

```

# Databricks notebook source
import pyspark
from pyspark.sql import SparkSession
from pyspark.sql.functions import expr
spark =
SparkSession.builder.appName('SparkByExamples.com').getOrCreate()

data = [("James", "Sales", 3000), \
        ("Michael", "Sales", 4600), \
        ("Robert", "Sales", 4100), \
        ("Maria", "Finance", 3000), \
        ("James", "Sales", 3000), \
        ("Scott", "Finance", 3300), \
        ("Jen", "Finance", 3900), \
        ("Jeff", "Marketing", 3000), \
        ("Kumar", "Marketing", 2000), \
        ("Saif", "Sales", 4100) \
    ]
columns= ["employee_name", "department", "salary"]
df = spark.createDataFrame(data = data, schema = columns)
df.printSchema()
df.show(truncate=False)

distinctDF = df.distinct()
print("Distinct count: "+str(distinctDF.count()))
distinctDF.show(truncate=False)

df2 = df.dropDuplicates()
print("Distinct count: "+str(df2.count()))
df2.show(truncate=False)

dropDisDF = df.dropDuplicates(["department","salary"])
print("Distinct count of department salary : "+str(dropDisDF.count()))
dropDisDF.show(truncate=False)

# COMMAND -----

```