

Asmita Porwal  
Data Engineering  
Batch-1  
12/02/2024

## Coding Challenge -3 Question-1

### Execute Manipulating

```
In [1]: from pyspark.sql import SparkSession

In [2]: import pyspark
import findspark
findspark.init()

In [3]: from pyspark import SparkContext
sc = SparkContext("local", "RDD Transformation")
sc

Out[3]: SparkContext

Spark UI
Version
v3.5.0
Master
local
AppName
RDD Transformation
```

### GroupBy DataFrames & Aggregations

#### SUM

```
In [4]: spark = SparkSession.builder.appName("Practice").getOrCreate()

In [5]: df_pyspark = spark.read.csv("E:\\downloads\\Marks_data.csv", header=True, inferSchema=True)
df_pyspark.show()

+---+-----+-----+---+
|Name|M1_Score|M2_Score|age|
+---+-----+-----+---+
|Alex|62|80|20|
|Brad|45|56|19|
|Joey|85|98|21|
|Abhi|54|79|20|
+---+-----+-----+---+

In [6]: df_pyspark.groupBy("age").sum("M2_Score").show()

+---+-----+
|age|sum(M2_Score)|
+---+-----+
|20|159|
|19|56|
|21|98|
+---+-----+
```

## Min and Max

```
In [7]: df_pyspark.groupBy("age").min("M2 Score").show()
```

```
+---+-----+
|age|min(M2 Score)|
+---+-----+
| 20|          79|
| 19|          56|
| 21|          98|
+---+-----+
```

```
In [8]: df_pyspark.groupBy("age").max("M2 Score").show()
```

```
+---+-----+
|age|max(M2 Score)|
+---+-----+
| 20|          80|
| 19|          56|
| 21|          98|
+---+-----+
```

## Count and Sum two columns

```
In [11]: df_pyspark.groupBy("age").count().show()
```

```
+---+-----+
|age|count|
+---+-----+
| 20|     2|
| 19|     1|
| 21|     1|
+---+-----+
```

```
In [12]: df_pyspark.groupBy("Name", "age").sum("M2 Score").show()
```

```
+---+---+-----+
|Name|age|sum(M2 Score)|
+---+---+-----+
|Alex| 20|          80|
|Brad| 19|          56|
|Abhi| 20|          79|
|Joey| 21|          98|
+---+---+-----+
```

## Aggregations using agg sum

### And pivot element

```
In [13]: df_pyspark.groupBy("age").agg(({ "M2 Score": "sum"})).show()
```

```
+-----+
|age|sum(M2 Score)|
+-----+
| 20|          159|
| 19|           56|
| 21|           98|
+-----+
```

```
In [14]: df_pyspark.agg(({ "M1 Score": "sum"})).show()
```

```
+-----+
|sum(M1 Score)|
+-----+
|          246|
+-----+
```

```
In [15]: df_pyspark.groupBy("age").pivot("Name").sum("M2 Score").show()
```

```
+-----+-----+
|age|Abhi|Alex|Brad|Joey|
+-----+-----+
| 20|  79|  80|NULL|NULL|
| 19|NULL|NULL|  56|NULL|
| 21|NULL|NULL|NULL|  98|
+-----+-----+
```

## Dropping

```
In [51]: df_pyspark.na.drop(how="all").show()
```

```
+-----+-----+-----+
|Name|M1 Score|M2 Score|age|
+-----+-----+-----+
|Alex|      62|      80| 20|
|Brad|      45|      56| 19|
|Joey|      85|      98| 21|
|Abhi|      54|      79| 20|
+-----+-----+-----+
```

```
In [52]: df_pyspark.na.drop(how="any", thresh=2).show()
```

```
+-----+-----+-----+
|Name|M1 Score|M2 Score|age|
+-----+-----+-----+
|Alex|      62|      80| 20|
|Brad|      45|      56| 19|
|Joey|      85|      98| 21|
|Abhi|      54|      79| 20|
+-----+-----+-----+
```

```
In [53]: df_pyspark.na.drop(how="any", subset=["M2 Score"]).show()
```

```
+-----+-----+-----+
|Name|M1 Score|M2 Score|age|
+-----+-----+-----+
|Alex|      62|      80| 20|
|Brad|      45|      56| 19|
|Joey|      85|      98| 21|
|Abhi|      54|      79| 20|
+-----+-----+-----+
```

## Sorting

### Sort by age in ascending order and descending order

```
In [54]: df_pyspark.sort("age").show()
```

```
+---+-----+-----+---+
|Name|M1 Score|M2 Score|age|
+---+-----+-----+---+
|Brad|    45|    56| 19|
|Alex|    62|    80| 20|
|Abhi|    54|    79| 20|
|Joey|    85|    98| 21|
+---+-----+-----+---+
```

```
In [55]: df_pyspark.sort(df_pyspark["age"].desc()).show()
```

```
+---+-----+-----+---+
|Name|M1 Score|M2 Score|age|
+---+-----+-----+---+
|Joey|    85|    98| 21|
|Alex|    62|    80| 20|
|Abhi|    54|    79| 20|
|Brad|    45|    56| 19|
+---+-----+-----+---+
```

## Sorting 2 columns at a time

```
In [21]: df_pyspark.sort("age", "Name").show()
```

```
+---+-----+-----+---+
|Name|M1 Score|M2 Score|age|
+---+-----+-----+---+
|Brad|    45|    56| 19|
|Abhi|    54|    79| 20|
|Alex|    62|    80| 20|
|Joey|    85|    98| 21|
+---+-----+-----+---+
```

## Joining

```
In [23]: emp = [(1,"Smith",-1,"2018","10","M",3000),(2, "Rose",1, "2010", "20","M", 4000),(3,"Williams",1,"2010","10","M",1000),(4, "Jones",2,"2010","10","M",2000)]
empColumns = ["emp_id","name","superior_emp_id","year_joined", "emp_dept_id","gender","salary"]
```

```
In [24]: empDF = spark.createDataFrame(data=emp, schema = empColumns)
empDF.printSchema()
```

```
root
 |-- emp_id: long (nullable = true)
 |-- name: string (nullable = true)
 |-- superior_emp_id: long (nullable = true)
 |-- year_joined: string (nullable = true)
 |-- emp_dept_id: string (nullable = true)
 |-- gender: string (nullable = true)
 |-- salary: long (nullable = true)
```

In [25]: empDF.show()

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary |
|--------|----------|-----------------|-------------|-------------|--------|--------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     |
| 6      | Brown    | 2               | 2010        | 50          |        | -1     |

In [26]: dept = [("Finance",10),("Marketing",20),("Sales",30),("IT",40)]  
deptColumns = ["dept\_name","dept\_id"]  
deptDF = spark.createDataFrame(data=dept, schema = deptColumns)  
deptDF.printSchema()  
deptDF.show()

```
root
|-- dept_name: string (nullable = true)
|-- dept_id: long (nullable = true)
```

| dept_name | dept_id |
|-----------|---------|
| Finance   | 10      |
| Marketing | 20      |
| Sales     | 30      |
| IT        | 40      |

## Inner join and outer join

In [27]: empDF.join(deptDF,empDF.emp\_dept\_id == deptDF.dept\_id,"inner").show()

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |

In [28]: empDF.join(deptDF,empDF.emp\_dept\_id == deptDF.dept\_id,"outer").show()

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| NULL   | NULL     | NULL            | NULL        | NULL        | NULL   | NULL   | Sales     | 30      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |
| 6      | Brown    | 2               | 2010        | 50          |        | -1     | NULL      | NULL    |

## Full Join

```
In [29]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"full").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| NULL   | NULL     | NULL            | NULL        | NULL        | NULL   | NULL   | Sales     | 30      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |
| 6      | Brown    | 2               | 2010        | 50          |        | -1     | NULL      | NULL    |

```
In [30]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"fullouter").show()
```

## Full outer join

```
In [30]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"fullouter").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| NULL   | NULL     | NULL            | NULL        | NULL        | NULL   | NULL   | Sales     | 30      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |
| 6      | Brown    | 2               | 2010        | 50          |        | -1     | NULL      | NULL    |

## Left and left outer

```
In [31]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"left").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 6      | Brown    | 2               | 2010        | 50          |        | -1     | NULL      | NULL    |
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |

```
In [32]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"leftouter").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 6      | Brown    | 2               | 2010        | 50          |        | -1     | NULL      | NULL    |
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |

## Right and Right outer

```
In [33]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"right").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| NULL   | NULL     | NULL            | NULL        | NULL        | NULL   | NULL   | Sales     | 30      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |

```
In [34]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"rightouter").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary | dept_name | dept_id |
|--------|----------|-----------------|-------------|-------------|--------|--------|-----------|---------|
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   | Finance   | 10      |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   | Finance   | 10      |
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   | Finance   | 10      |
| NULL   | NULL     | NULL            | NULL        | NULL        | NULL   | NULL   | Sales     | 30      |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   | Marketing | 20      |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     | IT        | 40      |

## Left semi and left anti

```
In [35]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"leftsemi").show()
```

| emp_id | name     | superior_emp_id | year_joined | emp_dept_id | gender | salary |
|--------|----------|-----------------|-------------|-------------|--------|--------|
| 1      | Smith    | -1              | 2018        | 10          | M      | 3000   |
| 3      | Williams | 1               | 2010        | 10          | M      | 1000   |
| 4      | Jones    | 2               | 2005        | 10          | F      | 2000   |
| 2      | Rose     | 1               | 2010        | 20          | M      | 4000   |
| 5      | Brown    | 2               | 2010        | 40          |        | -1     |

```
In [36]: empDF.join(deptDF,empDF.emp_dept_id == deptDF.dept_id,"leftanti").show()
```

| emp_id | name  | superior_emp_id | year_joined | emp_dept_id | gender | salary |
|--------|-------|-----------------|-------------|-------------|--------|--------|
| 6      | Brown | 2               | 2010        | 50          |        | -1     |