

06/02/2023

Date \_\_\_\_\_  
Page \_\_\_\_\_

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

Action in PySpark RDD:-

① collect() Action:-

Returns list of all elements of the RDD.

collect-rdd = sc.parallelize([1, 2, 3, 4, 5, 7])  
print(collect-rdd.collect())

② count() Action:-

Count() action of RDD is an operation that returns the number of elements of our RDD.

count-rdd = sc.parallelize([1, 2, 3])  
print(count-rdd.count())

$$[\because \text{OP} = 3]$$

③ The first() Action:-

It returns the first element of our RDD

first-rdd = sc.parallelize([1, 2, 3])  
print(first-rdd.first())

$$[\because \text{OP} = 1]$$

④ take() Action:-

take(n) action on RDD returns n no. of elements from RDD.

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

$$[\because \text{OP} = [1, 2]]$$

⑤ reduce() Action:-

It takes two elements from the given RDD and operates

reduce-rdd = sc.parallelize([1, 3, 4, 6])

print(reduce-rdd.reduce(lambda x, y: x+y))

$$[\because \text{OP} = 16]$$

## ⑥ Save As text file :-

It is used to save the resultant RDD as text file.

Transformation in Pyspark:-

### ① The map() Transformation:-

It maps a value to the element of an RDD

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

`print (my-rdd.map(lambda x: x+10).collect())`  $\{:\text{OP}=[11,12,13,14]\}$

### ② The filter() Transformation:-

It filtering the elements from a Pyspark RDD.

`filter-rdd = sc.parallelize([1,2,3,4,5,6])`

`print (filter-rdd.filter(lambda x: x%2 == 0).collect())`  $\{:\text{OP}=[2,4,6]\}$

Renaming Column in pandas Dataframe:-

METHOD-1:- Column Renamed()

SYNTAX:- `Dataframe.with column Renamed(existing, new)`

CODE :- `df.withColumnRenamed("DOB", "Date of Birth").show()`

METHOD-2:- Using select Expr()

Syntax:- `Dataframe.selectExpr(expression)`

Code :- `data = df.selectExpr("Name as name", "DOB", "Gender", "Salary").show()`

METHOD-3:- Using Select() method.

Syntax:- `Dataframe.select(cols)`

Code :- `data.select(col("Name"), col("DOB"), col("Gender"), col("Salary")).alias("Amount")`

METHOD-4:- Using to DFC()

Code :- `DataList = ["Emp Name", "Date of Birth", "Gender", "Salary"]`

`new-df = df.toDF(*DataList)`

`new-df.show()`

community.cloud.databricks.com/?o=1969337358064883#notebook/3867848309277915/command/3867848309277916

Renaming columns in Dataframes Python

File Edit View Run Help Last edit was yesterday New cell UI: OFF

Run all Terminated Share Publish

```
2 # Importing necessary libraries
3 from pyspark.sql import SparkSession
4
5 # Create a spark session
6 spark = SparkSession.builder.appName('pyspark - example join').getOrCreate()
7
8 # Create data in dataframe
9 data = [(['Ram'], '1991-04-01', 'M', 3000),
10      ('Mike', '2000-05-19', 'M', 4000),
11      ('Rohini', '1978-09-05', 'M', 4000),
12      ('Maria', '1967-12-01', 'F', 4000),
13      ('Jenis', '1980-02-17', 'F', 1200)]
14
15 # Column names in dataframe
16 columns = ["Name", "DOB", "Gender", "salary"]
17
18 # Create the spark dataframe
19 df = spark.createDataFrame(data=data,
20                           schema=columns)
21
22 # Print the dataframe
23 df.show()
```

(3) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [Name: string, DOB: string ... 2 more fields]

Name	DOB	Gender	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

# Renaming columns in Dataframes

Python

New cell UI: OFF

Run all

Terminated

Share

Publish

Command took 2.65 seconds -- by saisachin1616@gmail.com at 2/6/2024, 4:59:02 PM on My Cluster

Cmd 2

```
1 #USING withColumnRenamed()
2
3 # Rename the column name from DOB to DateOfBirth
4 # Print the dataframe
5 df.withColumnRenamed("DOB","DateOfBirth").show()
```

Python

(3) Spark Jobs

Name	DateOfBirth	Gender	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

Command took 1.40 seconds -- by saisachin1616@gmail.com at 2/6/2024, 5:00:46 PM on My Cluster

Cmd 3

## Databricks

### Renaming columns in Dataframes

Python

New cell UI: OFF

Run all

Terminated

Share

Publish

Command took 1.40 seconds -- by saisachin1616@gmail.com at 2/6/2024, 5:00:46 PM on My Cluster

Cmd 3

```
1 # Renaming multiple column names
2
3 # Rename the column name 'Gender' to 'Sex'
4 # Then for the returning dataframe
5 # again rename the 'salary' to 'Amount'
6 df.withColumnRenamed("Gender","Sex").withColumnRenamed("salary","Amount").show()
```

(3) Spark Jobs

Name	DOB	Sex	Amount
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

Command took 0.85 seconds -- by saisachin1616@gmail.com at 2/6/2024, 5:01:55 PM on My Cluster

databricks

## Renaming columns in Dataframes

Python

New cell UI: OFF

Run all

Terminated

Share

Publish

```
1 # USING selectExpr()
2
3 # Select the 'Name' as 'name'
4 # Select remaining with their original name
5 data = df.selectExpr("Name as name", "DOB", "Gender", "salary")
6
7 # Print the dataframe
8 data.show()
```

(3) Spark Jobs

data: pyspark.sql.dataframe.DataFrame = [name: string, DOB: string ... 2 more fields]

name	DOB	Gender	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



## Renaming columns in Dataframes

Python

New cell UI: OFF

Run all

Terminated

Share

Publish

```
1 # USING select() method
2
3 # Import col method from pyspark.sql.functions
4 from pyspark.sql.functions import col
5
6 # Select the 'salary' as 'Amount' using aliasing
7 # Select remaining with their original name
8 data = df.select(col("Name"),col("DOB"),col("Gender"),col("salary").alias('Amount'))
9
10 # Print the dataframe
11 data.show()
```

(3) Spark Jobs

data: pyspark.sql.dataframe.DataFrame = [Name: string, DOB: string ... 2 more fields]

Name	DOB	Gender	Amount
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

## Databricks

### Renaming columns in Dataframes

Python

New cell UI: OFF

Run all

Terminated

Share

Publish

Cmd 6

```
1 # USING toDF()
2
3 Data_list = ["Emp Name", "Date of Birth", "Gender-m/f", "Paid salary"]
4 new_df = df.toDF(*Data_list)
5 new_df.show()
```

(3) Spark Jobs

new\_df: pyspark.sql.dataframe.DataFrame = [Emp Name: string, Date of Birth: string ... 2 more fields]

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

Command took 1.11 seconds -- by saisachin1616@gmail.com at 2/6/2024, 5:07:02 PM on My Cluster

Cmd 7



Transformations in PySpark RDDs Python ▾

File Edit View Run Help Last edit was 6 minutes ago New cell UI: OFF ▾

Run all Demo Share Publish

1 3. The .union() Transformation

Cmd 6

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

▶ (1) Spark Jobs

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

Command took 0.77 seconds -- by saisachin1616@gmail.com at 2/7/2024, 5:35:22 PM on Demo

Cmd 7

1 4. The .flatMap() Transformation

Cmd 8

```
1 flatmap_rdd = sc.parallelize(["Hey there", "This is PySpark RDD Transformations"])
2 (flatmap_rdd.flatMap(lambda x: x.split(" ")).collect())
```

▶ (1) Spark Jobs

['Hey', 'there', 'This', 'is', 'PySpark', 'RDD', 'Transformations']

Command took 1.15 seconds -- by saisachin1616@gmail.com at 2/7/2024, 5:37:49 PM on Demo

Cmd 9





Actions in PySpark RDDs Python ★

File Edit View Run Help Last edit was 13 minutes ago New cell UI: OFF

Run all Demo Share Publish

1 5. The `.reduce()` Action

Cmd 11

```
1 from pyspark import SparkContext
2 sc = SparkContext.getOrCreate()
3 reduce_rdd = sc.parallelize([1,3,4,6])
4 print(reduce_rdd.reduce(lambda x, y : x + y))
```

▶ (1) Spark Jobs

14

Command took 0.61 seconds -- by saisachin1616@gmail.com at 2/7/2024, 5:29:11 PM on Demo

Cmd 12

1 6. The `.saveAsTextFile()` Action

Cmd 13

```
1 from pyspark import SparkContext
2 sc = SparkContext.getOrCreate()
3 save_rdd = sc.parallelize([1, 2, 3, 4, 5, 6])
4 save_rdd.saveAsTextFile('dbfs:/output_directory/file.txt')
5 sc.stop()
```

▶ (1) Spark Jobs

The spark context has stopped and the driver is restarting. Your notebook will be automatically reattached.

Cmd 14