**RIBA MAUBANE**
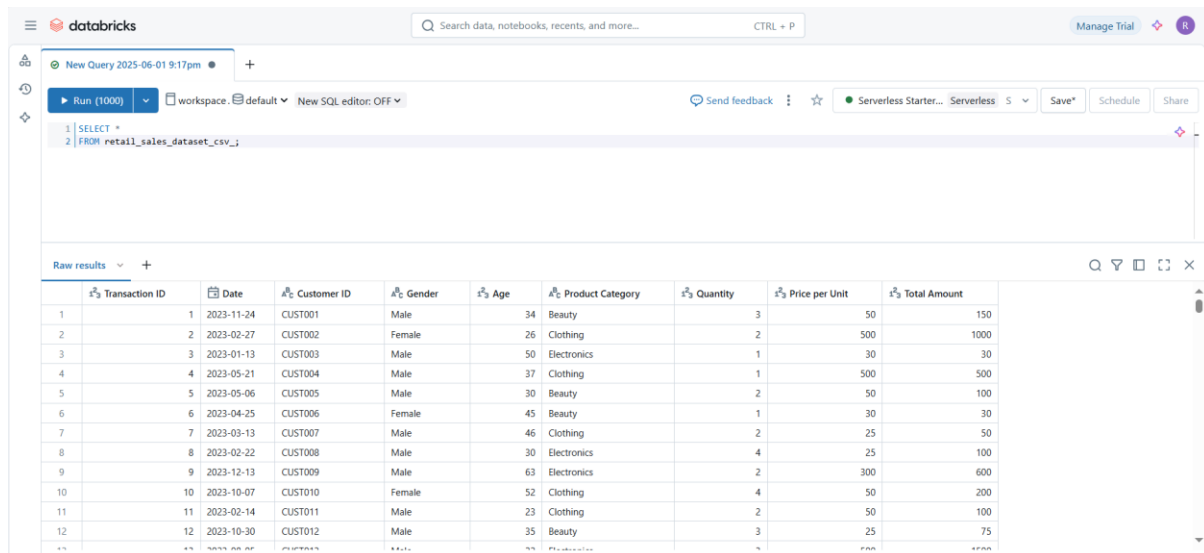**DATABRICKS PRACTICAL 1**

QUESTION 1

**SELECT Statement**

1. Display all columns for all transactions.
Expected output: All columns



2. Display only the Transaction ID, Date, and Customer ID for all records.
Expected output: Transaction ID, Date, Customer ID

**SELECT DISTINCT Statement**

3.   Display all the distinct product categories in the dataset.
Expected output: Product Category



4.   Display all the distinct gender values in the dataset.
Expected output: Gender

**WHERE Clause**

5.   Display all transactions where the Age is greater than 40.
Expected output: All columns



6. Display all transactions where the Price per Unit is between 100 and 500.
Expected output: All columns

7. Display all transactions where the Product Category is either 'Beauty' or 'Electronics'.
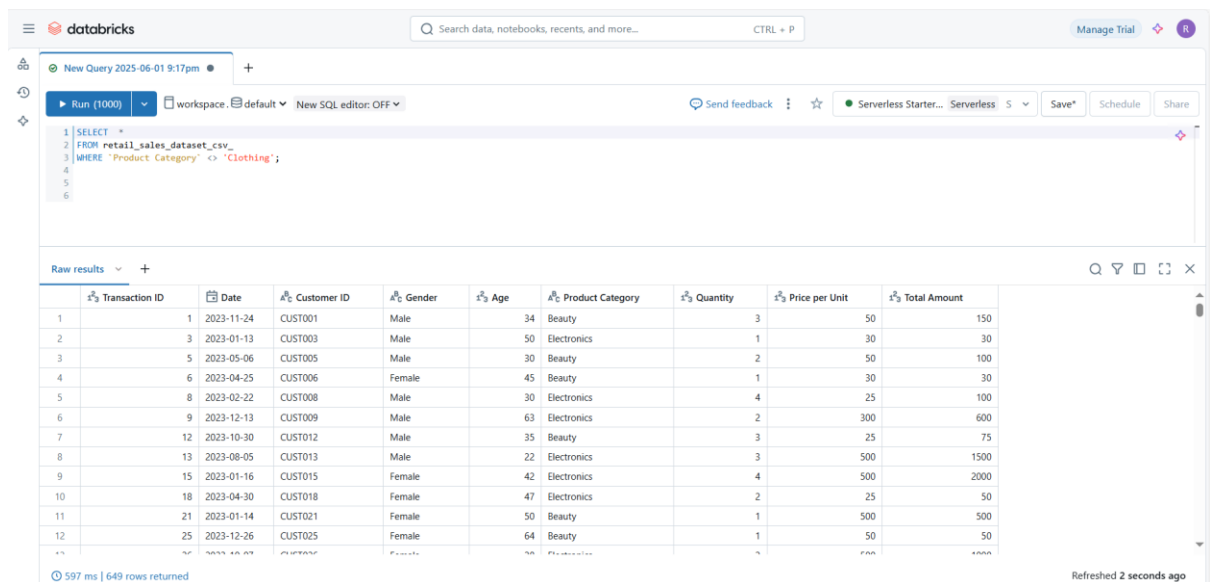Expected output: All columns



8. Display all transactions where the Product Category is not 'Clothing'.
Expected output: All columns

9. Display all transactions where the Quantity is greater than or equal to 3.
Expected output: All columns



**Aggregate Functions**

10. Count the total number of transactions.
Expected output: Total_Transactions

**11.** Find the average Age of customers.
Expected output: Average_Age



**12.** Find the total quantity of products sold.
Expected output: Total_Quantity



**13.** Find the maximum Total Amount spent in a single transaction.
Expected output: Max_Total_Amount

14. Find the minimum Price per Unit in the dataset.
Expected output: Min_Price_per_Unit



**GROUP BY Statement**

15. Find the number of transactions per Product Category.
Expected output: Product Category, Transaction_Count



16. <u>Find the total revenue (Total Amount) per gender.</u>
Expected output: Gender, Total_Revenue

17. <u>Find the average Price per Unit per product category.</u>
Expected output: Product Category, Average_Price



**HAVING Clause**

18. <u>Find the total revenue per product category where total revenue is greater than 10,000.</u> Expected output: Product Category, Total_Revenue
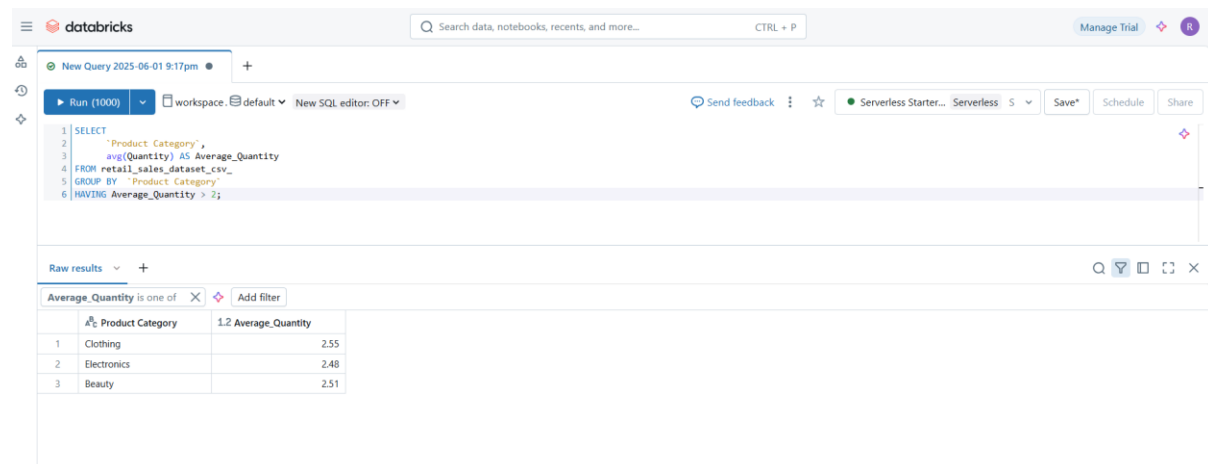


19. <u>Find the average quantity per product category where the average is more than 2.</u>
Expected output: Product Category, Average_Quantity

**CASE Statement**

20. Display a column called Spending_Level that shows 'High' if Total Amount > 1000, otherwise 'Low'.
Expected output: Transaction ID, Total Amount, Spending_Level



21. Display a new column called Age_Group that labels customers as: • 'Youth' if Age < 30 • 'Adult' if Age is between 30 and 59 • 'Senior' if Age >= 60
Expected output: Customer ID, Age, Age_Group