**BIKE STORE DATA ANALYSIS CHALLENGE**
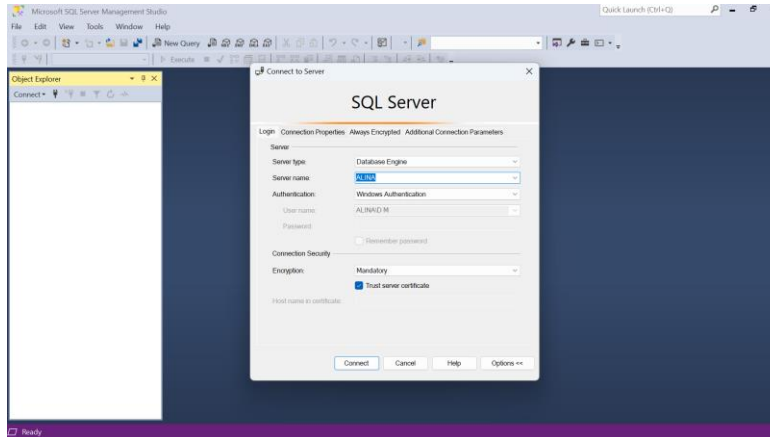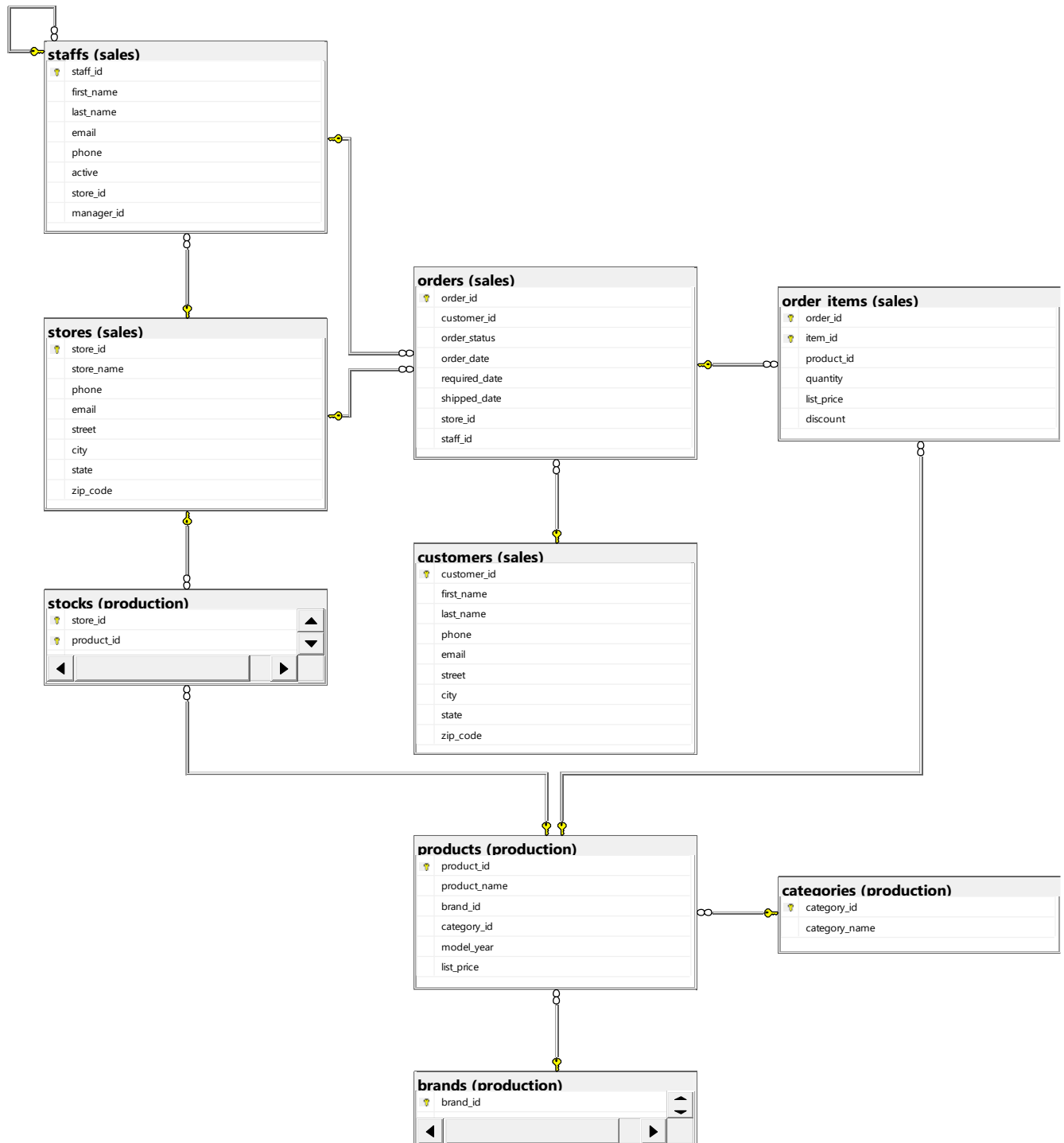
## 1. SETUP PHASE

Connect to SSMS



From the above BikeStores, we can understand the relationship between the tables in BikeStores database.

| Table Name | Description | Primary Key | Foreign Keys |
|---|---|---|---|
| Staffs (sales) | Information about staff members. | staff_id | store_id, manager_id |
| Stores (sales) | Details of stores in the bike shop. | store_id | None |
| Customers (sales) | Details of customers placing orders. | customer_id | None |
| Orders (sales) | Information about customer orders. | order_id | customer_id, store_id, staff_id |
| Order Items (sales) | Items included in each order. | order_id, item_id | order_id, product_id |
| Products (production) | Details of products available for sale. | product_id | brand_id, category_id |
| Stocks (production) | Inventory data for each store. | store_id, product_id | store_id, product_id |

Big Data and Data Mining

# BikeStores ERD diagram

**staffs (sales)**
- staff_id
- first_name
- last_name
- email
- phone
- active
- store_id
- manager_id

**stores (sales)**
- store_id
- store_name
- phone
- email
- street
- city
- state
- zip_code

**stocks (production)**
- store_id
- product_id

**orders (sales)**
- order_id
- customer_id
- order_status
- order_date
- required_date
- shipped_date
- store_id
- staff_id

**customers (sales)**
- customer_id
- first_name
- last_name
- phone
- email
- street
- city
- state
- zip_code

**order_items (sales)**
- order_id
- item_id
- product_id
- quantity
- list_price
- discount

**products (production)**
- product_id
- product_name
- brand_id
- category_id
- model_year
- list_price

**categories (production)**
- category_id
- category_name

**brands (production)**
- brand_id

Dipawoli Malla

Big Data and Data Mining

**Relationships**

We can also interpret relationships between the tables in ERD.

Production Relationships

1. categories and products

    o  Relationship: One-to-Many

    o  Explanation: A category can include multiple products, and a product belongs to a single category.

2. brands and products

    o  Relationship: One-to-Many

    o  Explanation: A brand can have multiple products, and each product belongs to one brand.

3. products and stocks

    o  Relationship: One-to-Many

    o  Explanation: A product can be stocked in multiple stores, and each stock record relates to a single product.

Sales Relationships

1. customers and orders

    o  Relationship: One-to-Many

    o  Explanation: A customer can place multiple orders, and each order is placed by a single customer.

2. orders and order_items

    o  Relationship: One-to-Many

    o  Explanation: An order can contain multiple order items, and each order item belongs to a single order.

3. order_items and products

    o  Relationship: Many-to-One

    o  Explanation: Multiple order items can refer to a single product, and each order item corresponds to one product.

4. staffs and orders

Dipawoli Malla

- o Relationship: One-to-Many
- o Explanation: A staff member can handle multiple orders, and each order is processed by a single staff member.

5. stores and stocks

- o Relationship: One-to-Many
- o Explanation: A store can have stocks of multiple products, and each stock record is for a specific product in a store.

6. stores and orders

- o Relationship: One-to-Many
- o Explanation: A store can process multiple orders, and each order is placed at a specific store.

## 2. Exploration

How complete is my dataset?

To assess the completeness of the dataset, we need to check for:

- Missing records
- Consistency
- Outliers

To check the missing records:

```sql
use BikeStores;

--Exploration of Data to check its completeness

--missing values in the data

SELECT 'production.brands' AS TableName, COUNT(*) AS MissingValues
FROM production.brands
WHERE brand_id IS NULL
    OR brand_name IS NULL;
```

# Big Data and Data Mining

| | TableName | MissingValues |
|---|---|---|
| 1 | production.brands | 0 |

| | TableName | MissingValues |
|---|---|---|
| 1 | production.categories | 0 |

| | TableName | MissingValues |
|---|---|---|
| 1 | production.products | 0 |

| | TableName | MissingValues |
|---|---|---|
| 1 | sales.customers | 1267 |

| | TableName | MissingValues |
|---|---|---|
| 1 | sales.order_items | 0 |

| | TableName | MissingValues |
|---|---|---|
| 1 | sales.orders | 170 |

| | TableName | MissingValues |
|---|---|---|
| 1 | sales.staffs | 1 |

| | TableName | MissingValues |
|---|---|---|
| 1 | sales.stores | 0 |

✅ Query executed successfully.

Ln 75    Col 4    Ch 4    IN

To check the date range of our data:

```
-- Date range of orders
SELECT
    MIN(order_date) AS EarliestOrderDate,
    MAX(order_date) AS LatestOrderDate
FROM sales.orders;

-- Date range of product availability in stocks
SELECT
    MIN(required_date) AS EarliestRequiredDate,
    MAX(required_date) AS LatestRequiredDate
FROM sales.orders;
```

| | order_id | customer_id | order_status | order_date | required_date | shipped_date | store_id | staff_id |
|---|---|---|---|---|---|---|---|---|
| 1 | 1 | 259 | 4 | 2016-01-01 | 2016-01-03 | 2016-01-03 | 1 | 2 |
| 2 | 2 | 1212 | 4 | 2016-01-01 | 2016-01-04 | 2016-01-03 | 2 | 6 |
| 3 | 3 | 523 | 4 | 2016-01-02 | 2016-01-05 | 2016-01-03 | 2 | 7 |
| 4 | 4 | 175 | 4 | 2016-01-03 | 2016-01-04 | 2016-01-05 | 1 | 3 |
| 5 | 5 | 1324 | 4 | 2016-01-03 | 2016-01-06 | 2016-01-06 | 2 | 6 |
| 6 | 6 | 94 | 4 | 2016-01-04 | 2016-01-07 | 2016-01-05 | 2 | 6 |
| 7 | 7 | 324 | 4 | 2016-01-04 | 2016-01-07 | 2016-01-05 | 2 | 6 |
| 8 | 8 | 1204 | 4 | 2016-01-04 | 2016-01-05 | 2016-01-05 | 2 | 7 |

| | EarliestOrderDate | LatestOrderDate |
|---|---|---|
| 1 | 2016-01-01 | 2018-12-28 |

| | EarliestRequiredDate | LatestRequiredDate |
|---|---|---|
| 1 | 2016-01-03 | 2018-12-28 |

Dipawoli Malla

Big Data and Data Mining

Potential Data quality issues can be

- Checking for duplicate products

```
dataRange.sql - AL...es (ALINA\D M (74))    SQLQuery6.sql - AL...s (ALINA\D M (55))*    X  SQLQuery5.sql - AL...s (ALINA\D M (54))*
    -- Duplicate products by name, brand, and category
  SELECT product_name, brand_id, category_id, COUNT(*) AS DuplicateCount
    FROM production.products
    GROUP BY product_name, brand_id, category_id
    HAVING COUNT(*) > 1;
```

121 %

Results    Messages

| product_name | brand_id | category_id | DuplicateCount |
|---|---|---|---|

Since the output is nothing, there are duplicate products.

- Checking for duplicate customers

```
dataRange.sql - AL...es (ALINA\D M (74))    SQLQuery6.sql - AL...s (ALINA\D M (55))*    X  SQLQuery5.sql - AL...s (ALINA\D M (54))*
    -- Duplicate customers by name and phone number
  SELECT first_name, last_name, phone, COUNT(*) AS DuplicateCount
    FROM sales.customers
    GROUP BY first_name, last_name, phone
    HAVING COUNT(*) > 1;


  SELECT customer_id, street, city, state
    FROM sales.customers
    WHERE first_name = 'Justina'
    and last_name = 'Jenkins';
```

121 %

Results    Messages

| | first_name | last_name | phone | DuplicateCount |
|---|---|---|---|---|
| 1 | Justina | Jenkins | NULL | 2 |

| | customer_id | street | city | state |
|---|---|---|---|---|
| 1 | 315 | 8236 Creek St. | Baldwin | NY |
| 2 | 1425 | 345 SE. Green Lane | Shirley | NY |

Since the output is Justina Jenkins with DuplicateCount of 2, it means the data is duplicated once. As we can see, it has two customer_id of same name but with different street, city and state which means they are two different people.

Dipawoli Malla

- Negative or zero pricing in products.

```
-- Products with zero or negative price
SELECT *
FROM production.products
WHERE list_price <= 0;
```

1 %

Results  Messages

| product_id | product_name | brand_id | category_id | model_year | list_price |
|---|---|---|---|---|---|

There was no data for pricing of the product with negative or zero.

- Orders with quantity = 0.

```
-- Orders with zero or negative quantity
SELECT *
FROM sales.order_items
WHERE quantity <= 0;
```

.21 %

Results  Messages

| order_id | item_id | product_id | quantity | list_price | discount |
|---|---|---|---|---|---|

The orders with quantity below 0 were found none.

## 3. Analysis

For analysis, I wanted to analyze the popular product using the total quantity and information related to stock management.

# Big Data and Data Mining

```sql
--Analysis for popular products with Total Quantity sold and other related information

SELECT
    p.product_name,
    SUM(oi.quantity) AS TotalQuantitySold,
    SUM((oi.list_price * oi.quantity) - oi.discount) AS TotalRevenue,
    c.category_name,
    b.brand_name,
    s.store_name
FROM sales.order_items oi
JOIN production.products p ON oi.product_id = p.product_id
JOIN production.categories c ON p.category_id = c.category_id
JOIN production.brands b ON p.brand_id = b.brand_id
JOIN sales.orders o ON oi.order_id = o.order_id
JOIN sales.stores s ON o.store_id = s.store_id
GROUP BY p.product_name, c.category_name, b.brand_name, s.store_name
ORDER BY TotalQuantitySold DESC, TotalRevenue DESC;
```

Result of the query:

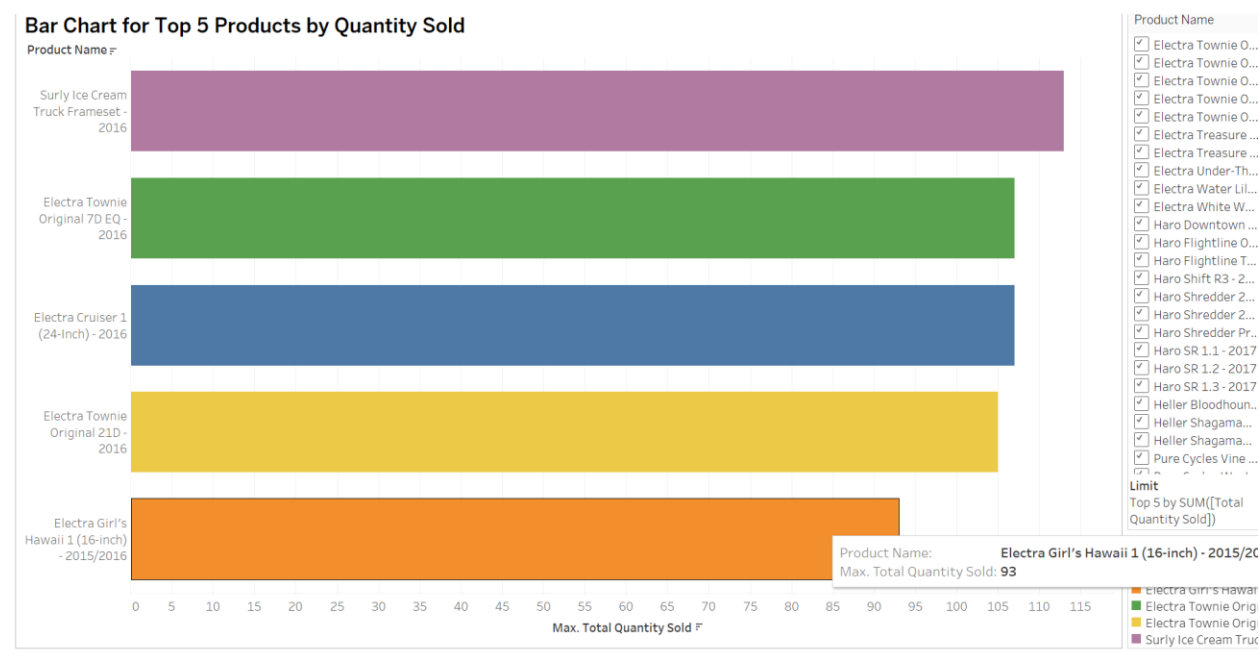| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | product_name | TotalQuantitySold | TotalRevenue | category_name | brand_name | store_name |
| 2 | Surly Ice Cream Truck Framese | 113 | 53101.07 | Mountain Bikes | Surly | Baldwin Bikes |
| 3 | Electra Girl's Hawaii 1 (20-inch | 111 | 33290.83 | Children Bicycles | Electra | Baldwin Bikes |
| 4 | Electra Townie Original 7D EQ - | 107 | 64191.8 | Cruisers Bicycles | Electra | Baldwin Bikes |
| 5 | Electra Cruiser 1 (24-Inch) - 20 | 107 | 28881.89 | Cruisers Bicycles | Electra | Baldwin Bikes |
| 6 | Electra Townie Original 21D - 2 | 105 | 57741.5 | Cruisers Bicycles | Electra | Baldwin Bikes |
| 7 | Electra Cruiser 1 (24-Inch) - 20 | 104 | 28073 | Children Bicycles | Electra | Baldwin Bikes |
| 8 | Electra Townie Original 7D - 20 | 103 | 51491.89 | Comfort Bicycles | Electra | Baldwin Bikes |
| 9 | Trek Slash 8 27.5 - 2016 | 101 | 403992.11 | Mountain Bikes | Trek | Baldwin Bikes |
| 10 | Surly Straggler 650b - 2016 | 101 | 169772.25 | Cyclocross Bicycles | Surly | Baldwin Bikes |
| 11 | Surly Straggler - 2016 | 100 | 154892.45 | Cyclocross Bicycles | Surly | Baldwin Bikes |
| 12 | Pure Cycles Western 3-Speed - | 99 | 44444.46 | Cruisers Bicycles | Pure Cycles | Baldwin Bikes |
| 13 | Electra Townie Original 7D EQ - | 98 | 58792.26 | Comfort Bicycles | Electra | Baldwin Bikes |
| 14 | Electra Townie Original 21D - 2 | 98 | 53892.55 | Comfort Bicycles | Electra | Baldwin Bikes |
| 15 | Trek Conduit+ - 2016 | 93 | 278992.68 | Electric Bikes | Trek | Baldwin Bikes |
| 16 | Electra Girl's Hawaii 1 (16-inch | 93 | 25102.48 | Children Bicycles | Electra | Baldwin Bikes |
| 17 | Trek Fuel EX 8 29 - 2016 | 91 | 263891.49 | Mountain Bikes | Trek | Baldwin Bikes |
| 18 | Heller Shagamaw Frame - 2016 | 90 | 118881.21 | Mountain Bikes | Heller | Baldwin Bikes |
| 19 | Electra Girl's Hawaii 1 (16-inch | 88 | 23752.06 | Cruisers Bicycles | Electra | Baldwin Bikes |
| 20 | Pure Cycles Vine 8-Speed - 201 | 87 | 37317.16 | Cruisers Bicycles | Pure Cycles | Baldwin Bikes |
| 21 | Electra Townie Original 7D EQ - | 85 | 50993.13 | Cruisers Bicycles | Electra | Baldwin Bikes |
| 22 | Trek Remedy 29 Carbon Frame | 84 | 151193.3 | Mountain Bikes | Trek | Baldwin Bikes |

Dipawoli Malla

4. **Visualization using Tableau**

Link to Tableu : Data Visualization | Tableau Public
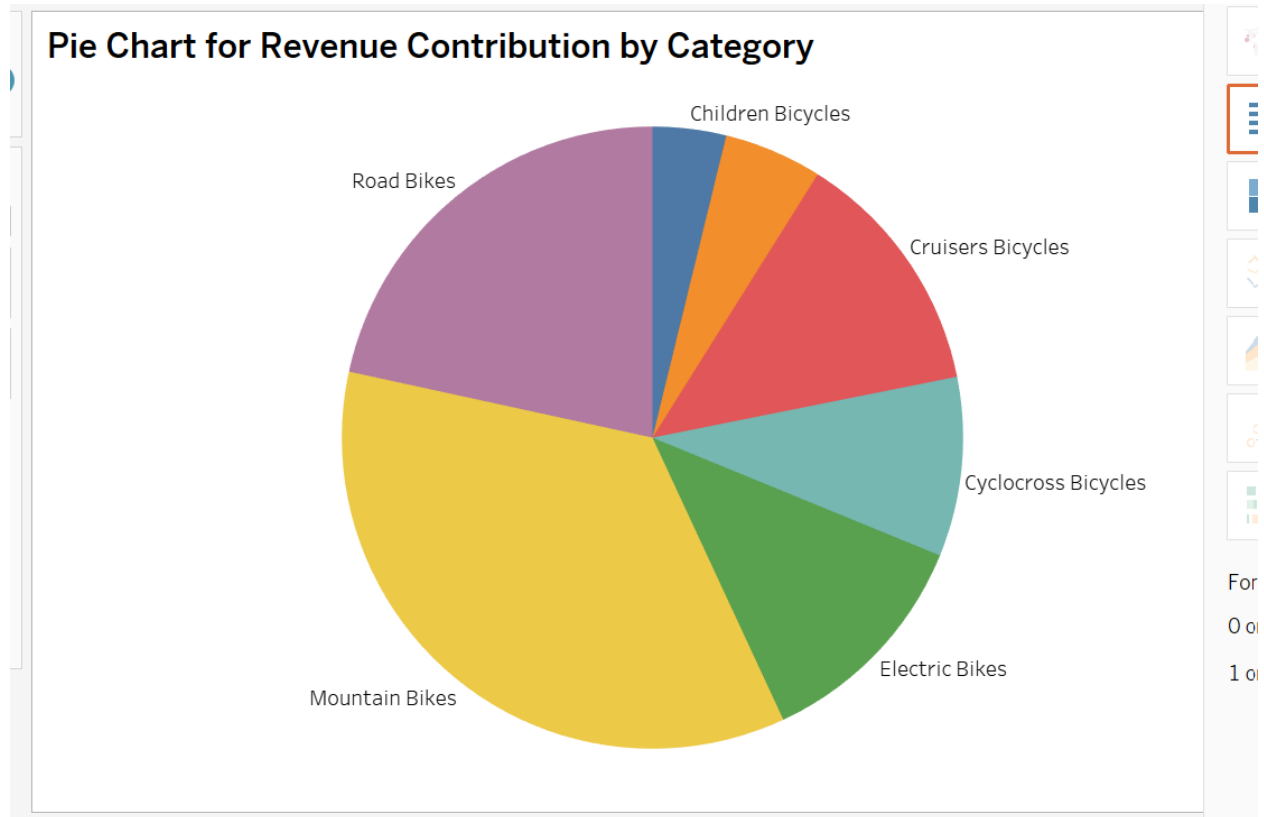
**Bar Chart: Top 5 Products by Quantity Sold**

A bar chart is ideal for comparing discrete items like products, showcasing their popularity based on sales volume (Total Quantity Sold). It ranks the products, making it easy to see which products customers prefer the most.



- The top product, "Surly Ice Cream Truck," sold the highest quantity (167 units), followed by "Electra Cruiser 1" (157 units).
- The chart helps businesses identify which products to stock more of and which to promote further.
- Products with higher quantities might be budget-friendly or have higher demand among a broader customer base.

**Pie Chart: Revenue Contribution by Category**

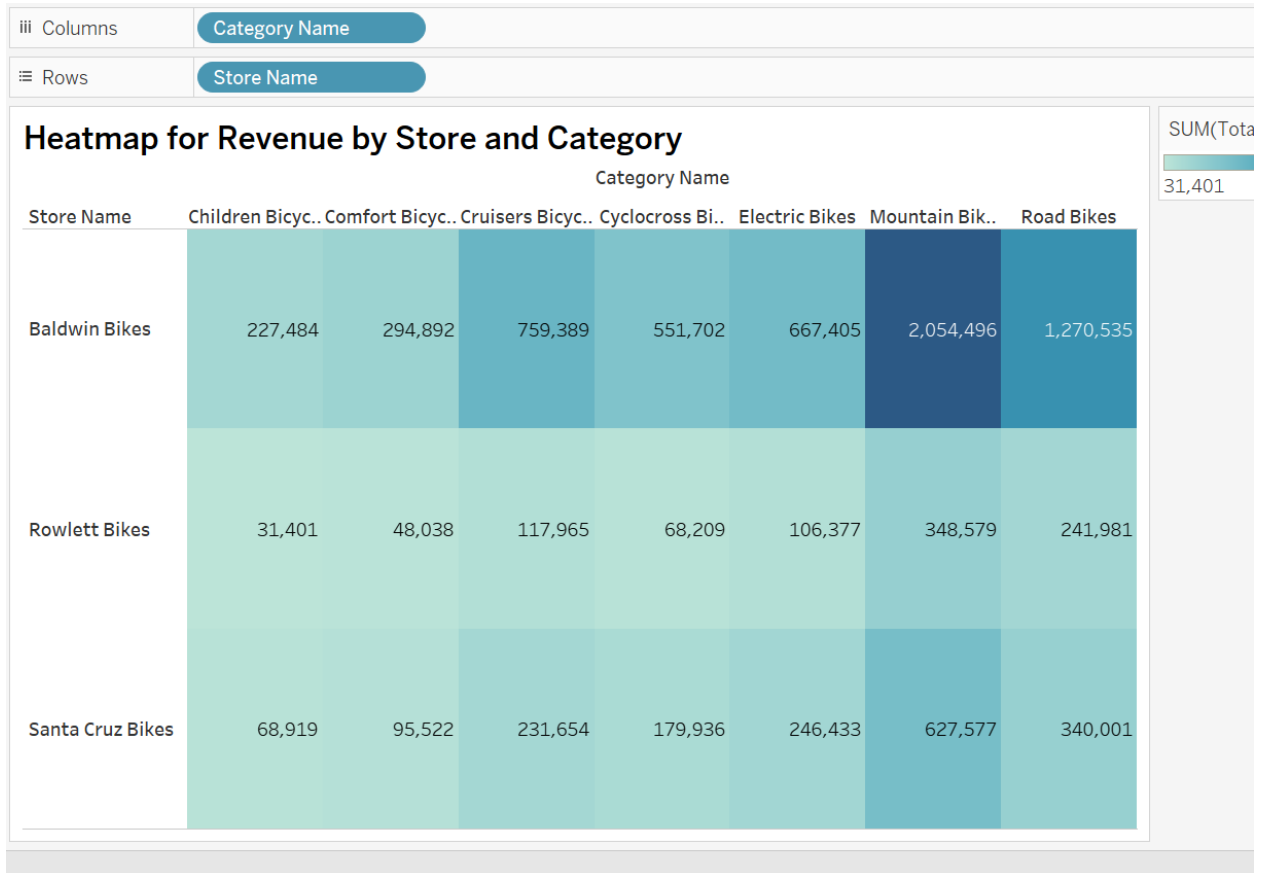Dipawoli Malla

Big Data and Data Mining

Pie charts are great for showing proportional data. In this case, it illustrates the share of each category in the total revenue. The visual focus is on understanding category dominance in revenue generation.



- "Mountain Bikes" contribute the highest revenue (3,030,651), followed by "Road Bikes" (1,852,516).
- Categories like "Children Bicycles" (327,804) and "Comfort Bicycles" (438,452) contribute significantly less.
- Businesses can use this to allocate resources and marketing budgets to higher-revenue categories and explore ways to boost sales in underperforming categories.

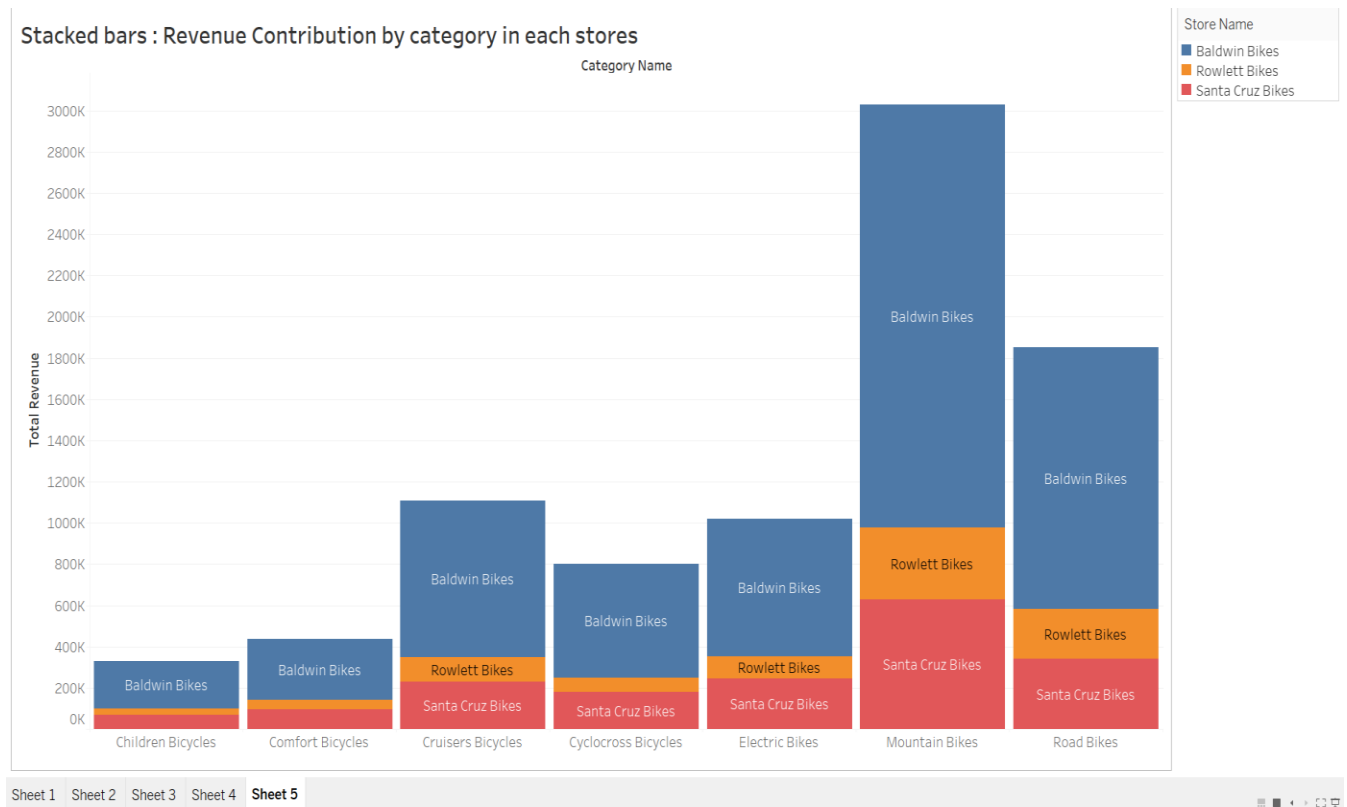Dipawoli Malla

**Heat Map: Revenue by Store and Category**

Heat maps visually represent data intensity using colors, making it easy to identify areas of high and low performance across multiple dimensions. This chart shows revenue for each category broken down by store.



- "Mountain Bikes" dominate in all stores, with Baldwin Bikes generating the highest revenue in this category (2,054,496).
- "Children Bicycles" and "Comfort Bicycles" perform poorly across all stores, indicating a need for reevaluation of inventory or marketing strategies for these categories.
- Baldwin Bikes is the highest-performing store overall, suggesting it might have a better location or customer base for premium products.

Dipawoli Malla

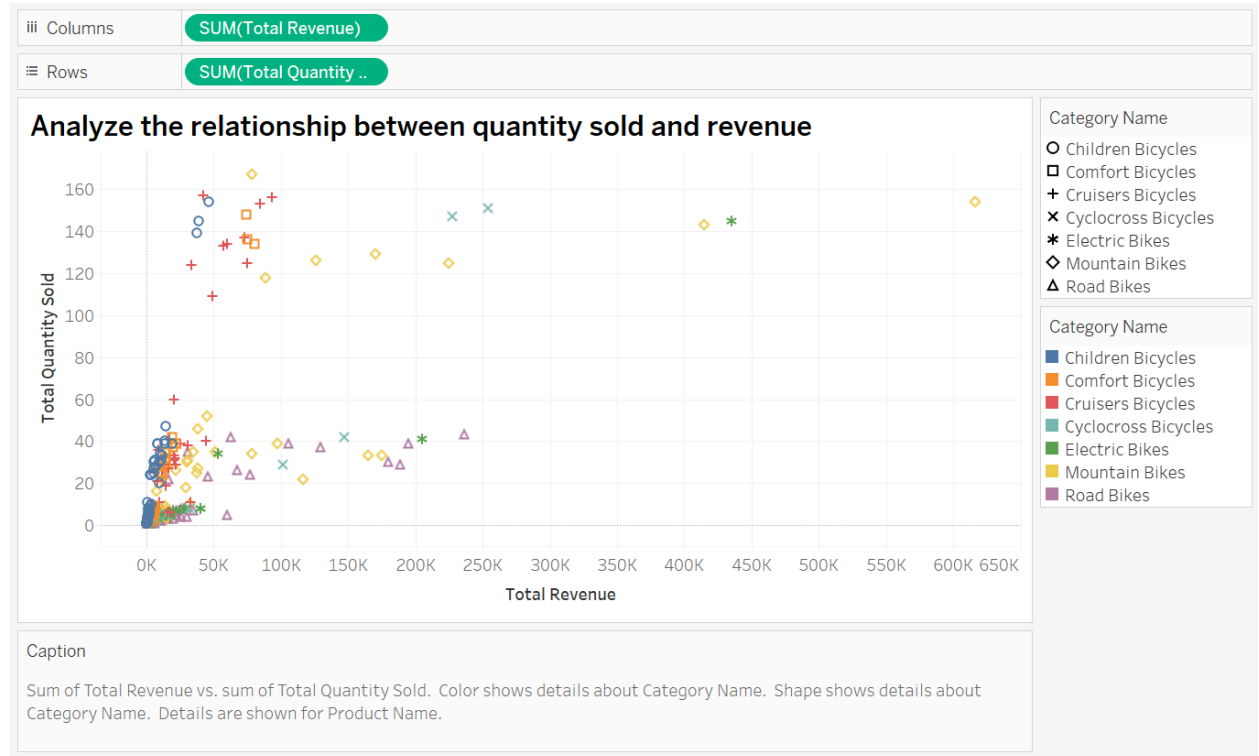**Stacked Chart: Revenue Contribution by Category with Store Names**

A stacked bar chart provides a clear breakdown of total revenue by category and further divides it by store. It enables businesses to analyze both the category-level performance and the store-specific contributions in one chart.



- "Mountain Bikes" generate the highest total revenue across all categories, with Baldwin Bikes contributing the largest share, followed by Rowlett Bikes and Santa Cruz Bikes.
- This indicates that Baldwin Bikes is a key driver for high-revenue products like Mountain Bikes, and efforts to boost sales at Rowlett Bikes and Santa Cruz Bikes could significantly increase overall revenue.

**Scatter Plot: Quantity Sold vs. Revenue**

Scatter plots are used to analyze relationships between two variables, here Total Quantity Sold and Total Revenue. The goal is to identify patterns, outliers, and insights about products.



| iii Columns | SUM(Total Revenue) |
| Rows | SUM(Total Quantity ..) |

**Analyze the relationship between quantity sold and revenue**

Category Name
○ Children Bicycles
□ Comfort Bicycles
+ Cruisers Bicycles
✕ Cyclocross Bicycles
✱ Electric Bikes
◇ Mountain Bikes
△ Road Bikes

Category Name
■ Children Bicycles
■ Comfort Bicycles
■ Cruisers Bicycles
■ Cyclocross Bicycles
■ Electric Bikes
■ Mountain Bikes
■ Road Bikes

Total Quantity Sold / Total Revenue

Caption

Sum of Total Revenue vs. sum of Total Quantity Sold. Color shows details about Category Name. Shape shows details about Category Name. Details are shown for Product Name.

- Products with high revenue but low quantity sold might be premium items with high margins.
- Products with high quantities but low revenue could be affordable items with thin margins but high demand.
- The scatter plot highlights opportunities to optimize pricing or focus on products that balance revenue and volume.

Dipawoli Malla