

# Assignment – Day 13

-Sarthak Niranjan Kulkarni (Maverick)

- [sarthakkul2311@gmail.com](mailto:sarthakkul2311@gmail.com)

- (+91) 93256 02791

20/11/2024 (Wednesday)

## Transaction & Action Practice:-

### 1. Creating RDDs and DataFrame in PySpark with a Schema: -

#to create rdds and dataframe

```
from pyspark import SparkContext
from pyspark.sql import SparkSession
```

```
sc = SparkContext.getOrCreate()
spark = SparkSession.builder.appName('pyspark first program').getOrCreate()
```

#create the rdd

```
rdd = sc.parallelize([('C',85,76,87,91), ('B',85,76,87,91), ('A', 85,78,96,92), ('A',
92,76,89,96)], 4)
```

```
mydata = ['Division','English','Mathematics','Physics','Chemistry']
```

```
marks_df = spark.createDataFrame(rdd, schema=mydata)
```

```
print(rdd.collect())
```

```
print(rdd) #---Transformation which gives rdd value
```

```
rdd.collect() #----Action gives non rdd value
```

▶ (3) Spark Jobs

```
▶ marks_df: pyspark.sql.dataframe.DataFrame = [Division: string, English: long ... 3 more fields]
[('C', 85, 76, 87, 91), ('B', 85, 76, 87, 91), ('A', 85, 78, 96, 92), ('A', 92, 76, 89, 96)]
ParallelCollectionRDD[181] at readRDDFromInputStream at PythonRDD.scala:435
Out[14]: [('C', 85, 76, 87, 91),
 ('B', 85, 76, 87, 91),
 ('A', 85, 78, 96, 92),
 ('A', 92, 76, 89, 96)]
```

### 2. Creating RDDs, Converting to DataFrame, and Performing Actions in PySpark: -

#to create rdds and dataframe

```
from pyspark import SparkContext
```

```

from pyspark.sql import SparkSession

sc = SparkContext.getOrCreate()
spark = SparkSession.builder.appName('pyspark first program').getOrCreate()

#create the rdd

rdd = sc.parallelize([(('C',85,76,87,91), ('B',85,76,87,91), ("A", 85,78,96,92), ("A",
92,76,89,96)], 4)
mydata = ['Division','English','Mathematics','Physics','Chemistry']
marks_df = spark.createDataFrame(rdd, schema=mydata)
print(rdd.count())

rdd.take(2) ##Action gives non rdd value

```

#### ▶ (4) Spark Jobs

```

▶ marks_df: pyspark.sql.dataframe.DataFrame = [Division: string, English: long ... 3 more fields]
4
Out[13]: [('C', 85, 76, 87, 91), ('B', 85, 76, 87, 91)]

```

### 3. Counting the Number of Elements in an RDD Using PySpark: -

```

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

```

#### ▶ (2) Spark Jobs

```

10
Out[3]: 10

```

#### 4. Using count() and first() Actions to Analyze an RDD in PySpark: -

```
from pyspark import SparkContext

sc = SparkContext.getOrCreate()

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

print(first_rdd.count())

first_rdd.count()

first_rdd.first() #First method is action
```

▶ (3) Spark Jobs

```
10
Out[4]: 1
```

---

#### 5. Filtering Elements of an RDD Based on a Condition in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
filter_rdd_2 = sc.parallelize(['Rahul', 'Sarthak', 'Rohan',
'Shreya', 'Priya'])
print(filter_rdd_2.filter(lambda x: x.startswith('S')).collect())
```

▶ (1) Spark Jobs

```
['Sarthak', 'Shreya']
```

---

#### 6. Applying Filters and Performing Union Operation on RDDs in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
union_inp = sc.parallelize([2,4,5,6,7,8,9,10])
union_rdd_1 = union_inp.filter(lambda x:x % 2 == 0)
union_rdd_2 = union_inp.filter(lambda x:x % 3 == 0)
print(union_rdd_1.union(union_rdd_2).collect())
```

▶ (1) Spark Jobs

```
[2, 4, 6, 8, 10, 6, 9]
```

---

## 7. Using flatMap() Transformation to Split and Flatten Data in an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
flatmap_rdd = sc.parallelize(["Hey there", "This is PySpark RDD Transformations"])
print(flatmap_rdd.flatMap(lambda x :x.split(" ").collect()))
flatmap_rdd.flatMap(lambda x :x.split(" ").collect())
```

```
PythonRDD[217] at RDD at PythonRDD.scala:58
Out[17]: PythonRDD[218] at RDD at PythonRDD.scala:58
```

---

## 8. Using flatMap() to Split and Flatten Strings into Words in an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
flatmap_rdd = sc.parallelize(["Hey there", "This is PySpark RDD Transformations"])
(flatmap_rdd.flatMap(lambda x: x.split(" ")).collect())
```

► (1) Spark Jobs

```
Out[18]: ['Hey', 'there', 'This', 'is', 'PySpark', 'RDD', 'Transformations']
```

---

## 9. Using reduce() Action to Aggregate Elements in an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
reduce_rdd = sc.parallelize([1,2,3])
print(reduce_rdd.reduce(lambda x,y:x+y))
```

► (1) Spark Jobs

6

---

## 10. Saving an RDD to a Text File Using saveAsTextFile() in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
save_rdd = sc.parallelize([1,2,3,4,5,5])
save_rdd.saveAsTextFile('file3.txt')
```

► (1) Spark Jobs

---

## 11. Using map() Transformation to Modify Elements in an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
map_rdd = sc.parallelize([1,2,3])
print(map_rdd.map(lambda x:x+10))
print(map_rdd.map(lambda x:x+10).collect())
```

► (1) Spark Jobs

```
PythonRDD[147] at RDD at PythonRDD.scala:58
[11, 12, 13]
```

---

## 12. Using filter() Transformation to Select Even Numbers from an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
filter_rdd = sc.parallelize([1,2,3])
print(filter_rdd.filter(lambda x:x%2==0))
print(filter_rdd.filter(lambda x:x%2==0).collect())
```

► (1) Spark Jobs

```
PythonRDD[150] at RDD at PythonRDD.scala:58
[2]
```

---

## 13. Creating an RDD, Converting it to a DataFrame, and Registering it as a Temp View in PySpark

```
from pyspark import SparkContext
from pyspark.sql import SparkSession
sc = SparkContext.getOrCreate()
spark = SparkSession.builder.appName('pyspark first
program').getOrCreate()

#create the rdd
rdd = sc.parallelize([( 'C',85,76,87,91), ( 'B',85,76,87,91), ( "A",
85,78,96,92), ( "A", 92,76,89,96)], 4)
mydata = [ 'Division', 'English', 'Mathematics', 'Physics', 'Chemistry' ]
marks_df = spark.createDataFrame(rdd, schema=mydata)
print(marks_df.createOrReplaceTempView("dataofmarks"))
```

▶ (1) Spark Jobs

▶ marks\_df: pyspark.sql.dataframe.DataFrame = [Division: string, English: long ... 3 more fields]

None

## 14. Using take() Action to Retrieve the First N Elements from an RDD in PySpark

```
from pyspark import SparkContext
sc = SparkContext.getOrCreate()
count_rdd = sc.parallelize([1,2,3,4,5,5,6,7,8,9])
print(count_rdd.take(2))
count_rdd.take(6)
```

▶ (4) Spark Jobs

[1, 2]

Out[10]: [1, 2, 3, 4, 5, 5]

## Summary of Transaction & Action:-

### 1. Creating RDDs:

- RDDs are created using `sc.parallelize()` with sample data.

### 2. Creating DataFrames:

- RDDs are converted to DataFrames using `spark.createDataFrame()` with a predefined schema.

### 3. RDD Transformations:

- `map()`: Applies a function to each element of the RDD.
- `filter()`: Filters elements based on a condition.
- `flatMap()`: Flattens the results of applied functions (e.g., splitting strings into words).
- `union()`: Combines two RDDs into one.

#### **4. RDD Actions:**

- `count()`: Counts the number of elements in an RDD.
- `first()`: Returns the first element of the RDD.
- `reduce()`: Reduces the RDD to a single value (e.g., summing elements).
- `collect()`: Collects all elements from the RDD to the driver.
- `saveAsTextFile()`: Saves the RDD to a text file.

#### **5. DataFrame Operations:**

- A DataFrame can be registered as a temporary SQL view using `createOrReplaceTempView()`, enabling SQL queries on the data.

#### **6. File Operations:**

- The code includes file I/O operations like saving an RDD to a text file using `saveAsTextFile()`.

#### **7. Lazy and Eager Execution:**

- Transformations are lazy (not executed immediately), while actions are eager (trigger computation and return results).
-