# Data Mining Project
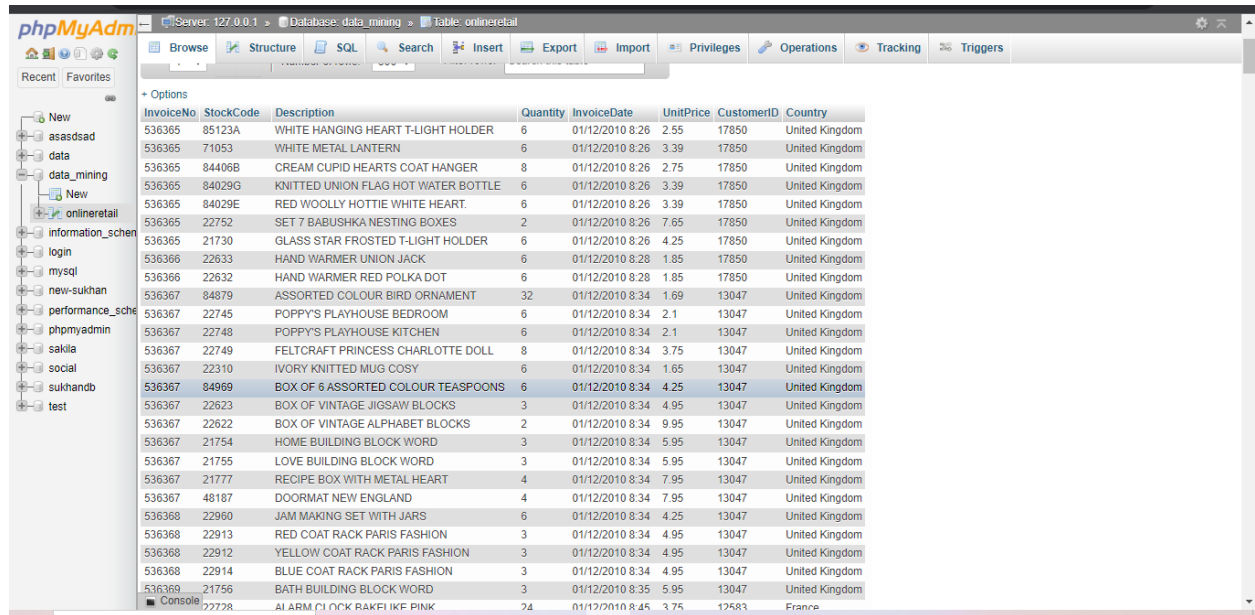
**Project Title**: Online Retail Segmentation.

Task1:

- Define meta data in mysql workbench
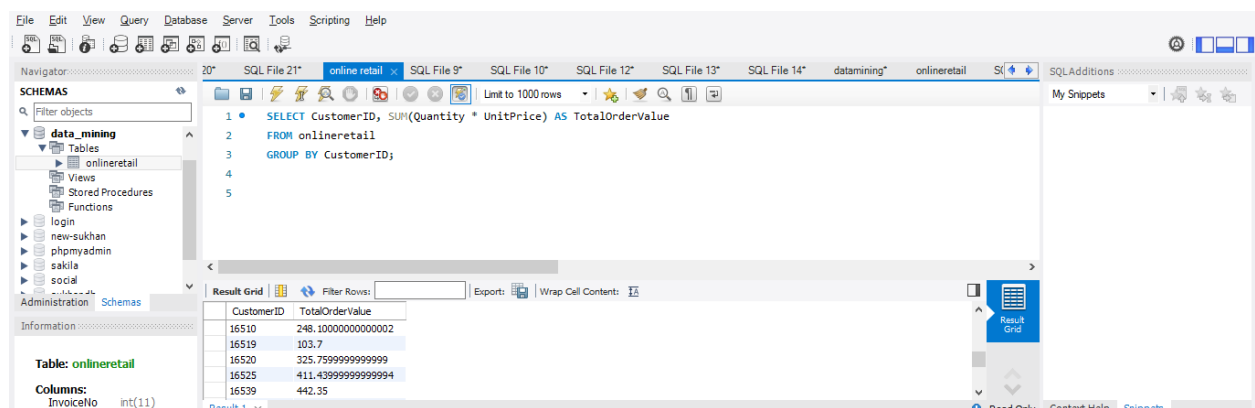


Task2:

- What is the distribution of order values across all customers in the dataset?

Task3:

- How many unique products has each customer purchased?



Task4:

- Which customers have only made a single purchase from the company?

Task5:

- Which products are most commonly purchased together by customers in the dataset?



I took this screenshot from localhost phpmyadmin because I was facing an issue in mysql workbench server.

## Advance Queries

### 1. Customer Segmentation by Purchase Frequency

Group customers into segments based on their purchase frequency, such as high, medium, and low frequency customers. This can help you identify your most loyal customers and those who need more attention.

Code:

SELECT

    CustomerID,

    CASE

        WHEN COUNT(DISTINCT InvoiceNo) >= StockCode THEN 'High'

        WHEN COUNT(DISTINCT InvoiceNo) >= StockCode THEN 'Medium'

        ELSE 'Low'

    END AS PurchaseFrequencySegment

FROM

    onlineretail

GROUP BY

    CustomerID;

## 2. Average Order Value by Country

Calculate the average order value for each country to identify where your most valuable customers are located.



## 3. Customer Churn Analysis

Identify customers who haven't made a purchase in a specific period (e.g., last 6 months) to assess churn.

### 4. Product Affinity Analysis

Determine which products are often purchased together by calculating the correlation between product purchases.



### 5. Time-based Analysis

Explore trends in customer behavior over time, such as monthly or quarterly sales patterns.



# The End