# Struct Type & StructField, Array type, Array type Column, Map

## StructType() & StructField()

- PySpark StructType & StructField classes are used to programmatically specify the schema to the DataFrame and create complex columns like nested struct, array, and map columns
- StructType is a collection of StructField's

```python
from pyspark.sql.types import StructType,StructField, StringType, IntegerType

data = [(1,'Maheer',3000),(2,'Wafa',4000)]
schema = StructType([StructField(name='id',dataType=IntegerType()),\
                        StructField(name='name',dataType=StringType()),\
                        StructField(name='salary',dataType=IntegerType())\
                    ])
df = spark.createDataFrame(data,schema)
df.show()
df.printSchema()
```

```python
1   from pyspark.sql.types import StructType, StructField, StringType, IntegerType
2
3   data = [(1,('Maheer','Shaik'),3000),(2,('Wafa','Shaik'),4000)]
4
5   structName = StructType([\
6                           StructField('firstName',StringType()),\
7                           StructField('lastName',StringType())])
8
9   schema = StructType([\
10                      StructField(name='id',dataType=IntegerType()),\
11                      StructField(name='name',dataType=structName),\
12                      StructField(name='salary',dataType=IntegerType())])
13
14  df = spark.createDataFrame(data, schema)
15  df.show()
16
17  df.printSchema()
```

## ArrayType Column

- Create a dataframe with ArrayType column

```python
data = [('abc',[1,2]),('mno',[4,5]),('xyz',[7,8])]
df = spark.createDataFrame(data,['id','numbers'])
df.show()
df.printSchema()
```

```
| id|numbers|
+---+-------+
|abc| [1, 2]|
|mno| [4, 5]|
|xyz| [7, 8]|
```

```python
from pyspark.sql.types import StructType,StructField, IntegerType, StringType, ArrayType

data = [('abc',[1,2]),('mno',[4,5]),('xyz',[7,8])]
schema = StructType([\
                    StructField('id',StringType()),\
                    StructField('numbers',ArrayType(IntegerType()))\
                    ])
df = spark.createDataFrame(data,schema)
df.show()
df.printSchema()
```

```
+---+-------+
| id|numbers|
+---+-------+
|abc| [1, 2]|
|mno| [4, 5]|
|xyz| [7, 8]|
+---+-------+

root
 |-- id: string (nullable = true)
 |-- numbers: array (nullable = true)
 |    |-- element: integer (containsNull = true)
```

## ArrayType Column

- Fetch Value from Array as new column

```python
df.withColumn('firstNumber', col('numbers')[0]).show()
```

```
+---+-------+-----------+
| id|numbers|firstNumber|
+---+-------+-----------+
|abc| [1, 2]|          1|
|mno| [4, 5]|          4|
|xyz| [7, 8]|          7|
+---+-------+-----------+
```

- Combine columns to Array

```python
df = spark.createDataFrame(
    [(33, 44), (55, 66)], ["num1", "num2"]
)
df.show()

df.withColumn("nums", array(df.num1, df.num2)).show()
```

# MapType Column

- PySpark MapType is used to represent map key-value pair similar to python Dictionary (Dict) *Dictionary is nothing but a json.*

```python
data = [('maheer',{'hair':'black','eye':'brown'}),('wafa',{'hair':'black','eye':'blue'})]
schema = ['name','properties']
df = spark.createDataFrame(data,schema)
df.show()
display(df)
df.printSchema()
```

```python
from pyspark.sql.types import StructType, StructField, StringType, MapType

data = [('maheer',{'hair':'black','eye':'brown'}),('wafa',{'hair':'black','eye':'blue'})]
schema = StructType([\
                StructField('name',StringType()),\
                StructField('properties',MapType(StringType(),StringType()))\
                ])
df = spark.createDataFrame(data,schema)
df.show(truncate=False)
display(df)
df.printSchema()
```

▸ (3) Spark Jobs

▸ ▦ df: pyspark.sql.dataframe.DataFrame = [name: string, properties: map]

```
+------+--------------------+
|  name|          properties|
+------+--------------------+
|maheer|{eye -> brown, ha...|
|  wafa|{eye -> blue, hai...|
+------+--------------------+

root
 |-- name: string (nullable = true)
 |-- properties: map (nullable = true)
 |    |-- key: string
 |    |-- value: string (valueContainsNull = true)
```

# Column

- PySpark Column class represents a single Column in a DataFrame.
- **pyspark.sql.Column** class provides several functions to work with DataFrame to manipulate the Column values, evaluate the boolean expression to filter rows, retrieve a value or part of a value from a DataFrame column
- One of the simplest ways to create a Column class object is by using PySpark **lit()** SQL function

```python
from pyspark.sql.functions import lit
col1 = lit("abcd")
print(type(col1))
```

Cmd 2

```python
1  df1 = df.withColumn('newCol',lit('newColVal'))
2  df1.show()
3  df1.printSchema()
```

▸ (2) Spark Jobs

▸ ▦ df1: pyspark.sql.dataframe.DataFrame = [name: string, gen...

```
+------+------+------+---------+
|  name|gender|salary|   newCol|
+------+------+------+---------+
|maheer|  male|  2000|newColVal|
|  wafa|  male|  4000|newColVal|
+------+------+------+---------+
```

**colNotebook**  Python ▾

File  Edit  View  Run  Help      Last edit was 3 minutes ago    Give feedback      ▶ Run all    ● SHAIK

```python
1   from pyspark.sql.functions import lit
2
3   data = [('maheer','male',2000),('wafa','male',4000)]
4
5   schema = ['name','gender','salary']
6
7   df = spark.createDataFrame(data,schema)
8
9   df.show()
10  df.printSchema()
```

*Create a dataframe*

▸ (2) Spark Jobs

▸ ▦ df: pyspark.sql.dataframe.DataFrame = [name: string, gender: string ... 1 more field]

```
+------+------+------+
|  name|gender|salary|
+------+------+------+
|maheer|  male|  2000|
|  wafa|  male|  4000|
+------+------+------+
```