# Spark by {Examples} (https://sparkbyexamples.com/)

Spark (https://sparkbyexamples.com/)

PySpark   (https://sparkbyexamples.com/pyspark-tutorial/)

Hive   (https://sparkbyexamples.com/apache-hive-tutorial/)

Delicious Cakes
Ferns N Petals

HBase   (https://sparkbyexamples.com/apache-hbase-tutorial/)

# PySpark UDF (User Defined Function)

Kafka   (https://sparkbyexamples.com/apache-kafka-tutorials-with-examples/)

NNK (https://sparkbyexamples.com/author/admin/)  ·  PySpark (https://sparkbyexamples.com/category/pyspark/)

FAQ's   (https://sparkbyexamples.com/spark-questions/)

More ⌄   (https://sparkbyexamples.com/)

Pyspark UDF Example

PySpark UDF (a.k.a User Defined Function) is the most useful feature of Spark SQL & DataFrame that is used to extend the PySpark build in capabilities. In this article, I will explain what is UDF? why do we need it and how to create and use it on DataFrame `select()`, withColumn() (https://sparkbyexamples.com/pyspark/pyspark-dataframe-withcolumn/) and SQL using PySpark (Spark with Python) examples.

**Note:** UDF's are the most expensive operations hence use them only you have no choice and when essential. In the later section of the article, I will explain why using UDF's is an expensive operation in detail.

**Table of contents**

- PySpark UDF Introduction
  - What is UDF?
  - Why do we need it?
- Create PySpark UDF (User Defined Function)
  - Create a DataFrame
  - Create a Python function
  - Convert python function to UDF
- Using UDF with DataFrame
  - Using UDF with DataFrame select()
  - Using UDF with DataFrame withColumn()
  - Registring UDF & Using it on SQL query
- Create UDF using annotation
- Special handling
  - Null check
  - Performance concern
- Complete Example

# 1. PySpark UDF Introduction

## 1.1 What is UDF?

UDF's a.k.a User Defined Functions, If you are coming from SQL background, UDF's are nothing new to you as most of the traditional RDBMS databases support User Defined Functions, these

functions need to register in the database library and use them on SQL as regular functions.

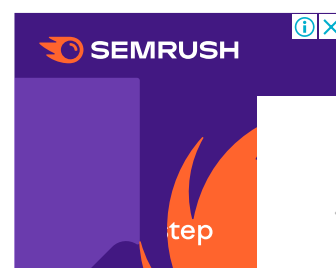PySpark UDF's are similar to UDF on traditional databases. In PySpark, you create a function in a Python syntax and wrap it with PySpark SQL udf() or register it as udf and use it on DataFrame and SQL respectively.

## 1.2 Why do we need a UDF?

UDF's are used to extend the functions of the framework and re-use these functions on multiple DataFrame's. For example, you wanted to convert every first letter of a word in a name string to a capital case; PySpark build-in features don't have this function hence you can create it a UDF and reuse this as needed on many Data Frames. UDF's are once created they can be re-used on several DataFrame's and SQL expressions.

Before you create any UDF, do your research to check if the similar function you wanted is already available in Spark SQL Functions (https://sparkbyexamples.com/spark/spark-sql-functions-understanding/). PySpark SQL provides several predefined common functions and many more new functions are added with every release. hence, It is best to check before you reinventing the wheel.

When you creating UDF's you need to design them very carefully otherwise you will come across optimization &

performance issues.

## 2. Create PySpark UDF

## 2.1 Create a DataFrame

Before we jump in creating a UDF, first
let's create a PySpark DataFrame
(https://sparkbyexamples.com/pyspark/
different-ways-to-create-dataframe-in-
pyspark/).

```
spark = SparkSession.builder.ap

columns = ["Seqno","Name"]
data = [("1", "john jones"),
    ("2", "tracey smith"),
    ("3", "amy sanders")]

df = spark.createDataFrame(data

df.show(truncate=False)
```

Yields below output.



```
+-----+------------+
|Seqno|Names       |
+-----+------------+
|1    |john jones  |
|2    |tracey smith|
|3    |amy sanders |
+-----+------------+
```

## 2.2 Create a Python

## Function

The first step in creating a UDF is creating a Python function. Below snippet creates a function `convertCase()` which takes a string parameter and converts the first letter of every word to capital letter. UDF's take parameters of your choice and returns a value.

```
ef convertCase(str):
    resStr=""
    arr = str.split(" ")
    for x in arr:
        resStr= resStr + x[0:1].
    return resStr
```

## 2.3 Convert a Python Function to PySpark UDF

Now convert this function `convertCase()` to UDF by passing the function to PySpark SQL udf(), this function is available at `.apache.spark.sql.functions.` package. Make sure you import package before using it.

Spark SQL udf() function returns `.apache.spark.sql.expression`... `serDefinedFunction` class object.

```
'" Converting function to UDF
onvertUDF = udf(lambda z: conv
```

Note: The default type of the udf() is StringType hence, you can also write the above statement without return type.

```
""" Converting function to UDF
StringType() is by default hence
convertUDF = udf(lambda z: conv
```

# 3. Using UDF with DataFrame

## 3.1 Using UDF with PySpark DataFrame select()

Now you can use `convertUDF()` on a DataFrame column as a regular build-in function.

```
df.select(col("Seqno"), \
    convertUDF(col("Name")).ali
    .show(truncate=False)
```

This results below output.

```
+-----+-------------+
|Seqno|Name         |
+-----+-------------+
|1    |John Jones   |
|2    |Tracey Smith |
|3    |Amy Sanders  |
+-----+-------------+
```

## 3.2 Using UDF with PySpark DataFrame withColumn()

You could also use udf on DataFrame `withColumn()` function, to explain this I will create another `upperCase()` function which converts the input string to upper case.

```
def upperCase(str):
    return str.upper()
```

Let's convert `upperCase()` python function to UDF and then use it with DataFrame `withColumn()`. Below

example converts the values of "Name" column to upper case and creates a new column "Curated Name"

```
upperCaseUDF = udf(lambda z:upp

df.withColumn("Cureated Name", 
    .show(truncate=False)
```

This yields below output.

```
+-----+-----------+----------
|Seqno|Name       |Cureated Na
+-----+-----------+----------
|1    |john jones |JOHN JONES
|2    |tracey smith|TRACEY SMIT
|3    |amy sanders|AMY SANDERS
+-----+-----------+----------
```

## 3.3 Registering PySpark UDF & use it on SQL

In order to use convertCase() function on PySpark SQL, you need to register the function with PySpark by using spark.udf.register().

```
""" Using UDF on SQL """
spark.udf.register("convertUDF"
df.createOrReplaceTempView("NAM
spark.sql("select Seqno, conver
    .show(truncate=False)
```

This yields the same output as 3.1 example.

## 4. Creating UDF using annotation

In the previous sections, you have learned creating a UDF is a 2 step process, first, you need to create a Python function, second convert function to UDF using SQL udf()

function, however, you can avoid these
two steps and create it with just a
single step by using annotations.

```python
@udf(returnType=StringType())
def upperCase(str):
    return str.upper()

df.withColumn("Cureated Name",
.show(truncate=False)
```

This results same output as section 3.2

# 5. Special Handling

## 5.1 Execution order

One thing to aware is in PySpark/Spark
does not guarantee the order of
evaluation of subexpressions meaning
expressions are not guarantee to
evaluated left-to-right or in any other
fixed order. PySpark reorders the
execution for query optimization and
planning hence, AND, OR, WHERE and
HAVING expression will have side
effects.

So when you are designing and using
UDF, you have to be very careful
especially with null handling as these
results runtime exceptions.

```python
"""
No guarantee Name is not null w
If convertUDF(Name) like '%John
you will get runtime error
"""
spark.sql("select Seqno, conver
        "where Name is not nul
    .show(truncate=False)
```

## 5.2 Handling null check

UDF's are error-prone when not
designed carefully. for example, when
you have a column that contains the
value null on some records

```
""" null check """

columns = ["Seqno","Name"]
data = [("1", "john jones"),
      ("2", "tracey smith"),
      ("3", "amy sanders"),
      ('4',None)]

df2 = spark.createDataFrame(dat
df2.show(truncate=False)
df2.createOrReplaceTempView("NAl

spark.sql("select convertUDF(Nal
      .show(truncate=False)
```

Note that from the above snippet, record with "Seqno 4" has value "None" for "name" column. Since we are not handling null with UDF function, using this on DataFrame returns below error. Note that in Python None is considered null.

```
AttributeError: 'NoneType' obje

    at org.apache.spark.api.pyth
    at org.apache.spark.sql.exec
    at org.apache.spark.sql.exec
    at org.apache.spark.api.pyth
    at org.apache.spark.Interrup
    at scala.collection.Iterator
```

Below points to remember

- Its always best practice to check for null inside a UDF function rather than checking for null outside.
- In any case, if you can't do a null check in UDF at lease use IF or CASE WHEN to check for null and call UDF conditionally.

```
spark.udf.register("_nullsafeUD

spark.sql("select _nullsafeUDF(
        .show(truncate=False)

spark.sql("select Seqno, _nulls
            " where Name is not n
        .show(truncate=False)
```

This executes successfully without errors as we are checking for null/none while registering UDF.

## 5.3 Performance concern using UDF

UDF's are a black box to PySpark hence it can't apply optimization and you will lose all the optimization PySpark does on Dataframe/Dataset. When possible you should use Spark SQL built-in functions (https://sparkbyexamples.com/spark/spark-sql-functions-understanding/) as these functions provide optimization. Consider creating UDF only when existing built-in SQL function doesn't have it.

## 6. Complete PySpark UDF Example

Below is complete UDF function example in Scala

```python
import pyspark
from pyspark.sql import SparkSe
from pyspark.sql.functions impo
from pyspark.sql.types import S

spark = SparkSession.builder.ap

columns = ["Seqno","Name"]
data = [("1", "john jones"),
    ("2", "tracey smith"),
    ("3", "amy sanders")]

df = spark.createDataFrame(data

df.show(truncate=False)

def convertCase(str):
    resStr=""
    arr = str.split(" ")
    for x in arr:
        resStr= resStr + x[0:1].
    return resStr

""" Converting function to UDF
convertUDF = udf(lambda z: conv

df.select(col("Seqno"), \
    convertUDF(col("Name")).ali
.show(truncate=False)

def upperCase(str):
    return str.upper()

upperCaseUDF = udf(lambda z:upp

df.withColumn("Cureated Name",
.show(truncate=False)

""" Using UDF on SQL """
spark.udf.register("convertUDF"
df.createOrReplaceTempView("NAM
spark.sql("select Seqno, conver
    .show(truncate=False)

spark.sql("select Seqno, conver
        "where Name is not nu
    .show(truncate=False)

""" null check """

columns = ["Seqno","Name"]
data = [("1", "john jones"),
    ("2", "tracey smith"),
    ("3", "amy sanders"),
    ('4',None)]
```

```
df2 = spark.createDataFrame(dat
df2.show(truncate=False)
df2.createOrReplaceTempView("NAI

spark.udf.register("_nullsafeUDF

spark.sql("select _nullsafeUDF(
      .show(truncate=False)

spark.sql("select Seqno, _nulls
          " where Name is not n
      .show(truncate=False)
```

This example is also available at [Spark GitHub project (https://github.com/spark-examples/pyspark-examples/blob/master/pyspark-udf.py)](https://github.com/spark-examples/pyspark-examples/blob/master/pyspark-udf.py) for reference.

## Conclusion

In this article, you have learned the following

- PySpark UDF is a User Defined Function that is used to create a reusable function in Spark.
- Once UDF created, that can be re-used on multiple DataFrames and SQL (after registering).
- The default type of the udf() is StringType.
- You need to handle nulls explicitly otherwise you will see side-effects.

## Reference

- [https://docs.databricks.com/spark/latest/spark-sql/udf-python.html (https://docs.databricks.com/spark/latest/spark-sql/udf-python.html)](https://docs.databricks.com/spark/latest/spark-sql/udf-python.html)
- [http://spark.apache.org/docs/latest/api/python/_modules/pyspark/sql/udf.html (https://spark.apache.org/docs/latest/api/python/_modules/pyspark/sql/udf.html)](http://spark.apache.org/docs/latest/api/python/_modules/pyspark/sql/udf.html)

**Share this:**

**TAGS:** **PYSPARK UDF
(HTTPS://SPARKBYEXAMPLES.COM/TAG/PYSPARK
-UDF/)**, **SPARK SQL UDF
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARK-
SQL-UDF/)**, **UDF
(HTTPS://SPARKBYEXAMPLES.COM/TAG/UDF/)**

## NNK

## (Https://Sparkbyexamples.Com/Author/Admin/)

(https://sp
arkbyexa
mples.co
m/author/
admin/)

SparkByExamples.com is a Big Data and Spark
examples community page, all examples are simple and
easy to understand and well tested in our development
environment Read more ..
(https://sparkbyexamples.com/about-sparkbyexamples/)

❯ **THIS POST HAS 4 COMMENTS**

**Anonymous**

30 JAN 2021      REPLY

Why are you showing the
whole example in Scala?
This is PySpark... it
should be in Python!

NNK  31 JAN 2021      REPLY

Hi, Complete example is in PySpark however, the Github link was pointing to Scala which I corrected now. Thanks for pointing it out.

## Anonymous

Do you know to make a UDF globally, means can a notebook calls the UDF defined in another notebook?

## Anonymous

If you use Zeppelin notebooks you can use the same interpreter in the several notebooks (change it in Intergpreter menu). Or search for precode option of Interpreter — in this optionn you can define any udf which will be created when the Interpreter started. You can setup the precode option in the same Interpreter menu

## Leave a Reply

## Categories

Apache Hadoop
(https://sparkbyexamples.com/category/hadoop/)

Apache Spark
(https://sparkbyexamples.com/category/spark/)

Apache Spark Streaming
(https://sparkbyexamples.com/category/spark/apache-spark-streaming/)

Apache Kafka
(https://sparkbyexamples.com/category/kafka/)

Apache HBase
(https://sparkbyexamples.com/category/hbase/)

Apache Cassandra
(https://sparkbyexamples.com/category/cassandra/)

Snowflake Database
(https://sparkbyexamples.com/catego

## Recent Posts

Spark regexp_replace() – Replace
String Value
(https://sparkbyexamples.com/spark/spark-regexp_replace-replace-string-value/)

How to Run a PySpark Script from
Python?
(https://sparkbyexamples.com/pyspark/run-pyspark-script-from-python-subprocess/)

Spark SQL like() Using Wildcard
Example
(https://sparkbyexamples.com/spark/spark-sql-like-using-wildcard-example/)

Spark isin() & IS NOT IN Operator
Example
(https://sparkbyexamples.com/spark/spark-isin-is-not-in-operator-example/)

Spark – Get Size/Length of Array & Map
Column
(https://sparkbyexamples.com/spark/spark-get-size-length-of-array-map-column/)

## About SparkByExamples.Com

SparkByExamples.com is a Big Data
and Spark examples community page,
all examples are simple and easy to
understand, and well tested in our
development environment Read more ..
(https://sparkbyexamples.com/about-sparkbyexamples/)

## Follow Us