

Step - 1: Problem Statement

26_groupby in pyspark

Write a pyspark code perform below function

- Write the query to get the department and department wise total(sum) salary from "EmployeeDetail" table.
- Write the query to get the department and department wise total(sum) salary, display it in ascending order according to salary.
- Write the query to get the department and department wise total(sum) salary, display it in descending order according to salary.
- Write the query to get the department, total no. of departments, total(sum) salary with respect to department from "EmployeeDetail" table.

Difficult Level: EASY

DataFrame:

Step - 2: Writing the pyspark code to solve the

```
# Creating Spark Session
from pyspark.sql import SparkSession
from pyspark.sql.types import
StructType,StructField,IntegerType,StringType
#creating spark session
spark = SparkSession. \
builder. \
config('spark.shuffle.useOldFetchProtocol', 'true'). \
config('spark.ui.port','0'). \
config("spark.sql.warehouse.dir", "/user/itv008042/warehouse"). \
enableHiveSupport(). \
master('yarn'). \
getOrCreate()
# Create a list of rows from the image
      [1, "Vikas", "Ahlawat", 600000.0, "2013-02-15 11:16:28.290", "IT", "Male"],
      [2, "nikita", "Jain", 530000.0, "2014-01-09 17:31:07.793", "HR", "Female"],
      [3, "Ashish", "Kumar", 1000000.0, "2014-01-09 10:05:07.793", "IT", "Male"],
      [4, "Nikhil", "Sharma", 480000.0, "2014-01-09 09:00:07.793", "HR", "Male"],
      [5, "anish", "kadian", 500000.0, "2014-01-09 09:31:07.793", "Payroll", "Male"],
```

emp_df=spark.createDataFrame(data,schema)

```
# 42. Write the query to get the department and department wise # total(sum) salary from "EmployeeDetail" table.

from pyspark.sql.functions import sum

emp_df.groupby(col('Department'))\
    .agg(sum('Salary').alias("sum_of_salary")).show()
```

```
+----+
|Department|sum_of_salary|
+----+
| HR| 1010000.0|
| Payroll| 500000.0|
| IT| 1600000.0|
```

```
+-----+
|Department|sum_of_salary|
+-----+
| Payroll| 500000.0|
| HR| 1010000.0|
| IT| 1600000.0|
```

```
| Department | sum_of_salary | 
| IT | 1600000.0 | 
| HR | 1010000.0 | 
| Payroll | 500000.0 |
```

