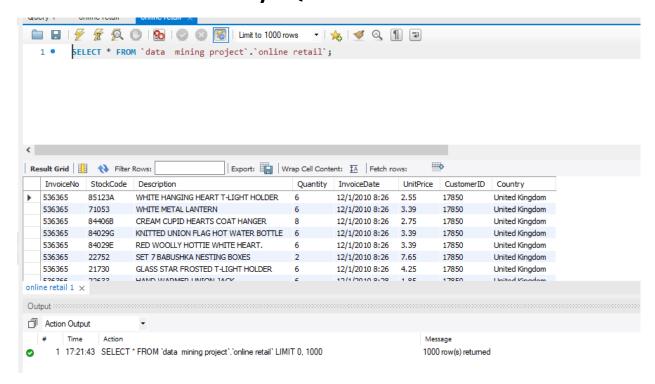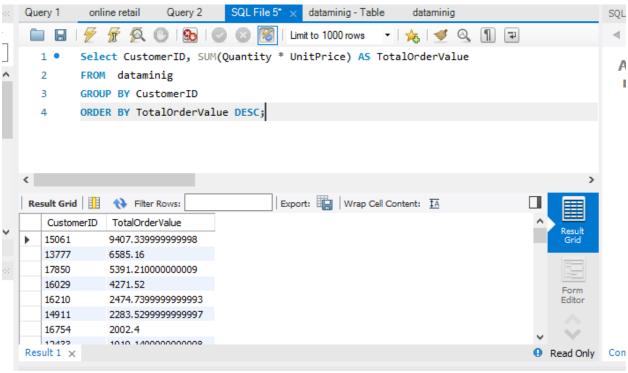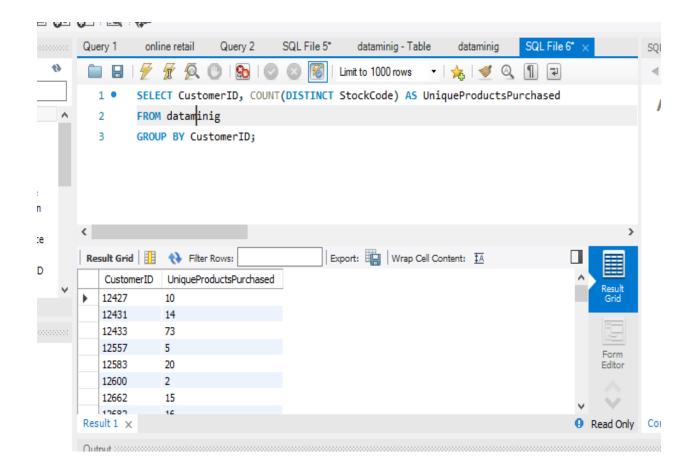# DATA MINING PROJECT

## 1-METADATA

### a. Define metadata in my SQL workbench ?

## b. What is the distribution of order values across all customers in the dataset?

- c. How many unique products has each customer purchased?

## d. Which customers have only made a single purchase from the company ?

```sql
1   SELECT CustomerID
2   FROM dataminig
3   GROUP BY CustomerID
4   HAVING COUNT(DISTINCT InvoiceNo) = 1;
```

Result Grid | Filter Rows: | Export: | Wrap Cell Content: ĀA

| CustomerID |
|---|
| 12427 |
| 12431 |
| 12433 |
| 12557 |
| 12583 |
| 12600 |
| 12662 |
| 12682 |

Result Grid

Form Editor

dataminig 1 ×

Read Only    Cc

Output

# e. Which products are most commonly purchased together by customers in the dataset?

# Advance Queries

## 1. Customer Segmentation by Purchase Frequency

```sql
  1 ● SELECT CustomerID,
  2       CASE
  3           WHEN COUNT(DISTINCT InvoiceDate) >= 10 THEN 'High Frequency'
  4           WHEN COUNT(DISTINCT InvoiceDate) >= 5 THEN 'Medium Frequency'
  5           ELSE 'Low Frequency'
  6       END AS PurchaseFrequencySegment
  7  FROM dataminig
  8  GROUP BY CustomerID;
  9
```

| CustomerID | PurchaseFrequencySegment |
|------------|--------------------------|
| 12431 | Low Frequency |
| 12433 | Low Frequency |
| 12557 | Low Frequency |
| 12583 | Low Frequency |
| 12600 | Low Frequency |
| 12662 | Low Frequency |

Result 1 ✕

# 2. Average Order Value by Country

# 3. Customer Churn Analysis



```sql
SELECT CustomerID
FROM dataminig
WHERE InvoiceDate <= DATE_SUB(NOW(), INTERVAL 6 MONTH)
GROUP BY CustomerID;
```

| CustomerID |
|------------|
| 17850 |
| 13047 |
| 12583 |
| 13748 |
| 15100 |
| 15291 |
| 14688 |
| 17809 |

# 4. Product Affinity Analysis

```
1 ●  SELECT DATE_FORMAT(InvoiceDate, '%Y-%m') AS Month,
2          COUNT(DISTINCT CustomerID) AS UniqueCustomers,
3          SUM(Quantity * UnitPrice) AS TotalSales
4      FROM dataminig
5      GROUP BY Month
6      ORDER BY Month;
7
```

Limit to 1000 rows

Result Grid    Filter Rows:                Export:    Wrap Cell Content:

| Month | UniqueCustomers | TotalSales |
|-------|-----------------|------------|
| NULL  | 229             | 113902.76999999864 |

Result Grid
Form Editor
Field Types

Result 1 ✕

🛈 Read Only

# 5. Time-based Analysis

```sql
SELECT a.StockCode AS Product1, b.StockCode AS Product2, COUNT(*) AS Frequency
FROM dataminig a
JOIN dataminig b ON a.InvoiceNo = b.InvoiceNo AND a.StockCode < b.StockCode
GROUP BY Product1, Product2
ORDER BY Frequency DESC;
```

| Product1 | Product2 | Frequency |
|----------|----------|-----------|
| 22632 | 22633 | 26 |
| 84029E | 84029G | 21 |
| 84029E | 85123A | 20 |
| 84029G | 85123A | 20 |
| 21730 | 85123A | 18 |
| 71053 | 85123A | 18 |
| 22752 | 85123A | 17 |
| 21730 | 71053 | 17 |
| 21730 | 84029G | 17 |
| 22752 | 84029G | 17 |