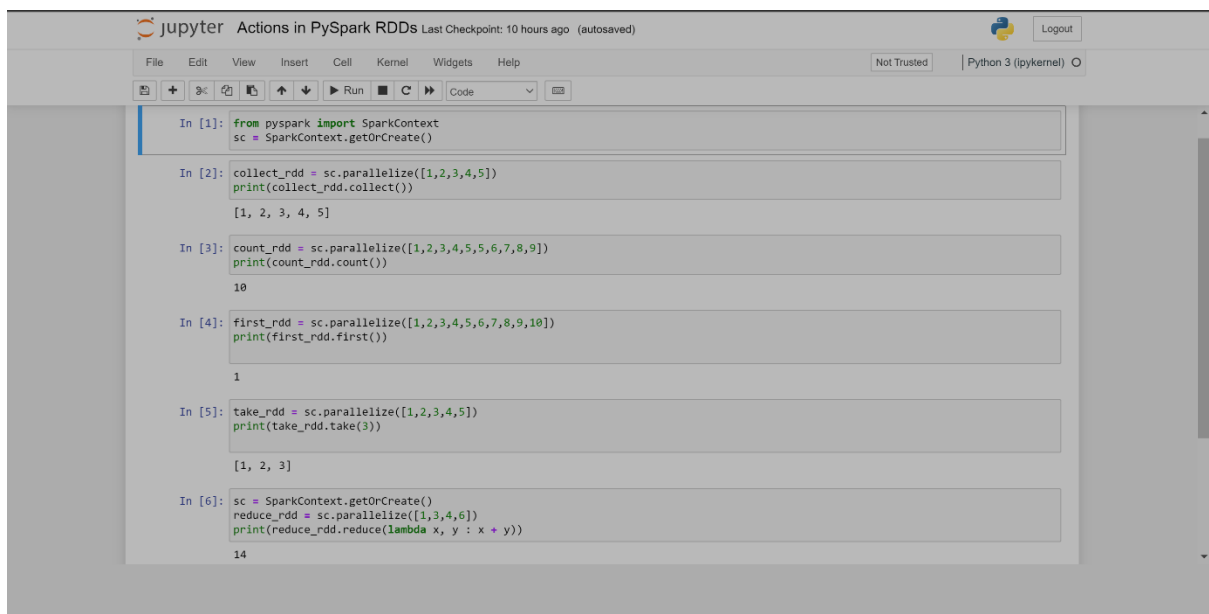**Name: Pradip Bochare**

# 🔧 PySpark RDD Operations

- Transformations: Transformations are a kind of operation that takes an RDD as input and produces another RDD as output. Once a transformation is applied to an RDD, it returns a new RDD, the original RDD remains the same and thus are immutable.

- Actions: Actions are a kind of operation which are applied on an RDD to produce a single value. These methods are applied on a resultant RDD and produces a non-RDD value

## ❖ Actions in PySpark RDDs

## ❖ Transformations in PySpark RDDs



## ✚ Selecting, Renaming, Filtering Data in a Pandas DataFrame

Method 1: Using withColumnRenamed()

We will use of withColumnRenamed() method to change the column names of pyspark data frame.

Syntax: DataFrame.withColumnRenamed(existing, new)

Parameters

- existingstr: Existing column name of data frame to rename.
- newstr: New column name.

Returns type: Returns a data frame by renaming an existing column.


Method 2: Using selectExpr()

Renaming the column names using selectExpr() method

Syntax : DataFrame.selectExpr(expr)

Parameters :

expr : It's an SQL expression.

Method 3: Using select() method
Syntax: DataFrame.select(cols)
Parameters :
cols: List of column names as strings.
Return type: Selects the cols in the dataframe and returns a new DataFrame.

Method 4: Using toDF()
This function returns a new DataFrame that with new specified column names.
Syntax: toDF(*col)
Where, col is a new column name

File   Edit   View   Insert   Cell   Kernel   Widgets   Help                                                   Trusted    | Python 3 (ipykernel) ○

```
|   Ram|1991-04-01|     M|  3000|
|  Mike|2000-05-19|     M|  4000|
|Rohini|1978-09-05|     M|  4000|
| Maria|1967-12-01|     F|  4000|
| Jenis|1980-02-17|     F|  1200|
+------+----------+------+------+
```

In [13]:
```python
Data_list = ["Emp Name","Date of Birth",
             " Gender-m/f","Paid salary"]

new_df = df.toDF(*Data_list)
new_df.show()
```

```
+--------+-------------+-----------+-----------+
|Emp Name|Date of Birth| Gender-m/f|Paid salary|
+--------+-------------+-----------+-----------+
|     Ram|   1991-04-01|          M|       3000|
|    Mike|   2000-05-19|          M|       4000|
|  Rohini|   1978-09-05|          M|       4000|
|   Maria|   1967-12-01|          F|       4000|
|   Jenis|   1980-02-17|          F|       1200|
+--------+-------------+-----------+-----------+
```

In [ ]:

**Notes**

\* PySpark RDD operations.

① Transformations

takes an RDD as input and produces another RDD as output. original RDD remains same and thus are immutable.

② Actions

Applied on an RDD to produce single value. produces a non-RDD value.

Transformations are applied to RDD to give another RDD While Actions are performed on RDD to give a non-RDD value.

⚹ Actions in PySpark RDD.

from pyspark import SparkContext.
sc = SparkContext.getOrCreate()

① • Collect() Action.
Returns list of all elements of RDD

collect_rdd = sc. parallelize ([1,2,3,4,5])
print (collect_rdd. collect())

② •count() Action.
returns no. of elements of RDD

count_rdd = sc. parallelize ([1,2,3,4,5,6,7,8,9])
print (count_rdd. count())

③ • first() Action

Returns first element of RDD.

```
first-rdd = sc.parallelize ([1,2,3,4,5,6,7,8,9,10])
print (first-rdd.first())
```

④ • take() Action

• take(n) returns n number of elements of RDD.

```
take-rdd = sc.parallelize ([1,2,3,4,5])
print (take-rdd.take())
```

⑤ • reduce() Action

takes two elements from RDD and operates.

```
reduce-rdd = sc.parallelize ([1,3,4,6]).
print (reduce-rdd.reduce(lambda x,y: x+y))
```

⑥ • saveAs Textfile() Action

```
save-rdd = sc.parallelize ([1,2,3,4,5,6])
sade-rdd.saveAs Textfile('file.txt')
```

\* Transformations in PySpark RDD

① .map()
maps value to elements of RDD

```
my-rdd = sc.parallelize([1,2,3,4])
print (my-rdd.map (lambda x : x+10). collect())
```

② .filter()
filtering elements from RDD

```
filter-rdd = sc. parallelize( [2,3,4,5,6,7])
print ( filter-rdd.filter(lambda x : x.1.2==0). collect())
```

③ .union()

```
union-inp= sc. parallelize ([2,4,5,6,7,8,9])
union-rdd-1 = union-inp.filter (lambda x : x.1.2==0 )
union-rdd.2 = union-inp. filter( lambda x : x.1.3==0)
print (union-rdd-1.union (union-rdd-2). collect())
```

④ ..flatMap()

```
flatmap-rdd = sc. parallelize (["Hey there", "this is PySpark"])
(flatmap-rdd.flatMap(lambda x : x.split(" ")).collect() )
```