

```
In [1]: 1 # Initialization
2 import os
3 import sys
4
5 os.environ["SPARK_HOME"] = "/home/talentum/spark"
6 os.environ["PYLIB"] = os.environ["SPARK_HOME"] + "/python/lib"
7 # In below two lines, use /usr/bin/python2.7 if you want to use Python 2
8 os.environ["PYSPARK_PYTHON"] = "/usr/bin/python3.6"
9 os.environ["PYSPARK_DRIVER_PYTHON"] = "/usr/bin/python3"
10 sys.path.insert(0, os.environ["PYLIB"] + "/py4j-0.10.7-src.zip")
11 sys.path.insert(0, os.environ["PYLIB"] + "/pyspark.zip")
12
13 # NOTE: Whichever package you want mention here.
14 # os.environ['PYSPARK_SUBMIT_ARGS'] = '--packages com.databricks'
15 # os.environ['PYSPARK_SUBMIT_ARGS'] = '--packages org.apache.spark'
16 os.environ['PYSPARK_SUBMIT_ARGS'] = '--packages com.databricks:spark-sql-kafka-0-10_2.4:3.0.0'
17 # os.environ['PYSPARK_SUBMIT_ARGS'] = '--packages com.databricks'
```

```
In [2]: 1 #Entry point 2.x
2 from pyspark.sql import SparkSession
3 spark = SparkSession.builder.appName("Spark SQL basic example").getOrCreate()
4
5 # On yarn:
6 # spark = SparkSession.builder.appName("Spark SQL basic example")
7 # specify .master("yarn")
8
9 sc = spark.sparkContext
```

In [20]:

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StructType, StructField, StringType
3
4
5 # Initialize Spark Session
6 spark = SparkSession.builder \
7     .appName("UberRidesAnalysis") \
8     .getOrCreate()
9
10 # Define the schema for the Uber rides data
11 schema = StructType([
12     StructField("ride_id", StringType(), True),
13     StructField("date", StringType(), True),
14     StructField("pickup_location", StringType(), True),
15     StructField("drop_location", StringType(), True),
16     StructField("distance_km", DoubleType(), True),
17     StructField("fare_amount", DoubleType(), True),
18     StructField("payment_type", StringType(), True),
19     StructField("vehicle_type", StringType(), True)
20 ])
21
22 df = spark.read \
23     .option("header", "true") \
24     .schema(schema) \
25     .csv("file:///home/talentum/test-jupyter/bdt/uber_rides.csv")
26 # Display the dataframe
27 print("Data Preview:")
28 df.show()
29
30 print("\nSchema:")
31 df.printSchema()
32
33
```

Data Preview:

```
+-----+-----+-----+-----+-----+
|ride_id|      date|pickup_location|drop_location|distance_km|fare_
amount|payment_type|vehicle_type|
+-----+-----+-----+-----+-----+
|   R001|2024-01-10|      Downtown|       Suburb|    12.3|
|320.5|        Cash|          Sedan|          |
|   R002|2024-01-11|      Airport|  CityCenter|     8.5|
|250.0|        Card| Hatchback|          |
|   R003|2024-01-12|      Uptown|       Downtown|   15.0|
|410.75|        Card|          Sedan|          |
|   R004|2024-01-13|      Airport|       Suburb|   22.1|
|590.0|        UPI|          SUV|          |
|   R005|2024-01-14|      Downtown|      Uptown|    9.8|
|270.25|        Cash|          Sedan|          |
+-----+-----+-----+-----+-----+
```

Schema:

```
root
 |-- ride_id: string (nullable = true)
 |-- date: string (nullable = true)
```

```
|-- pickup_location: string (nullable = true)
|-- drop_location: string (nullable = true)
|-- distance_km: double (nullable = true)
|-- fare_amount: double (nullable = true)
|-- payment_type: string (nullable = true)
|-- vehicle_type: string (nullable = true)
```

In [19]:

```
1 #find the vehicle type with highest average fare
2 Q1_df=fare1_df.groupBy("vehicle_type") \
3     .avg("fare_amount") \
4     .orderBy("avg(fare_amount)", ascending=False)
5 Q1_df.show(1)
```

```
+-----+-----+
|vehicle_type|avg(fare_amount)|
+-----+-----+
|          SUV|      590.0|
+-----+-----+
only showing top 1 row
```

In [18]:

```
1 #find the top3 pickup locations by number of rides
2 Q2_df=fare1_df.groupBy("pickup_location") \
3     .count() \
4     .orderBy("count", ascending=False)
5 Q2_df.show(3)
```

```
+-----+-----+
|pickup_location|count|
+-----+-----+
|      Downtown|    2|
|        Airport|    2|
|       Uptown|    1|
+-----+-----+
```

In [21]:

```
1 #calculate the average distance per payment_type
2 Q3_df=fare1_df.groupBy("payment_type") \
3     .avg("distance_km")
4 Q3_df.show()
```

```
+-----+-----+
|payment_type|avg(distance_km)|
+-----+-----+
|      Card|      11.75|
|     Cash|      11.05|
|      UPI|      22.1|
+-----+-----+
```

```
In [24]: 1 #Find total earnings per day  
2  
3 Q4_df=df.groupBy("date") \  
4     .sum("fare_amount") \  
5     .orderBy("date") \  
6     .show()  
7
```

date	sum(fare_amount)
2024-01-10	320.5
2024-01-11	250.0
2024-01-12	410.75
2024-01-13	590.0
2024-01-14	270.25

```
In [ ]: 1
```