

```

# Databricks notebook source
import pyspark
from pyspark.sql import SparkSession
from pyspark.sql.functions import col,sum,avg,max

spark =
SparkSession.builder.appName('SparkByExamples.com').getOrCreate()

simpleData = [("James","Sales","NY",90000,34,10000),
  ("Michael","Sales","NY",86000,56,20000),
  ("Robert","Sales","CA",81000,30,23000),
  ("Maria","Finance","CA",90000,24,23000),
  ("Raman","Finance","CA",99000,40,24000),
  ("Scott","Finance","NY",83000,36,19000),
  ("Jen","Finance","NY",79000,53,15000),
  ("Jeff","Marketing","CA",80000,25,18000),
  ("Kumar","Marketing","NY",91000,50,21000)
]

schema = ["employee_name","department","state","salary","age","bonus"]
df = spark.createDataFrame(data=simpleData, schema = schema)
df.printSchema()
df.show(truncate=False)

df.groupBy("department").sum("salary").show(truncate=False)

df.groupBy("department").count().show(truncate=False)

df.groupBy("department","state") \
  .sum("salary","bonus") \
  .show(truncate=False)

df.groupBy("department") \
  .agg(sum("salary").alias("sum_salary"), \
    avg("salary").alias("avg_salary"), \
    sum("bonus").alias("sum_bonus"), \
    max("bonus").alias("max_bonus") \
  ) \
  .show(truncate=False)

df.groupBy("department") \
  .agg(sum("salary").alias("sum_salary"), \
    avg("salary").alias("avg_salary"), \
    sum("bonus").alias("sum_bonus"), \
    max("bonus").alias("max_bonus")) \
  .where(col("sum_bonus") >= 50000) \
  .show(truncate=False)

# COMMAND -----

```

