# User-Defined Functions (UDF)

```
@author: Mahmoud Parsian
         Ph.D. in Computer Science
         email: mahmoud.parsian@yahoo.com
```

This short article shows how to use Python user-defined functions in PySpark applications. To use a UDF, we need to do some basic tasks:

1. create a UDF in Python
2. register UDF
3. use UDF in Spark SQL.

## 1. Define a UDF IN Python

Consider a function which triples its input:

```python
# n : integer
def tripled(n):
  return 3 * n
#end-def
```

## 2. Register UDF

To register a UDF, we can use `SparkSession.udf.register()`. The `register()` function takes 3 parameters:

- 1st: the desired name for UDF to be used in SQL
- 2nd: the name of Python UDF function
- 3rd: the return data type of Python UDF function (if this parameter is missing, then it is assumed that it is `StringType()`

```python
# "tripled_udf" : desired name to use in SQL
# tripled : defined Python function
# the last argument is the return type of UDF function
from pyspark.sql.types import IntegerType
spark.udf.register("tripled_udf", tripled, IntegerType())
```

Now, lets create a DataFrame and then apply the created UDF.

Create a sample DataFrame:

```
>>> data = [('alex', 20, 12000), ('jane', 30, 45000),
            ('rafa', 40, 56000), ('ted', 30, 145000),
            ('xo2', 10, 1332000), ('mary', 44, 555000)]
>>>
>>> column_names = ['name', 'age', 'salary']
>>> df = spark.createDataFrame(data, column_names)
>>>
>>> df
DataFrame[name: string, age: bigint, salary: bigint]
>>> df.printSchema()
root
 |-- name: string (nullable = true)
 |-- age: long (nullable = true)
 |-- salary: long (nullable = true)

>>>
>>> df.show()
+----+---+-------+
|name|age| salary|
+----+---+-------+
|alex| 20|  12000|
|jane| 30|  45000|
|rafa| 40|  56000|
| ted| 30| 145000|
| xo2| 10|1332000|
|mary| 44| 555000|
+----+---+-------+

>>> df.count()
6
>>> df2 = spark.sql("select * from people where salary > 67000")
>>> df2.show()
+----+---+-------+
|name|age| salary|
+----+---+-------+
| ted| 30| 145000|
| xo2| 10|1332000|
|mary| 44| 555000|
+----+---+-------+
```

## 3. Use UDF in SQL Query

```
>>> df.createOrReplaceTempView("people")
>>> df2 = spark.sql("select name, age, salary, tripled_udf(salary) as tripled_sala
>>> df2.show()
+----+---+-------+--------------+
|name|age| salary|tripled_salary|
+----+---+-------+--------------+
|alex| 20|  12000|         36000|
|jane| 30|  45000|        135000|
|rafa| 40|  56000|        168000|
| ted| 30| 145000|        435000|
| xo2| 10|1332000|       3996000|
|mary| 44| 555000|       1665000|
+----+---+-------+--------------+

>>>
```