstrate wrapping this control flow logic in a function that will be reusable anywhere we can execute SQL.

Cmd 12

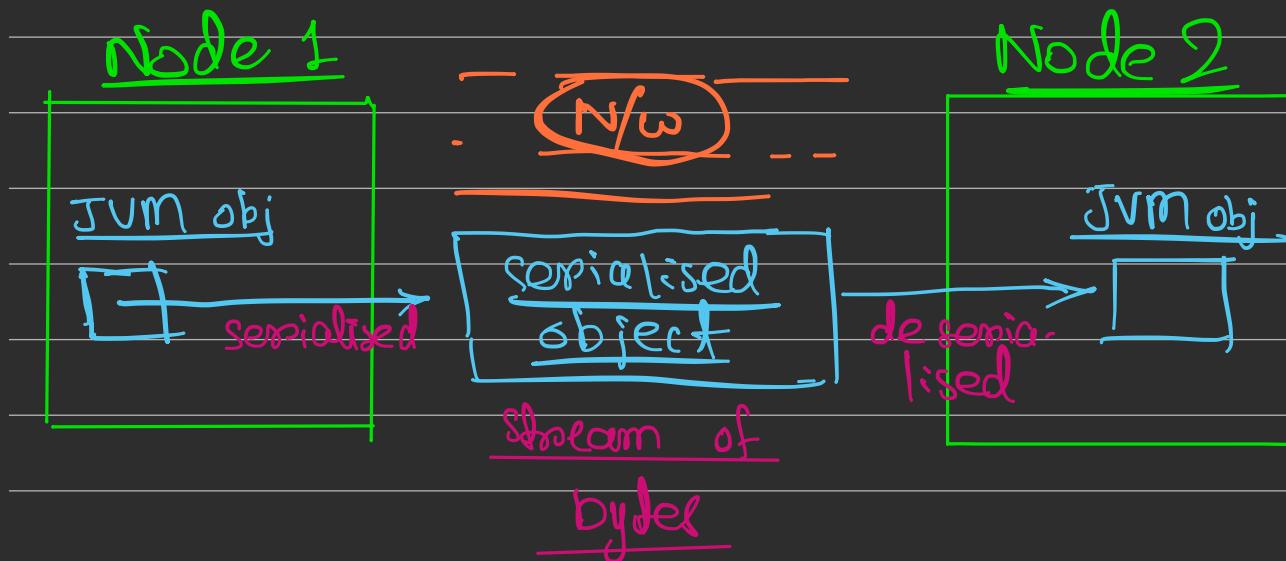
```
1 CREATE OR REPLACE FUNCTION item_preference(name STRING, price INT)
2 RETURNS STRING
3 RETURN CASE
4   WHEN name = "Standard Queen Mattress" THEN "This is my default mattress"
5   WHEN name = "Premium Queen Mattress" THEN "This is my favorite mattress"
6   WHEN price > 100 THEN concat("I'd wait until the ", name, " is on sale for $", round(price * 0.8, 0))
7   ELSE concat("I don't need a ", name)
8 END;
9
10 SELECT *, item_preference(name, price) FROM item_lookup
```

⇒ Serialisation :-

- ⇒ RDDs, Dataframes, Datasets, are all called JVM objects. These JVM objects exists as partitions in the nodes.
- ⇒ In a Cluster, there are multiple nodes, and each node has their own partitions. When wide Transformations are applied, there is shuffling of data, i.e, the data in the form of JVM objects are transferred from one node to other, through network.

Now, JVM objects cannot be transferred through n/w as it is.

To transfer JVM objects b/w nodes, they need to be Serialized first.



Serialization :-

Process of conversion of a JVM object into a Stream of Bytes.

Deserialization :-

Construction of a JVM object from the received Stream of Bytes.

⇒ Python UDF :-

⇒ We created a normal function :-

```
[def first_letter(var):  
    return var[0]]
```

⇒ This function can be registered as a PySpark UDF, using udf() function as :-

```
[first_letter_udf = udf(.first_letter)  
↓                                 ↓  
registered UDF             python function name]
```

This serialises the function and sends it to executors to be able to transform DF records.

Q Applying UDF on columns :-

≥ Registering UDF to use in SQL :-

Syntax:

python_udf = spark.udf.register("sql_udf", first_letter)

%sql
SELECT sql.udf(username) FROM users;

Column on which it is to be applied.

Use Decorator Syntax (Python Only)

Markdown

Alternatively, you can define and register a UDF using Python decorator syntax. The `@udf` decorator parameter is the Column datatype the function returns.

You will no longer be able to call the local Python function (i.e., `first_letter_udf("annagray@kaufman.com")` will not work).

Note: This example also uses [Python type hints](#), which were introduced in Python 3.5. Type hints are not required for this example, but instead serve as "documentation" to help developers use the function correctly. They are used in this example to emphasize that the UDF processes one record at a time, taking a single `str` argument and returning a `str` value.

Cmd 18

```
1 # Our input/output is a string
2 @udf("string")
3 def first_letter_udf(email: str) -> str:
4     return email[0]
```

Python type hints

Cmd 19

And let's use our decorator UDF here.

Cmd 20

```
1 from pyspark.sql.functions import col
2
3 sales_df = spark.table("sales")
4 display(sales_df.select(first_letter_udf(col("email"))))
```

Important Points :-

≥ A UDF

⇒ Cannot be optimised by Catalyst optimiser.

⇒ Is serialised and sent to executors.

⇒ Row data is first deserialised, and then UDF process deserialised data (objects). Then again, the result is serialized to Spark's native format.

⇒ For Python UDFs, additional interprocess communication overhead occurs b/w Python Interpreter & Executor