

databricks

# 14 DAYS

## AI CHALLENGE

---

**DAY 03**

---

**Topic:**

PySpark Transformations Deep Dive

**Challenge:**

1. Load full e-commerce dataset
2. Perform complex joins
3. Calculate running totals with window functions
4. Create derived features

+ New



Databricks day 0

Day 1

Databricks day 2

Databricks Day 3 ×



Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses

Data Engineering

Load Data of the full e-commerce into Spark so it can be analysed

▶ ✓ 12:08 PM (36s) 2

```
1 events = spark.read.csv(  
2     "/Volumes/workspace/e-commerce/e-commerce_data/2019-Oct.csv",  
3     header=True,  
4     inferSchema=True  
5 )  
6  
7 events.count()  
8 events.printSchema()
```

&gt; See performance (1)

&gt; events: pyspark.sql.connect.DataFrame = [event\_time: timestamp, event\_type: string ... 7 more fields]

root

```
|-- event_time: timestamp (nullable = true)  
|-- event_type: string (nullable = true)  
|-- product_id: integer (nullable = true)
```



+ New



Databricks day 0

Day 1

Databricks day 2

Databricks Day 3 ×

+



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 5 minutes ago



Run all

Serverless ▾

Schedule

Share



Advanced Full(Outer join)



12:24 PM (&lt;1s)

10: Cell 10

Python



```
1 import pyspark.sql.functions as F
2
3 full_sales = purchases.alias("p") \
4     .join(products.alias("pr"), on="product_id", how="outer") \
5     .withColumn("data_quality_flag",
6                 F.when(F.col("p.product_id").isNull(), "Missing in Purchases")
7                     .when(F.col("pr.product_id").isNull(), "Missing in Products")
8                     .otherwise("Valid")) \
9     .withColumn("revenue_flag",
10                F.when(F.col("price") > 1000, "High Value")
11                    .when(F.col("price") > 500, "Medium Value")
12                    .otherwise("Low Value")) \
13     .withColumn("event_hour", F.hour("event_time")) \
14     .withColumn("processing_date", F.current_date())
```

&gt; full\_sales: pyspark.sql.connect.DataFrame = [product\_id: integer, event\_time: timestamp ... 13 more fields]

+ New

Home

Workspace

Recents

Catalog

Jobs &amp; Pipelines

Compute

Marketplace

SQL

SQL Editor

Queries

Dashboards

Genie

Alerts

Query History

SQL Warehouses

Data Engineering



Databricks day 0

Day 1

Databricks day 2

Databricks Day 3 ×



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



Last edit was 6 minutes ago



Run all

Serverless ▾

Schedule

Share



File

Edit

View

Run

Help

Python ▾

Tabs: ON ▾



