**Pyspark3** · Python ⌄ ☆
File  Edit  View  Run  Help     Last edit was 5 minutes ago     New cell UI: OFF ⌄

▶ Run all   ● My Cluster ⌄   Share   Publish

Cmd 1
```
1
```

Cmd 2                                                                    Python ▶▼ ⌄ − ✕
```
1
```

Cmd 3                                                                    Python ▶▼ ⌄ − ✕
```
1    # Joins
2
```
Command took 1.47 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 4
```
1    from pyspark.sql import SparkSession
2    spark =SparkSession.builder.appName('pyspark - example join').getOrCreate()
3    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 ,"2005","10","F",2000),(5,"Brown",2,"2010",
     "40","",-1),(6, "Brown", 2, "2010","50","",-1)]
4    empColumns = ["emp_id","name","superior_emp_id","year_joined", "emp_dept_id","gender","salary"]
5
6    empDF = spark.createDataFrame(data=emp, schema = empColumns)
7    empDF.printSchema()
8    empDF.show()
9
10   dept = [("Finance",10),("Marketing",20),("Sales",30),("IT",40)]
11   deptColumns = ["dept_name","dept_id"]
12   deptDF = spark.createDataFrame(data=dept, schema = deptColumns)
13   deptDF.printSchema()
14   deptDF.show()
```

▶ (6) Spark Jobs

▶ ▦ empDF:  pyspark.sql.dataframe.DataFrame = [emp_id: long, name: string ... 5 more fields]
▶ ▦ deptDF:  pyspark.sql.dataframe.DataFrame = [dept_name: string, dept_id: long]

```
|    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|
+------+--------+---------------+-----------+----------+------+------+

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

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

Command took 15.36 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 5
```
1    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"inner") .show()
```

▶ (3) Spark Jobs

```
+------+--------+---------------+-----------+-----------+------+------+---------+-------+
|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|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+
```

Command took 6.38 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 6
```
1    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"outer").show()
2    #Or
3    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"full").show()
4    #Or
5    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"fullouter").show()
```

▶ (9) Spark Jobs

```
|    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|
+------+--------+---------------+-----------+----------+------+------+---------+------+
```

```
|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|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+
```

Command took 5.58 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 7

```
1    # Left Join
2    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"left").show()
3    #Or
4    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"leftouter").show()
```

▶ (12) Spark Jobs

```
|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|
|     2|    Rose|              1|       2010|         20|     M|  4000|Marketing|     20|
|     3|Williams|              1|       2010|         10|     M|  1000|  Finance|     10|
|     4|   Jones|              2|       2005|         10|     F|  2000|  Finance|     10|
|     5|   Brown|              2|       2010|         40|      |    -1|       IT|     40|
|     6|   Brown|              2|       2010|         50|      |    -1|     null|   null|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+


|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|
|     2|    Rose|              1|       2010|         20|     M|  4000|Marketing|     20|
|     3|Williams|              1|       2010|         10|     M|  1000|  Finance|     10|
|     4|   Jones|              2|       2005|         10|     F|  2000|  Finance|     10|
|     5|   Brown|              2|       2010|         40|      |    -1|       IT|     40|
|     6|   Brown|              2|       2010|         50|      |    -1|     null|   null|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+
```

Command took 5.46 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 8

```
1    # Right Join
2    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"right").show()
3    #Or
4    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"rightouter").show()
5
```

▶ (12) Spark Jobs

```
|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|
|     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|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+


|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|
|     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|
+------+--------+---------------+-----------+-----------+------+------+---------+-------+
```

Command took 4.91 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 9

```
1    # Left semi Join
2    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"leftsemi").show()
3
```

▶ (3) Spark Jobs

```
+------+--------+---------------+-----------+-----------+------+------+
|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|
+------+--------+---------------+-----------+-----------+------+------+
```

Command took 1.90 seconds -- by abhishek7621cse@gmail.com at 2/12/2024, 11:38:21 AM on My Cluster

Cmd 10

```
1    # Leftanti Join
2    empDF.join(deptDF,empDF.emp_dept_id ==  deptDF.dept_id,"leftanti").show()
3
```

▶ (6) Spark Jobs

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

Cmd 11

1

[Shift+Enter] to run
[Shift+Ctrl+Enter] to run selected text

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