

# Final Project – Data Warehousing

## 1. Executive Summary

In today's data-driven retail environment, organizations must leverage extensive datasets to gain deeper insights into their operations, customers, and products. This project explores a scenario inspired by a mid-sized online retail company seeking to enhance profitability, understand customer buying behaviors, and drive strategic decision-making. By integrating transactional (OLTP) and analytical (OLAP) models, as well as performing dimensionally-modeled data warehousing, we develop a set of analyses and visualizations that inform executive-level business strategies.

The work focuses on several key areas of inquiry. First, we examine whether a customer's commute distance to a store location correlates with their average spending, potentially guiding location-based promotions or personalized marketing. Next, we investigate month-over-month sales growth trends to identify seasonal patterns and operational improvements. We also assess the impact of customer demographics—such as income or marital status—on overall sales, enabling more targeted marketing efforts. Further, the analysis delves into customer lifetime value to identify the most profitable customer segments for long-term retention strategies. Lastly, we explore which products yield the highest profit margins and segment customers by their yearly spending patterns, guiding product selection, inventory management, and personalized recommendations.

Through dimensional modeling, analytic SQL queries, and potential visualization platforms, this project offers a coherent strategy for transforming raw data into actionable business insights. The ultimate goal is to improve decision-making, allocate marketing resources effectively, and drive sustained profitability.

## 2. Problem Statement

As online retail competition intensifies, decision-makers face a complex web of challenges: maintaining growth, enhancing profitability, and consistently meeting customer expectations. The complexity is compounded by varied customer demographics, dynamic product lines, and evolving market conditions.

This project addresses the core challenge of turning disparate, large-scale transactional data into meaningful insights that inform strategic actions. Specifically, the organization seeks to answer critical questions:

**Customer Behavior & Geography:** Does the distance customers commute to shop correlate with their average sales amounts, and how can this inform tailored marketing or logistics strategies?

**Revenue Growth Patterns:** Can identifying month-over-month sales trends help forecast demand, optimize inventory, and schedule targeted promotions to drive steady growth?

**Customer Demographics:** How do demographic factors—such as marital status, gender, income, or education level—influence purchasing behaviors and overall sales, and what does this mean for segmentation and personalization?

**Customer Lifetime Value (CLV):** Which customers generate the most long-term value, and how can retention strategies be refined to nurture these relationships over time?

**Product Profitability:** Which products yield the highest margins, and how can product mix adjustments or supplier negotiations improve profitability?

**Customer Spending Segmentation:** How can we categorize customers by their yearly spending patterns to inform differentiated marketing approaches and loyalty programs?

By addressing these questions through a robust data warehousing architecture, dimensional modeling, and analytic queries, the organization aims to reduce uncertainty, identify growth opportunities, and ultimately make data-driven decisions that enhance long-term success.

### 3. Literature Review

Implementing data warehousing and business intelligence solutions centers on key concepts like dimensional modeling, data cube design, and the integration of OLTP and OLAP systems. Foundational authors such as Ralph Kimball and Bill Inmon have established best practices that inform essential architectural decisions for advanced analytics.

Dimensional modeling, as advocated by Kimball and Ross (2013) in *The Data Warehouse Toolkit*, differentiates analytical storage structures from operational relational schemas. Star and snowflake schemas organize data around conformed dimensions and fact tables, enhancing clarity and query efficiency. Key principles include defining the "grain" of fact tables, identifying relevant "dimensions" (e.g., customers, products, time), and specifying "facts"

(quantitative measures like sales amounts). These choices facilitate complex queries and integrated analysis across fact tables.

Data cube design extends dimensional modeling by providing a framework for multi-dimensional analysis. Cubes aggregate facts across multiple dimensions, enabling operations like drilling down, rolling up, and slicing data to uncover patterns and trends. This simplifies navigation of complex datasets and supports diverse decision-making scenarios, from growth analyses to customer segmentation. Gray et al. (1997) highlight that pre-aggregated data cubes enhance OLAP performance by speeding up query responses.

Effective SQL aggregation is crucial for leveraging dimensional models and data cubes. Fact tables optimized for SQL aggregate functions (SUM, AVG, COUNT) enable the generation of key metrics like monthly revenue trends and customer lifetime value. Well-structured schemas and appropriate indexing facilitate efficient SQL-based aggregation, supporting rapid responses for BI dashboards and ad-hoc queries.

The distinction between OLTP and OLAP environments is critical. OLTP systems manage day-to-day transactions using normalized schemas for data integrity and minimal redundancy. In contrast, OLAP systems handle read-heavy analytical workloads, employing denormalized schemas like star and snowflake designs to reduce join complexity and accelerate queries (Inmon, 2005; Kimball & Ross, 2013). This separation allows data warehouses to support enterprise-wide analytics while OLTP systems maintain operational efficiency.

Modern BI platforms continue to evolve, incorporating techniques like columnar storage and Massively Parallel Processing (MPP) architectures in cloud-based data warehouses such as BigQuery, Snowflake, and Azure Synapse. While foundational principles remain relevant, these platforms simplify aspects like indexing and partitioning. Nonetheless, core concepts of dimensional modeling, cube structures, and SQL aggregations persist as guiding practices for robust, scalable BI solutions.

In conclusion, dimensional modeling, data cube design, SQL aggregation strategies, and the differentiation between OLAP and OLTP form a comprehensive framework for data warehousing and business intelligence. Adhering to these principles enables efficient data exploration and transforms transactional data into actionable insights, facilitating informed and strategic decision-making.

## 4. Data Collection and Preparation

The data collection process for this project involves sourcing, integrating, and refining datasets

that capture customer attributes, product details, sales transactions, and temporal information. Data was sourced from:

Below is an overview of the tables in our dataset,

**DIM\_Customer:**

Contains customer-related information such as CustomerKey (unique identifier), demographics (gender, marital status), location, and purchase behavior (income, number of cars, children). This table helps to link customers to sales data for customer segmentation and analysis.

**DIM\_Date:**

Stores date-related information, including DateKey (unique identifier), specific day, month, and year details, along with fiscal periods and calendar breakdowns. Useful for time-based analysis such as month-on-month growth, sales trends, and time-series analysis.

**DIM\_Product:**

Provides product-related details such as ProductKey (unique identifier), name, pricing, inventory data (safety stock, reorder points), and specifications like size, weight, and color. It helps in product-level analysis such as sales, profitability, and stock levels.

**FACT\_Product:**

Tracks product movements and quantities over time, including attributes like ProductKey, DateKey, and transaction details (e.g., units in, units out). This table supports analysis related to inventory flow and product performance over time.

**FACT\_Sales\_Internet:**

Contains sales transaction data from internet-based orders. It includes ProductKey, CustomerKey, sales amounts, order quantities, discounts, and tax information. This table is essential for sales and customer analytics such as total revenue, profit margins, and customer spending patterns.

**OLTP\_Product:**

Provides detailed product information such as ProductID, name, product number, cost, price, inventory attributes (safety stock, reorder point), and product specifications (size, color, weight). This table is used for real-time product management and helps to monitor product performance and inventory health.

**OLTP\_Sales:**

Contains transactional sales data, including SalesOrderID, SalesOrderDetailID, OrderQty, ProductID, SpecialOfferID, unit price, and discounts. This table is used for real-time sales reporting, tracking sales per order, and monitoring order quantities for each product.

## **Migrating Data to Google Cloud Platform**

We have used GCP as a **Cloud Data Warehouse Solution** for storage, processing and performing analytics. Below are the steps we followed,

### **1. Sign In to GCP Console**

- Open the Google Cloud ([console.cloud.google.com](https://console.cloud.google.com)).
- Sign in with the Google account and select or create a project.

### **2. Access BigQuery Studio**

- In the Google Cloud Console, use the navigation menu on the left to find BigQuery under the "Big Data" section.
- Click on BigQuery to open the BigQuery Console.

### **4. Create a Dataset (available under Project ID)**

- In the BigQuery console, click the Project Name on the left panel.
- Click Create Dataset at the top.
- Enter a name for the dataset and configure settings as needed.
- Click Create Dataset.

The screenshot displays the Google Cloud BigQuery Studio interface. On the left, the 'Explorer' pane shows a project named 'able-armor-443800-m5' with various resources like Queries, Notebooks, Data canvases, Workflows, and External connections. The main editor shows a SQL query titled 'Untitled query' with the following code:

```

1 --Product Sales, Profit Margin, and Inventory Health Report
2
3 WITH ProductSales AS (
4     SELECT
5         p.ProductID,
6         p.Name AS ProductName,
7         SUM(s.OrderQty) AS TotalQuantitySold,
8         SUM(s.UnitPrice * s.OrderQty) AS TotalSalesRevenue,
9         SUM(s.OrderQty * (s.UnitPrice - p.StandardCost)) AS TotalProfit
10    FROM dw_project.OLTP_Sales s
11   JOIN dw_project.OLTP_Product p
12     ON s.ProductID = p.ProductID
13   GROUP BY p.ProductID, p.Name
14 ),
15 InventoryStatus AS (
16     SELECT
17         p.ProductID,
18         p.SafetyStockLevel,
19         p.ReorderPoint,

```

Below the query editor, the 'Query results' section shows a table with 3 rows and 5 columns: Row, quantitySold, TotalSalesRevenue, TotalProfit, and ProfitMargin. The data is as follows:

Row	quantitySold	TotalSalesRevenue	TotalProfit	ProfitMargin
1	2977	4406151.266199...	679002.9361000...	15.41034102275...
2	2664	4014067.799899...	678789.6167000...	16.91026785140...
3	2394	3696486.472599...	666593.3895999...	18.03316188334...

On the right, the 'Create dataset' dialog is open. It shows the 'Project ID' as 'able-armor-443800-m5' and a 'Dataset ID' field. The 'Location type' is set to 'Multi-region' with a dropdown showing 'US (multiple regions in United States)'. There are checkboxes for 'Link to an external dataset', 'Default table expiration', and 'Enable table expiration'. At the bottom, there are 'CREATE DATASET' and 'CANCEL' buttons.

## 6. Upload Your Data

As for this project, our data was small – we preferred a bit easier approach by directly uploading the CSV file into the BigQuery Dataset. For Big Data, the proper approach would be to use Google Cloud Storage (Storage Bucket) followed by ETL operations. Eventually, the final step is importing the data into BigQuery.

- In BigQuery Studio, click on the + Create button.
- Select Upload Data from the options.
- Choose the file you want to upload.
- Select the dataset you want to upload the file to or create a new one.
- Specify the file format (CSV, JSON, etc.), delimiter (for CSV), and other options (e.g., schema, headers).

**Create table**

**Source**

Create table from  
Empty table

**Destination**

Project \*  
able-armor-443800-m5 [BROWSE](#)

Dataset \*  
dw\_project

Table \*  
Maximum name size is 1,024 UTF-8 bytes. Unicode letters, marks, numbers, connectors, dashes, and spaces are allowed.

Table type  
Native table

☐ Create a BigQuery table for Apache Iceberg [PREVIEW](#)

**Schema**

☐ Edit as text

[+ ADD FIELD](#)

**Partition and cluster settings**

Partitioning  
No partitioning

[CREATE TABLE](#) [CANCEL](#)

### Tools and Platforms:

Data may be collected and stored within a cloud-based environment, such as Google BigQuery, Snowflake, or Azure Synapse, or housed on-premises in a data warehouse appliance. Extract, Transform, Load (ETL) or Extract, Load, Transform (ELT) processes can be orchestrated using tools like Apache Airflow, dbt, or native cloud services.

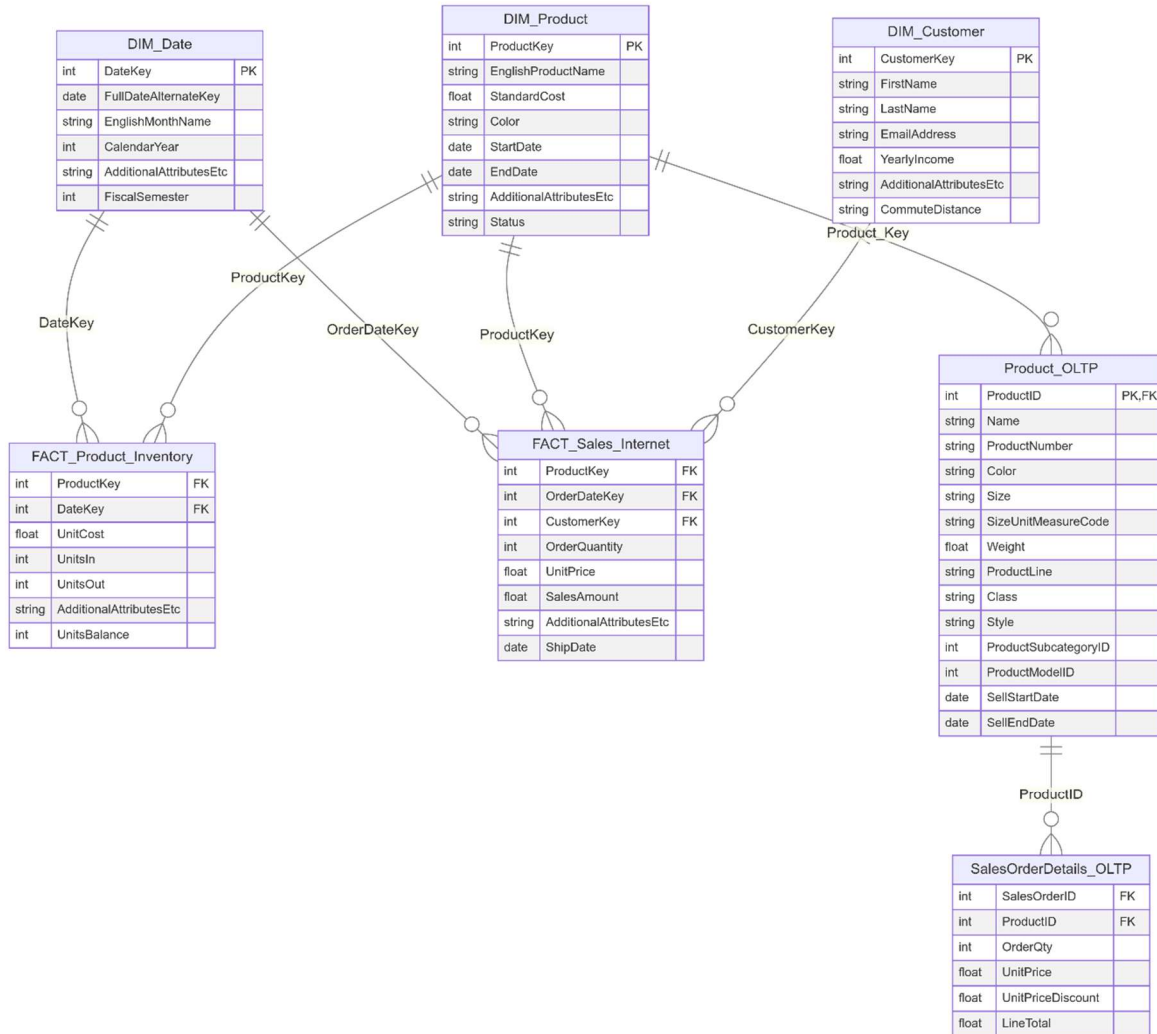
### Quality Assurance and Governance:

Throughout the data collection and preparation phases, quality checks (row counts, referential integrity checks, duplication tests) and metadata documentation ensure the reliability of the final data warehouse environment. This governance framework helps maintain trust in the analytics produced downstream.

By meticulously collecting and preparing the data, the project sets the stage for robust dimensional modeling, efficient OLAP-style queries, and the advanced analytics required to address the research questions—ranging from commute distance correlations and monthly growth trends to demographic-based customer segmentation and product profitability.

## 5. Database Design

# Dimensional Modelling using ERD



- **Fact Tables:**

- **Sales\_Table:**

- Stores measurable business metrics (e.g., SalesAmount, TaxAmt, OrderQuantity) with keys to link to dimensions like ProductKey, CustomerKey, and DateKey.

- **Product\_Inventory\_Table:**

- Captures metrics related to inventory or product movement (e.g., UnitsIn, UnitsOut, UnitCost).



- **Dimension Tables:**

- DIM\_Product:
  - Includes descriptive attributes like ProductLine, ModelName, and multilingual descriptions, enhancing slicing and dicing capabilities.
- DIM\_Date:
  - Provides time-based granularity, enabling analysis across multiple time hierarchies such as days, weeks, quarters, and fiscal periods.
- DIM\_Customer:
  - Offers customer-level details like YearlyIncome, Gender, and GeographyKey to analyze sales by demographics or location.

- **Hierarchies and Granularity:**

- Dimensions like DIM\_Date have natural hierarchies (e.g., Day → Month → Quarter → Year).
- DIM\_Product includes hierarchies like ProductSubcategory → ProductLine.
- Fact tables provide transaction-level granularity.

### **Appropriating Data Cubes Design through our ERD.**

- Data cubes are often used in OLAP for summarizing data in multiple dimensions. The cube for our data include measures like total sales or quantities, broken down by dimensions like products, time, and customer.
- To design a data cube, we have to define measures (numeric values we want to aggregate) and dimensions (categories for summarizing data).
- These cubes allow slicing (e.g., filtering by product category), dicing (e.g., comparing sales across regions), and rolling up (e.g., aggregating by year)

# Creating Data Cubes in Google Big Query Studio:

## 1. Sales Cube

Goal: Analyze SalesAmount by ProductLine, Year, and GeographyKey.

The screenshot shows the Google BigQuery Studio interface. On the left is a navigation sidebar with options like Orchestration, Migration, Administration, Monitoring, Jobs explorer, Capacity management, BI Engine, Recommendations, Disaster recovery, Policy tags, and Release notes. The main area displays a query titled 'Untitled query' with the following SQL code:

```
1 SELECT
2   p.ProductLine,
3   d.CalendarYear,
4   c.GeographyKey,
5   SUM(f.SalesAmount) AS TotalSalesAmount,
6   SUM(f.OrderQuantity) AS TotalOrderQuantity,
7   AVG(f.UnitPrice) AS AvgUnitPrice
8 FROM
9   dw_project.FACT_Sales_Internet f
10  JOIN
```

Below the query editor, the 'Query results' section is active, showing a table with 7 columns: Row, ProductLine, CalendarYear, GeographyKey, TotalSalesAmount, TotalOrderQuantity, and AvgUnitPrice. The table contains 5 rows of data.

Row	ProductLine	CalendarYear	GeographyKey	TotalSalesAmount	TotalOrderQuantity	AvgUnitPrice
1	R	2011	301	85878.48	24	3578.27
2	R	2011	39	82301.3228	27	3048.197140740...
3	R	2011	4	81518.33279999...	26	3135.320492307...
4	M	2013	298	79260.32000000...	158	501.6475949367...
5	R	2013	4	72666.21000000...	91	798.5297802197...

### Explanation:

- **Dimensions:** ProductLine, CalendarYear, GeographyKey.
- **Measures:** SUM(SalesAmount), SUM(OrderQuantity), AVG(UnitPrice).
- **Analysis:** This query generates a cube showing total sales and order quantities for each product line in each geography for every year.

## 2. Customer Insights Cube

**Goal:** Analyze SalesAmount by YearlyIncome, Gender, and Year.

The screenshot displays the Google Cloud BigQuery console. On the left, a sidebar lists various BigQuery features like Migration, Administration, Monitoring, and BI Engine. The main area shows a query editor with a SQL query and a results table.

**Query:**

```
1 SELECT
2   c.YearlyIncome,
3   c.Gender,
4   d.CalendarYear,
5   SUM(f.SalesAmount) AS TotalSalesAmount,
6   COUNT(DISTINCT f.CustomerKey) AS TotalCustomers
7 FROM
8   `dw_project.FACT_Sales_Internet` f
9 JOIN
10  `dw_project.DIM_Customer` c
11  ON f.CustomerKey = c.CustomerKey
12 JOIN
```

**Query results:**

Row	YearlyIncome	Gender	CalendarYear	TotalSalesAmount	TotalCustomers
1	70000	F	2013	1444668.799999...	1177
2	70000	M	2013	1312756.839999...	1123
3	40000	F	2013	1226811.029999...	1272
4	40000	M	2013	1192899.179999...	1261

### Explanation:

- **Dimensions:** YearlyIncome, Gender, CalendarYear.
- **Measures:** SUM(SalesAmount), COUNT(DISTINCT CustomerKey).
- **Analysis:** This query builds a cube for analyzing spending trends by income bracket and gender over time.

## 6. Exploratory Data Analysis

## 7. Reporting, Modeling and Storytelling

- Queries:

OLAP systems are primarily used for data analysis and business intelligence purposes. They allow users to query large datasets quickly and perform complex aggregations, calculations, and trend analyses. Here are some common use cases of OLAP:

- Business Reporting and Dashboards
- Data Warehousing and Trend Analysis
- Customer Segmentation
- Financial Planning and Forecasting

Our Project mainly consists of an OLAP Data Mart with 2 FACTs and 3 Dims. Below are some analytical SQL queries using this DataMart,

### 1. Relation Between Commute Distance and Sales

The screenshot shows the Google Cloud BigQuery console interface. On the left is the Explorer pane with a tree view of resources including 'dw\_project' and its dimensions and facts. The main area displays an 'Untitled query' with the following SQL code:

```
1 --Relation b/w commute distance
2
3 SELECT
4   c.CommuteDistance,
5   AVG(s.SalesAmount) AS AvgSalesAmount,
6   SUM(s.OrderQuantity) AS TotalItemsPurchased,
7   COUNT(DISTINCT s.CustomerKey) AS UniqueCustomers
8 FROM dw_project.FACT_Sales_Internet s
9 JOIN dw_project.DIM_Customer c ON s.CustomerKey = c.CustomerKey
10 GROUP BY c.CommuteDistance
11 ORDER BY AvgSalesAmount DESC;
```

Below the query is the 'Query results' section, which includes a table with 5 rows of data. The table has columns for CommuteDistance, AvgSalesAmount, TotalItemsPurchased, and UniqueCustomers. The results are as follows:

Row	CommuteDistance	AvgSalesAmount	TotalItemsPurchased	UniqueCustomers
1	0-1 Miles	526.0051731731...	21307	6310
2	2-5 Miles	492.4151563070...	10084	3234
3	5-10 Miles	460.9654301742...	10615	3214
4	10+ Miles	456.1924554366...	8222	2494
5	1-2 Miles	446.5691396066...	10170	3232

At the bottom of the results section, there is a 'Job history' tab and a 'Results per page' dropdown set to 50.

Purpose:

- This query examines how the commute distance of customers relates to their purchasing behavior.

- Metrics like average sales amount, total items purchased, and unique customers are calculated for each commute distance category.

#### Use Cases:

- Business Targeting: Identify customer groups based on commute distance and create targeted marketing campaigns.
- Resource Allocation: Align resources and product availability for regions with high-spending commuters.

#### Insights:

- Uncovers if customers with longer or shorter commutes spend more on average.
- Determines if certain commute groups purchase more items, which could suggest lifestyle patterns.

## 2. Month-on-Month (MoM) Growth in Sales

The screenshot displays the Google Cloud BigQuery console interface. The top navigation bar shows the Google Cloud logo, the project name 'My First Project', and a search bar. Below the navigation bar, a message indicates a free trial status with \$300.00 credit and 89 days remaining. The main content area is divided into three sections: Explorer, Query Editor, and Query Results.

**Explorer:** The left sidebar shows the project structure, including 'Queries', 'Notebooks', 'Data canvases', 'Data preparations', 'Workflows', 'External connections', and 'dw\_project'. The 'dw\_project' folder is expanded, showing subfolders like 'DIM\_Customer', 'DIM\_Date', 'DIM\_Product', 'FACT\_Product', 'FACT\_Sales\_Internet', 'OLTP\_Product', and 'OLTP\_Sales'.

**Query Editor:** The central pane shows a SQL query titled 'Untitled query'. The query is as follows:

```
--Month on month growth
WITH MonthlySales AS (
  SELECT
    d.CalendarYear,
    d.EnglishMonthName,
    d.MonthNumberOfYear,
    SUM(s.SalesAmount) AS MonthlyRevenue
  FROM dw_project.FACT_Sales_Internet s
  JOIN dw_project.DIM_Date d ON s.OrderDateKey = d.DateKey
  GROUP BY d.CalendarYear, d.EnglishMonthName, d.MonthNumberOfYear
),
RevenueWithLag AS (
  SELECT
    CalendarYear,
    EnglishMonthName,
    MonthNumberOfYear, -- Include MonthNumberOfYear here
    LAG(MonthlyRevenue, 1) OVER (
      PARTITION BY CalendarYear, EnglishMonthName
      ORDER BY MonthNumberOfYear
    ) AS PrevMonthRevenue
  FROM MonthlySales
)
SELECT
  CalendarYear,
  EnglishMonthName,
  MonthNumberOfYear,
  MonthlyRevenue,
  PrevMonthRevenue,
  (MonthlyRevenue - PrevMonthRevenue) / PrevMonthRevenue * 100 AS MoMGrowth
FROM RevenueWithLag
```

**Query Results:** The bottom section displays the results of the query. The table has 6 columns: Row, CalendarYear, EnglishMonthName, MonthlyRevenue, PrevMonthRevenue, and MoMGrowth. The results show data for the years 2011, 2010, and 2009, with months from January to May.

Row	CalendarYear	EnglishMonthName	MonthlyRevenue	PrevMonthRevenue	MoMGrowth
1	2011	January	469823.9148000...	43421.03639999...	982.0191173511...
2	2011	February	466334.9030000...	469823.9148000...	-0.74262115871...
3	2011	March	485198.6594000...	466334.9030000...	4.045109272037...
4	2011	April	502073.8458000...	485198.6594000...	3.477995265046...
5	2011	May	561681.4758000...	502073.8458000...	11.87228342974...

The bottom of the console shows a 'Job history' section with a 'SUMMARY' tab and a 'Job history' tab. The 'Job history' tab is currently selected, showing a list of jobs with their status and execution details.

### Full Query:

```
1  --Month on month growth
2
3  WITH MonthlySales AS (
4      SELECT
5          d.CalendarYear,
6          d.EnglishMonthName,
7          d.MonthNumberOfYear,
8          SUM(s.SalesAmount) AS MonthlyRevenue
9      FROM dw_project.FACT_Sales_Internet s
10     JOIN dw_project.DIM_Date d ON s.OrderDateKey = d.DateKey
11     GROUP BY d.CalendarYear, d.EnglishMonthName, d.MonthNumberOfYear
12 ),
13 RevenueWithLag AS (
14     SELECT
15         CalendarYear,
16         EnglishMonthName,
17         MonthNumberOfYear, -- Include MonthNumberOfYear here
18         MonthlyRevenue,
19         LAG(MonthlyRevenue) OVER (ORDER BY CalendarYear, MonthNumberOfYear) AS PrevMonthRevenue
20     FROM MonthlySales
21 )
22 SELECT
23     CalendarYear,
24     EnglishMonthName,
25     MonthlyRevenue,
26     PrevMonthRevenue,
27     ((MonthlyRevenue - PrevMonthRevenue) / PrevMonthRevenue) * 100 AS MoMGrowth
28 FROM RevenueWithLag
29 WHERE PrevMonthRevenue IS NOT NULL
30 ORDER BY CalendarYear, MonthNumberOfYear; -- Ensure proper ordering
31
```

### Purpose:

- Tracks sales growth or decline by comparing the monthly revenue to the previous month.
- Calculates the percentage growth or drop (MoM Growth).

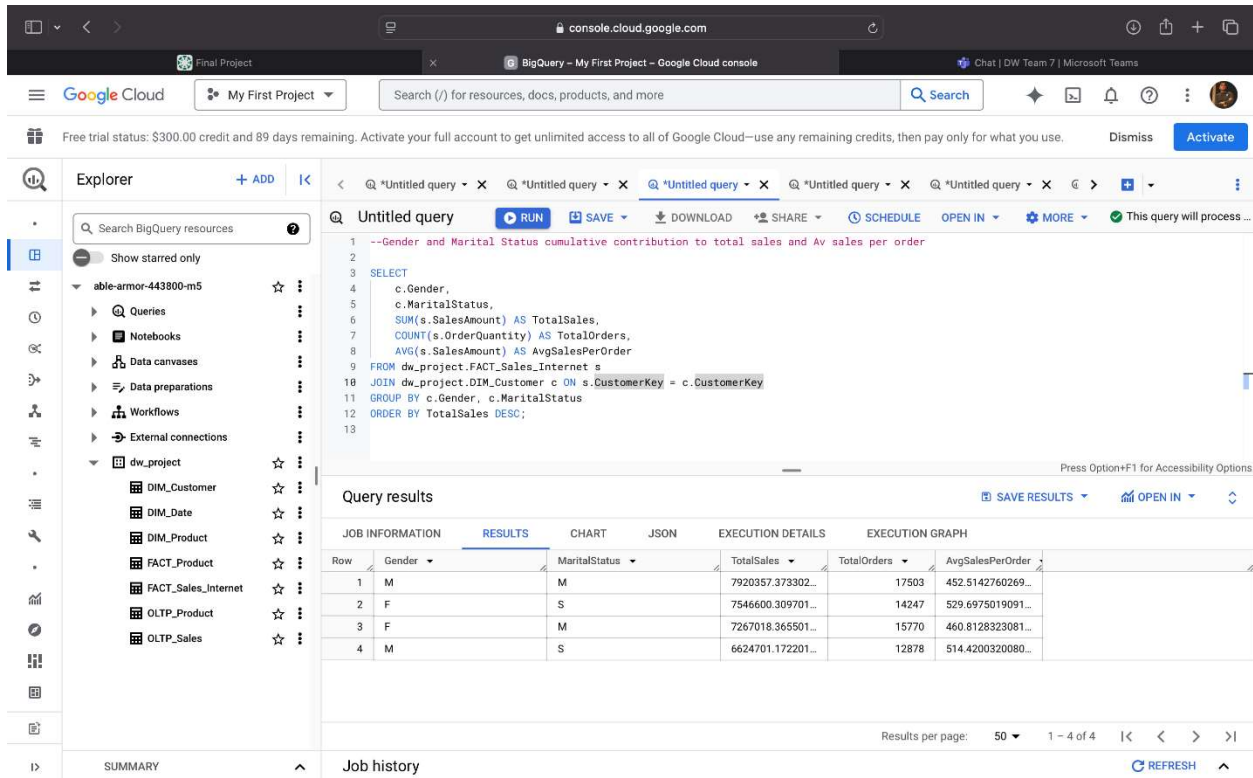
### Use Cases:

- Trend Analysis: Evaluate business performance over time and identify seasonal patterns.
- Strategy Adjustment: Pinpoint months with poor performance and investigate causes, such as marketing efforts, product availability, or customer demand.

### Insights:

- Shows months with significant growth or decline in revenue.
- Identifies trends to forecast future performance and plan marketing campaigns or sales initiatives.

### 3. Gender and Marital Status Contribution



The screenshot shows the Google Cloud BigQuery console interface. On the left is the Explorer pane with a tree view of resources including 'dw\_project' and its tables. The main area displays a SQL query titled 'Untitled query' and its results. The query calculates cumulative sales and average sales per order, grouped by gender and marital status. The results table shows four rows of data.

**Query:**

```

1  --Gender and Marital Status cumulative contribution to total sales and Av sales per order
2
3  SELECT
4    c.Gender,
5    c.MaritalStatus,
6    SUM(s.SalesAmount) AS TotalSales,
7    COUNT(s.OrderQuantity) AS TotalOrders,
8    AVG(s.SalesAmount) AS AvgSalesPerOrder
9  FROM dw_project.FACT_Sales_Internet s
10 JOIN dw_project.DIM_Customer c ON s.CustomerKey = c.CustomerKey
11 GROUP BY c.Gender, c.MaritalStatus
12 ORDER BY TotalSales DESC;
13

```

**Query results:**

Row	Gender	MaritalStatus	TotalSales	TotalOrders	AvgSalesPerOrder
1	M	M	7920357.373302...	17503	452.5142760269...
2	F	S	7546600.309701...	14247	529.6975019091...
3	F	M	7267018.365501...	15770	460.8128323081...
4	M	S	6624701.172201...	12878	514.4200320080...

#### Purpose:

- Breaks down sales data by customer gender and marital status to understand their contributions to overall sales.

#### Use Cases:

- Demographic Insights: Help businesses tailor products, offers, and communication strategies for specific demographic groups.
- Product Recommendations: Identify differences in spending habits between groups to recommend products.

#### Insights:

- Indicates which demographic groups generate the most revenue.
- Reveals purchasing patterns, such as whether single or married customers spend more.

## 4. Customer Lifetime Value Analysis

The screenshot displays the Google Cloud BigQuery console interface. The top navigation bar shows the project name 'My First Project' and a search bar. Below the navigation bar, the 'Explorer' panel on the left lists various BigQuery resources, including 'Queries', 'Notebooks', 'Data canvases', 'Data preparations', 'Workflows', 'External connections', and a folder named 'dw\_project'. The main panel shows an 'Untitled query' with a SQL query for Customer Lifetime Value Analysis. The query calculates total revenue, total orders, and customer tenure days for the top customers. The query results are displayed in a table with columns: Row, CustomerKey, FirstName, LastName, TotalRevenue, TotalOrders, and CustomerTenureDays. The results show the top 7 customers, with the highest total revenue being 13295.38 for customer 12301 (Nichole Nara).

```
1 --Customer Lifetime Value Analysis - Can target top customers
2
3 SELECT
4   c.CustomerKey,
5   c.FirstName,
6   c.LastName,
7   SUM(s.SalesAmount) AS TotalRevenue,
8   COUNT(s.SalesOrderNumber) AS TotalOrders,
9   MAX(s.OrderDateKey) - MIN(s.OrderDateKey) AS CustomerTenureDays
10 FROM dw_project.FACT_Sales_Internet s
11 JOIN dw_project.DIM_Customer c ON s.CustomerKey = c.CustomerKey
12 GROUP BY c.CustomerKey, c.FirstName, c.LastName
13 ORDER BY TotalRevenue DESC;
14
```

Row	CustomerKey	FirstName	LastName	TotalRevenue	TotalOrders	CustomerTenureDays
1	12301	Nichole	Nara	13295.38	13	20610
2	12132	Kaitlyn	Henderson	13294.269999999...	14	20887
3	12308	Margaret	He	13269.269999999...	14	20516
4	12131	Randall	Dominguez	13265.99	11	20789
5	12300	Adriana	Gonzalez	13242.699999999...	10	20493
6	12321	Rosa	Hu	13215.649999999...	15	20519
7	12124	Brandi	Gill	13195.639999999...	12	20878

Purpose:

- Calculates customer lifetime value (CLV) by aggregating total spending, order count, and tenure (days active).

Use Cases:

- Customer Retention: Focus on high-value customers by offering personalized rewards or discounts.
- Profitability Analysis: Identify which customers contribute the most to revenue.

Insights:

- Highlights the most loyal and high-revenue customers.
- Provides data to improve customer retention strategies and increase repeat purchases.



## 5. Product Profitability

The screenshot displays the Google Cloud BigQuery console interface. On the left, the Explorer pane shows a project named 'dw\_project' with various tables like 'DIM\_Customer', 'DIM\_Date', 'DIM\_Product', 'FACT\_Product', 'FACT\_Sales\_Internet', 'OLTP\_Product', and 'OLTP\_Sales'. The main editor shows a SQL query titled 'Product Profitability' that calculates revenue, total cost, profit, and profit margin for different products. The query results are displayed in a table with columns: Row, EnglishProductName, Revenue, TotalCost, Profit, and ProfitMargin. The results show data for products like 'Mountain-200 Black, 46', 'Mountain-200 Black, 42', 'Mountain-200 Silver, 38', etc.

```
--Product Profitability
SELECT
  p.EnglishProductName,
  SUM(s.SalesAmount) AS Revenue,
  SUM(s.TotalProductCost) AS TotalCost,
  SUM(s.SalesAmount - s.TotalProductCost) AS Profit,
  (SUM(s.SalesAmount - s.TotalProductCost) / SUM(s.TotalProductCost)) * 100 AS ProfitMargin
FROM dw_project.FACT_Sales_Internet s
JOIN dw_project.DIM_Product p
  ON SAFE_CAST(p.ProductKey AS INT64) = s.ProductKey -- Use SAFE_CAST to prevent errors
WHERE REGEXP_CONTAINS(p.ProductKey, r'\d+') -- Only allow rows where ProductKey is numeric
GROUP BY p.EnglishProductName
ORDER BY Profit DESC;
```

Row	EnglishProductName	Revenue	TotalCost	Profit	ProfitMargin
1	Mountain-200 Black, 46	1373469.548199...	746847.9746999...	626621.5735000...	83.90215876955...
2	Mountain-200 Black, 42	1363142.093399...	741382.4850999...	621759.6083000...	83.86489036305...
3	Mountain-200 Silver, 38	1339462.790399...	728598.3555999...	610864.4348000...	83.84103945677...
4	Mountain-200 Silver, 46	1301100.098399...	707609.6255999...	593490.4728000...	83.87258331834...
5	Mountain-200 Black, 38	1294866.141199...	704388.6807999...	590477.4604000...	83.82835733949...
6	Mountain-200 Silver, 42	1257434.572799...	683922.6351999...	573511.9376000...	83.85625918526...
7	Road-150 Red, 48	1205876.990000...	731726.1454000...	474150.8446000...	64.79894801911...

Purpose:

- Evaluates product profitability by comparing revenue to costs and calculating profit margins.

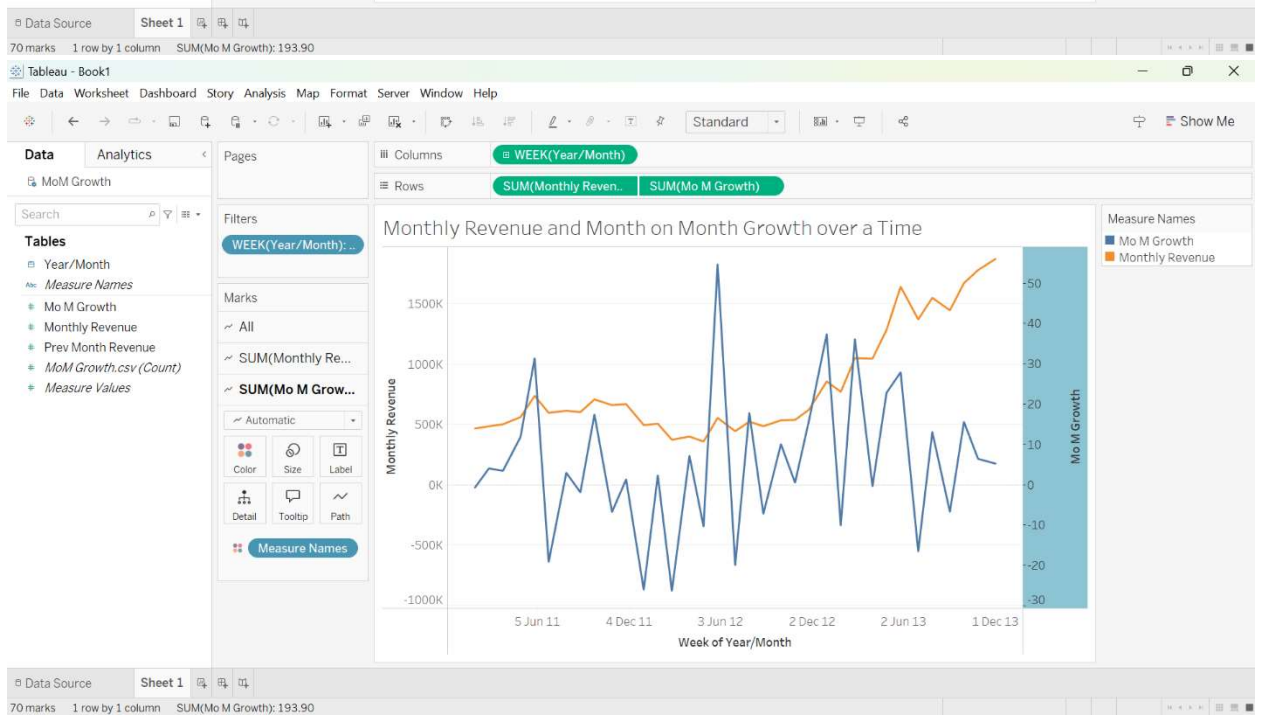
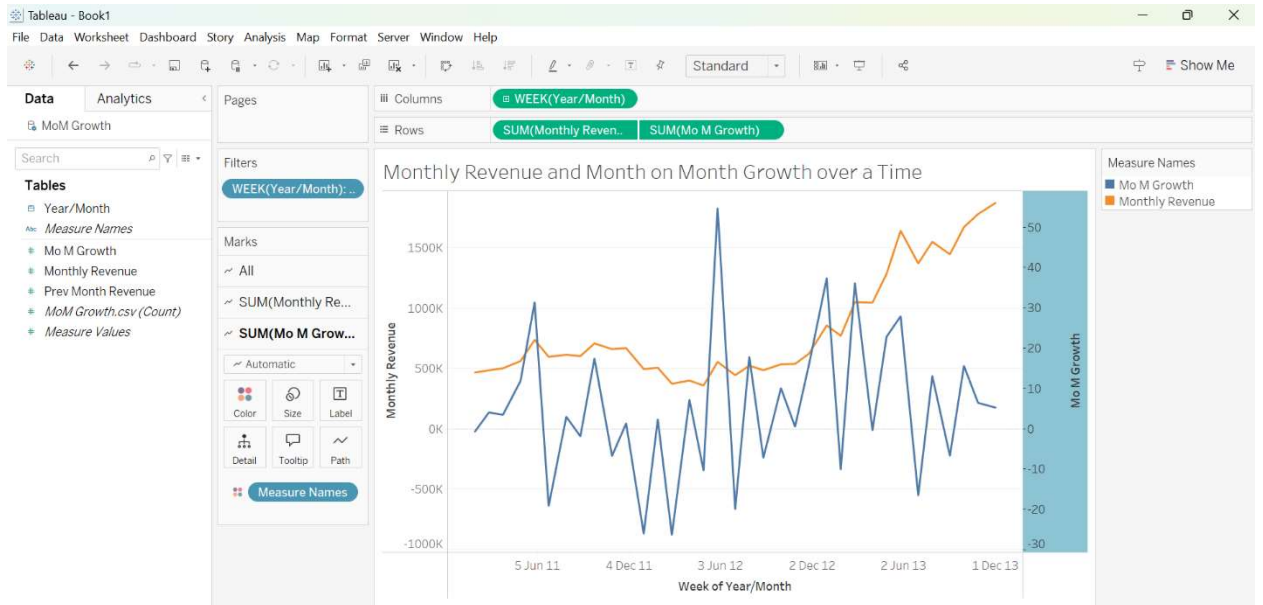
Use Cases:

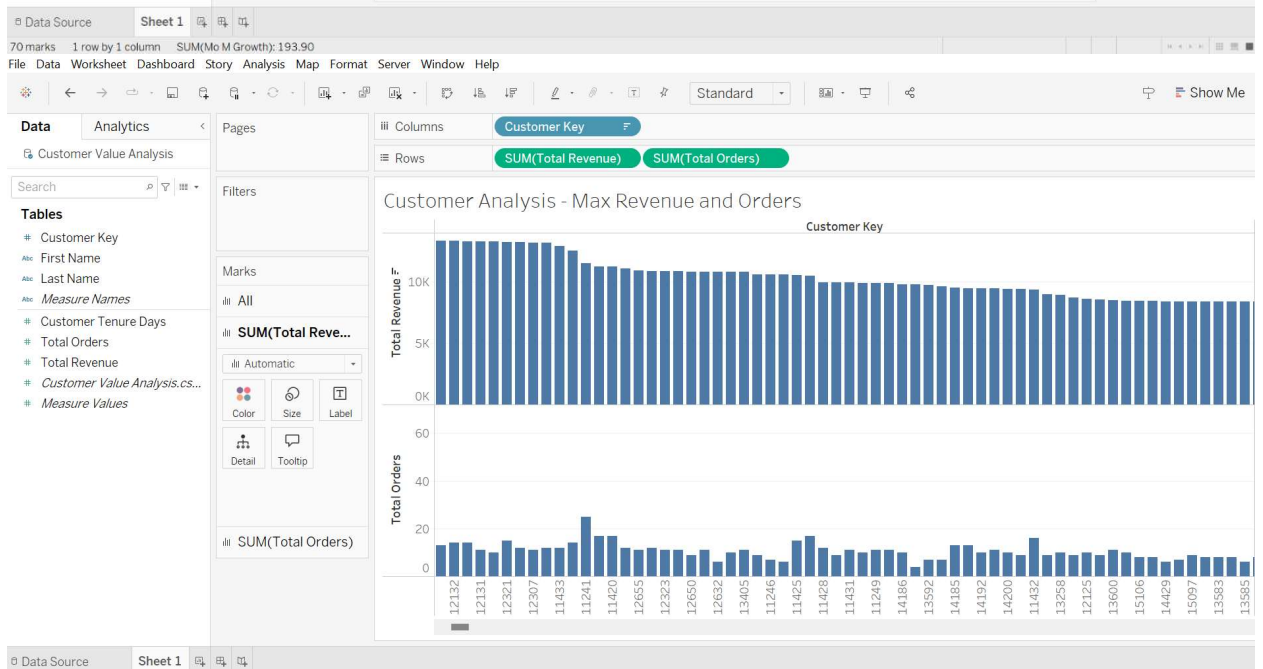
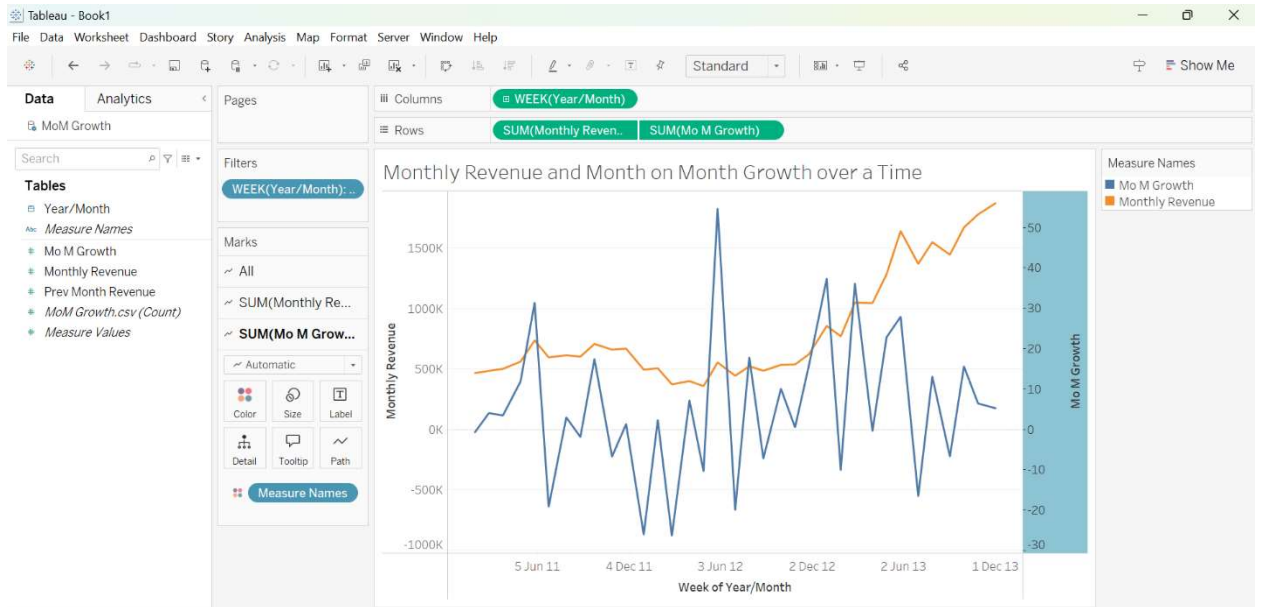
- Inventory Decisions: Identify top-performing products to prioritize in inventory and promotions.
- Pricing Strategies: Adjust prices for low-margin products to improve profitability.

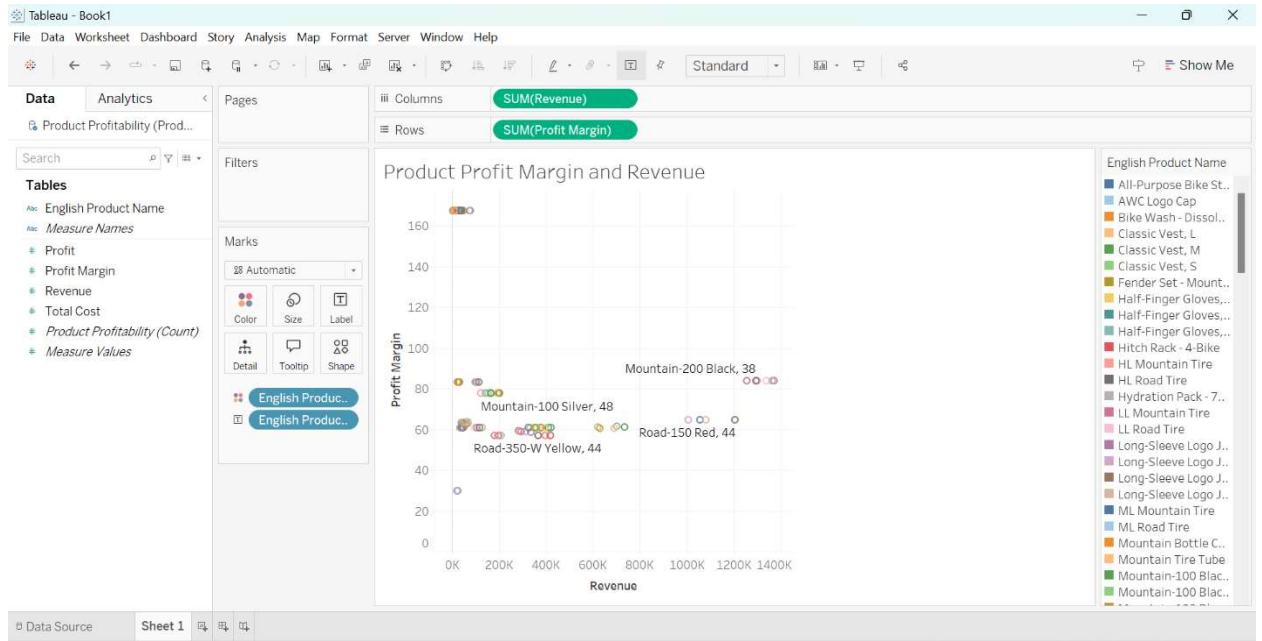
Insights:

- Highlights the most profitable products.
- Identifies products with low profitability for further investigation.

## 6. Customer Segmentation by Spending Category







## 8. Conclusion