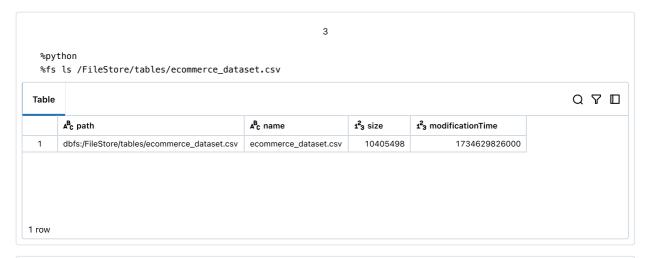


FileStore/tables/ecommerce_dataset.csv



4

ecommercedatadf=spark.read.format("csv").option("header", "true").option("inferSchema", "true").load("/FileStore/t

▶ ■ ecommercedatadf: pyspark.sql.dataframe.DataFrame = [Product ID: string, Product Name: string ... 14 more fields]

				5					
ecommercedat	tadf.show()								
	+	+	+	+	+ -+		+-		+ -+
	Product Name ion Customer Gend		•		•		•	-	
				-					
P6879	+ Jacket	+ Apparel	+ 53 . 85	5	-+ 15	150	+ S535		-+ 35-44
ew York, USA	Male	23.32	Stan	dard	4.49	Yes		56	
P5132	Camera	Electronics	761.26	10	15	224	S583		25-34
ondon, UK	Female	20.88	0vernigh	t :	16.11	No		79	
P2941	Sneakers	Footwear	1756.76	5	8	468	S118		25-34
okyo, Japan	Non-Binary	16.43	Stand	ard	4.93	No		40	
P8545	Cookbooks	Books	295.24	10	15	25	S104		18-24
aris, France	Female	27.49	Stan	dard	1.31	No		93	
P4594	Camera	Electronics	832.0	10	12	340	S331		55
okyo, Japan	Male	45.93	0verni	ght	4.37	No		56	
P1388	Non-Fiction	Books	584.19	15	8	204	S523		45-54
ingapore	Female	40.12	Express	19	9.03	No	9	1	
P7313 I	Running Shoes	Footwear	1343.95	0	10	493	S878		18-24
ydney, Austra	lia Fema	le 35	.91 0	vernight	17.73	i Ye	es	41	
P1060	Blender Hom	e Appliances!	1873.521	0	15	349	S3961		18-24

```
6
  # Replace invalid characters in column names
  cleaned_columns = [col.strip().replace(" ", "_").replace(".", "_") for col in ecommercedatadf.columns]
  ecommercedatadf = ecommercedatadf.toDF(*cleaned_columns)
  # Show the updated schema
  ecommercedatadf.printSchema()
  # Save raw data as a Delta table
  ecommercedatadf.write.format("delta").mode("overwrite").saveAsTable("ecommerce_bronze")
  print("Raw data saved to Bronze Layer!")
▶ ■ ecommercedatadf: pyspark.sql.dataframe.DataFrame = [Product_ID: string, Product_Name: string ... 14 more fields]
root
|-- Product_ID: string (nullable = true)
|-- Product_Name: string (nullable = true)
|-- Category: string (nullable = true)
|-- Price: double (nullable = true)
|-- Discount: integer (nullable = true)
|-- Tax_Rate: integer (nullable = true)
 |-- Stock_Level: integer (nullable = true)
 |-- Supplier_ID: string (nullable = true)
 |-- Customer_Age_Group: string (nullable = true)
 |-- Customer_Location: string (nullable = true)
 |-- Customer_Gender: string (nullable = true)
|-- Shipping_Cost: double (nullable = true)
 |-- Shipping_Method: string (nullable = true)
|-- Return_Rate: double (nullable = true)
|-- Seasonality: string (nullable = true)
|-- Popularity_Index: integer (nullable = true)
```

Raw data saved to Bronze Layer!

Query the Bronze Layer Delta table and display the first 10 rows bronze_layer_data = spark.sql("SELECT * FROM ecommerce_bronze LIMIT 10")

bronze_layer_data.show()

▶ ■ bronze_laye	er_data: pyspark.sq	.dataframe.DataFrame =	[Product_ID: st	ring, Product_N	ame: string 1	4 more fie	elds]		
·	+		+	+	+	+-		+-	
Product_ID Pr ustomer_Locatio 	roduct_Name on Customer_Gend	+	Discount Ta	x_Rate Stock_ d Return_Rate 	Level Suppl Seasonalit	ier_ID C y Popula	Customer_Agarity_Index	e_Group +	
P6879	•	++ Apparel 53.85		•	150	-+ S535		+ 35-44	
lew York, USA	Male	23.32	Standard	4.49	Yes		56		
P5132	Camera	Electronics 761.26	10	15	224	S583		25-34	
ondon, UK	Female	20.88 0ve	ernight	16.11	No		79		
P2941	Sneakers	Footwear 1756.76	5	8	468	S118		25-34	
okyo, Japan	Non-Binary	16.43	Standard	4.93	No		40		
P8545	Cookbooks	Books 295.24	10	15	25	S104		18-24	
aris, France	Female	27.49	Standard	1.31	No		93		
P4594	Camera	Electronics 832.0	10	12	340	S331		55	
okyo, Japan	Male	45.93)vernight	4.37	No		56		
P1388 N	on-Fiction	Books 584.19	15	8	204	S523		45-54	
ingapore	Female	40.12 Ex	press	19.03	No		91		
P7313 Run	ning Shoes	Footwear 1343.95	0	10	493	S878		18-24	
dney, Australi	.a Fema	le 35.91	Overnigh	t 17.73	Ye	s	41	1	
P1060	Blender Home	Appliances 1873.52	0	15	349	S396		18-24	

```
from pyspark.sql.functions import col
  # Clean the data
  cleaned_data = ecommercedatadf \
      .filter((col("Price").isNotNull()) & (col("Stock_Level").isNotNull())) \
      .filter((col("Price") > 0) & (col("Stock_Level") > 0)) \
      .withColumn("Price", col("Price").cast("double")) \
      .withColumn("Stock_Level", col("Stock_Level").cast("int"))
  # Show the cleaned data
  cleaned_data.show()
  # Save cleaned data as Delta table
  cleaned_data.write.format("delta").mode("overwrite").saveAsTable("ecommerce_silver")
  print("Cleaned data saved to Silver Layer!")
▶ ■ cleaned_data: pyspark.sql.dataframe.DataFrame = [Product_ID: string, Product_Name: string ... 14 more fields]
                              Category| Price|Discount|Tax_Rate|Stock_Level|Supplier_ID|Customer_Age_Group|
|Product ID| Product Name|
Customer_Location|Customer_Gender|Shipping_Cost|Shipping_Method|Return_Rate|Seasonality|Popularity_Index|
     P6879|
                 Jacket|
                              Apparel| 53.85| 5|
                                                                     150|
                                                                                 S535|
                                                                                                  35-44|
New York, USA|
                   Male|
                                                                        Yes|
                                 23.32| Standard|
                                                             4.491
                                                                                          56 I
                                                           15|
                 Camera| Electronics| 761.26| 10|
    P5132|
                                                                       224|
                                                                                 S583|
                                                                                                  25-34|
                Female|
                                                                                       79|
London, UK|
                             20.88|
                                          Overnight| 16.11|
                                                                       Nol
  P2941|
                Sneakers|
                              Footwear|1756.76| 5| 8|
                                                                       468|
                                                                                 S118|
                                                                                                  25-34|
                               16.43| Standard|
Books| 295.24| 10|
27.49| Standard|
Tokyo, Japan|
                Non-Binary|
                                                            4.93|
                                                                        No|
                                                                                         40|
   P8545|
                Cookbooks|
                                                                       25|
                                                                                 S104|
                                                                                                  18-24|
                                                            15|
Paris, France|
                   Female|
                                                             1.31|
                                                                        No |
                                                                                         93|
                                                                       340|
    P4594|
                   Camera|
                             Electronics| 832.0| 10|
                                                             12|
                                                                                 S331|
                   Male|
                             45.93| Overnight|
                                                            4.371
                                                                       Nol
                                                                                         56 I
Tokyo, Japan|
                                                            8|
              Non-Fiction|
    P1388|
                                 Books | 584.19 | 15 |
                                                                       204|
                                                                                 S523|
                                                                                                  45-54|
Singapore|
               Female|
                               40.12|
                                           Express
                                                        19.03|
                                                                      No|
                                                                                      91|
                                                                                 S878|
| P7313| Running Shoes|
                              Footwear|1343.95|
                                                     0| 10|
                                                                                                  18-24|
                                                                       493 I
Sydney, Australia| Female| 35.91|
                                                 Overnight|
                                                              17.73|
| P1060|
                  Blender|Home Appliances|1873.52|
                                                     0|
                                                                       349|
                                                                                 S396|
                                                             15|
                                                                                                  18-24|
```

Gold Layer

Perform aggregations to create meaningful datasets and save them as Delta tables.

Total Sales by Category

```
10
   # Total sales by category
   total_sales_by_category = cleaned_data.groupBy("Category") \
       .agg({"Price": "sum"}) \
       .withColumnRenamed("sum(Price)", "Total_Sales")
   # Save as Delta table
   total_sales_by_category.write.format("delta").mode("overwrite").saveAsTable("total_sales_by_category")
   print("Total Sales by Category saved to Gold Layer!")
   total_sales_by_category.show()
 ▶ ■ total_sales_by_category: pyspark.sql.dataframe.DataFrame = [Category: string, Total_Sales: double]
Total Sales by Category saved to Gold Layer!
        Category|
                           Total_Sales|
        Apparel | 2.0063856269999962E7 |
     Electronics | 2.010186938999998E7|
        Footwear | 2.0161514950000018E7 |
           Books | 2.0276547950000018E7 |
|Home Appliances|1.9960686890000023E7|
```

Average Price by Category

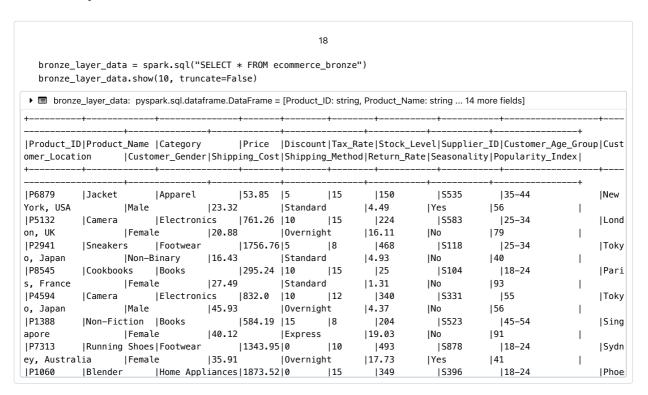
```
12
  # Average price by category
  avg_price_by_category = cleaned_data.groupBy("Category") \
       .agg({"Price": "avg"}) \
       .withColumnRenamed("avg(Price)", "Avg_Price")
  # Save as Delta table
  avg_price_by_category.write.format("delta").mode("overwrite").saveAsTable("avg_price_by_category")
  avg_price_by_category.show()
▶ ■ avg_price_by_category: pyspark.sql.dataframe.DataFrame = [Category: string, Avg_Price: double]
       Category|
                        Avg_Price|
        Apparel|1001.6902780828738|
    Electronics | 1009.3833487321103 |
        Footwear | 1004.7099691035041 |
           Books | 1009.6374022805367 |
|Home Appliances|1012.9757366150735|
```

Total Quantity Sold by Shipping Method

Monthly Sales Trends

Verifying Layers

Bronze Layer



Verifying Layers

Silver Layer

```
silver_layer_data = spark.sql("SELECT * FROM ecommerce_silver")
silver_layer_data.show(10, truncate=False)

| silver_layer_data: pyspark.sql.dataframe.DataFrame = [Product_ID: string, Product_Name: string ... 14 more fields]
| truncate=False |
```

Product_I omer_Locat	D Product_Nam ion Cus	, ,		•				er_ID Customer_A ity Popularity_I	
+	-+	+		+	+	+		+	+ +
P6879	Jacket	Apparel	53.85	i 5	15	150	S535	35–44	New
York, USA	Mal	e	23.32	Stand	ard	4.49	Yes	56	1
P5132	Camera	Electron	ics 761.2	26 10	15	224	S583	25-34	Lon
on, UK	Fem	ale	20.88	Overn	ight	16.11	No	79	1
P2941	Sneakers	Footwear	1756.	76 5	8	468	S118	25-34	Tok
o, Japan	Nor	-Binary	16.43	Stand	ard	4.93	No	40	1
P8545	Cookbooks	Books	295.2	24 10	15	25	5104	18-24	Par
s, France	Fen	ale	27.49	Stand	ard	1.31	No	93	I
P4594	Camera	Electron	ics 832.6	10	12	340	S331	55	Tok
o, Japan	Mal	e	45.93	Overn	ight	4.37	No	56	1
P1388	Non-Fiction	Books	584.1	9 15	8	204	S523	45-54	Sin
apore	Fen	ale	40.12	Expre	SS	19.03	No	91	1
P7313	Running Sho	es Footwear	1343	95 0	10	493	5878	18-24	Syd
ey, Austra	lia Fem	ale	35.91	Overn	ight	17.73	Yes	41	1
I D1 060	IRlander	IHome Ann	liances 1973	5210	115	1340	16306	112_24	I Pho

Verifying Layers

Gold Layers

Machine learning Model to Predict Stock Levels (Regression model)

Objective: To Predict the Stock_Level for a product based on features like Price, Discount, and Tax_Rate.

```
26
   from pyspark.ml.feature import VectorAssembler
   from pyspark.ml.regression import LinearRegression
   from pyspark.ml.evaluation import RegressionEvaluator
   # Prepare data for regression
   assembler = VectorAssembler(inputCols=["Price", "Discount", "Tax_Rate"], outputCol="features")
   regression_data = assembler.transform(cleaned_data).select("features", "Stock_Level")
   # Split data into training and testing sets
   train_data, test_data = regression_data.randomSplit([0.8, 0.2], seed=42)
   # Train Linear Regression model
   lr = LinearRegression(featuresCol="features", labelCol="Stock_Level")
   lr_model = lr.fit(train_data)
   # Evaluate model
   predictions = lr_model.transform(test_data)
   evaluator = RegressionEvaluator(labelCol="Stock_Level", predictionCol="prediction", metricName="rmse")
   rmse = evaluator.evaluate(predictions)
   print(f"Root Mean Squared Error (RMSE): {rmse}")
   predictions.show(10, truncate=False)
 ▶ ■ predictions: pyspark.sql.dataframe.DataFrame
 ▶ ■ regression_data: pyspark.sql.dataframe.DataFrame
 ▶ ■ test_data: pyspark.sql.dataframe.DataFrame
 ▶ ■ train_data: pyspark.sql.dataframe.DataFrame
 Loading widget. This should take less than 30 seconds.
 Loading widget. This should take less than 30 seconds.
Root Mean Squared Error (RMSE): 144.1967092226598
|features
                 |Stock Level|prediction
|[10.23,25.0,8.0] |207
                               |249.11485704564956|
                              |249.11522916082723|
|[10.46,25.0,8.0] |114
                             |248.56769425997706|
|[10.58,5.0,8.0] |480
[10.78,25.0,10.0]|210
                            |249.72145366904982|
|249.64322719336815|
|249.91762562170322|
|[11.23,0.0,12.0] |5
|[11.56,10.0,12.0]|312
|[11.77,15.0,8.0] |293
|[12.14,15.0,12.0]|128
                               |248.84348407549822|
                               |250.05549626108262|
|[12.64,15.0,5.0] |225
                               |247.93633146746808|
|[12.7,15.0,12.0] |165
                               |250.05640228064567|
only showing top 10 rows
```

Machine Learning Model: TO Classify Products by Popularity (Classification Model)

Objective: Classify products as "Popular" or "Not Popular" based on their Popularity_Index.

```
from pyspark.ml.classification import LogisticRegression
   from pyspark.sql.functions import when
   from\ pyspark.ml.evaluation\ import\ Multiclass Classification Evaluator
   # Add binary column for popularity classification
   {\tt classification\_data = cleaned\_data.withColumn("is\_popular", when (col("Popularity\_Index") >= 50, 1).otherwise (0))}
   # Prepare features for classification
   assembler = VectorAssembler(inputCols=["Price", "Discount", "Tax_Rate", "Stock_Level"], outputCol="features")
  classification_data = assembler.transform(classification_data).select("features", "is_popular")
  # Split data into training and testing sets
  train_data, test_data = classification_data.randomSplit([0.8, 0.2], seed=42)
   # Train Logistic Regression model
   lr = LogisticRegression(featuresCol="features", labelCol="is_popular", maxIter=10)
   lr_model = lr.fit(train_data)
  # Evaluate model
   predictions = lr_model.transform(test_data)
   # Use MulticlassClassificationEvaluator for accuracy
   evaluator = MulticlassClassificationEvaluator(labelCol="is_popular", predictionCol="prediction", metricName="accur
   accuracy = evaluator.evaluate(predictions)
  print(f"Accuracy: {accuracy}")
  predictions.show(10, truncate=False)
▶ ■ classification_data: pyspark.sql.dataframe.DataFrame
 ▶ ■ predictions: pvspark.sql.dataframe.DataFrame
 ▶ ■ test_data: pyspark.sql.dataframe.DataFrame
 ▶ ■ train_data: pyspark.sql.dataframe.DataFrame
 Loading widget. This should take less than 30 seconds.
 Loading widget. This should take less than 30 seconds.
Accuracy: 0.5021143777688281
                        |is_popular|rawPrediction
                                                                                     |probability
|features
|prediction|
|[10.23,25.0,8.0,207.0] |1
                                   |[0.009160758245145367,-0.009160758245145367] |[0.5022901735454588,0.4977098264
545412] |0.0
                                   |[0.009531934371093056,-0.009531934371093056] |[0.5023829655502221,0.4976170344
|[10.46,25.0,8.0,114.0] |0
497779] |0.0
|[10.58,5.0,8.0,480.0] |1
                                    |[-0.010561687068939042,0.010561687068939042] |[0.49735960277724345,0.502640397
2227566111.0
|[10.78,25.0,10.0,210.0]|1
                                    [0.00739821539632827,-0.00739821539632827]
                                                                                    |[0.5018495454130679,0.4981504545
869321] |0.0
|[11.23,0.0,12.0,5.0] |1
                                    |[-0.016798815313672933,0.016798815313672933] |[0.4958003949318983,0.5041996050
681017] |1.0
                                    |[-0.008729998488339009,0.008729998488339009] |[0.49781751423902365,0.502182485
|[11.56,10.0,12.0,312.0]|1
7609763] | 1.0
                                    |[-5.162029514218184E-4,5.162029514218184E-4] |[0.49987094926501013,0.500129050
|[11.77,15.0,8.0,293.0] |1
7349898]|1.0
```

Machine Learining Model: Customer Segmentation Using K-Means Clustering

Objective: Group customers into segments based on features like total revenue, order count, and average order value.

Loading widget. This should take less than 30 seconds.

Loading widget. This should take less than 30 seconds.

Customer_Location	Total Povenue	lordor Count	+ Avg_Order_Value	Cluster
+	+	+	+	+
Singapore	1.7024310824700003E9	6695	 1010.3822001493661	0
Toronto, Canada	1.672599277999999E9	6585	1013.0271556567948	1
Mumbai, India	1.669635879249999E9	6767	1005.132818087779	1
Chicago, USA	1.693613506369999E9	6598	1007.4677826614122	0
Sydney, Australia	1.654938673480001E9	6527	1002.6644921096992	1
Dubai, UAE	1.692078486000001E9	6759	1005.3609675987568	0
Phoenix, USA	1.6739677950999987E9	6643	1010.5593632394999	1
London, UK	1.653133941290001E9	6514	1006.0021998771882	1
Berlin, Germany	1.667030874180001E9	6638	1002.7849789093117	1
Los Angeles, USA	1.6816780004499977E9	6635	1013.1032313489078	0
+	+	+	+	+

only showing top 10 rows

Silhouette Score: 0.7126998999023788