



Analysis of E-commerce Dataset

Using Python and MySQL Workbench

—Sathwika Reddy

Introduction

- ▶ The primary purpose of this project is to perform a comprehensive analysis of an e-commerce dataset to extract valuable insights that can drive business decisions.
- ▶ The dataset, downloaded from Kaggle, consists of multiple CSV files containing various aspects of e-commerce operations such as orders, customers, products, Order items, Payments, Geographical locations, Sellers.



Steps undertaken in the project include:

1. **Data Collection:** Downloaded the dataset from Kaggle and organized it into a folder named "Ecommerce."
2. **Data Loading:** Utilized Python and Pandas to read the CSV files and load the data into MySQL Workbench.
3. **Data Analysis:** Conducted detailed data analysis using SQL to uncover trends, customer behavior, and product performance.

Purpose of Loading the Dataset Using Python for Process Optimization

Efficient Data Handling:

Python, with libraries like Pandas, handles large datasets efficiently and provides powerful data manipulation capabilities.

Automated Workflow:

Python scripts automate data loading and preprocessing, making the process repeatable, reducing human error, and saving time.



Integration with Databases:

Python easily connects to databases like MySQL, ensuring smooth data integration and enabling complex, customized data analysis.

Scalability and Flexibility:

Python solutions can scale with increasing data volumes and adapt to different data sources and formats.

Data Analysis Approach

- `create database ecommerce;`
- `use ecommerce;`
 - Basic Queries
 - -1. List all unique cities where customers are located.
- `select * from customers;`
- `select distinct customer_city from customers;`
 - 2. Count the number of orders placed in 2017.
- `select * from orders;`
- `select count(order_id) as total_orders_in_2017 from orders
where year(order_purchase_timestamp)=2017;`

```
14
15     -- 3. Find the total sales per category.
16
17 • select * from products;
18 • select * from order_items;
19
20 • select p1.product_category, sum(o.price) as total_price from products p1
21     join order_items o
22     on p1.product_id=o.product_id
23     group by p1.product_category;
24
25     -- 4. Calculate the percentage of orders that were paid in installments.
26 • select * from payments;
27
28 • select (sum(case when payment_installments>=1 then 1 else 0 end)/count(*) )*100 as percentage from payments;
29
```

-- 5. Count the number of customers from each state.

select * from customers;

select customer_state ,count(customer_id)
from customers group by customer_state;

-- Calculate the number of orders per month in 2018.

select * from orders;

select month(order_purchase_timestamp) as month,count(order_id) as total_orders from orders
where year(order_purchase_timestamp)=2018
group by month(order_purchase_timestamp)
order by month(order_purchase_timestamp) asc;


```
44 -- Find the average number of products per order, grouped by customer city.
45
46 • select * from order_items;
47 • select * from orders;
48 • select * from customers;
49
50 • ⊖ with count_per_order as(
51     select orders.order_id, orders.customer_id, count(order_items.order_id) as oc
52     from orders join order_items
53     on orders.order_id = order_items.order_id
54     group by orders.order_id, orders.customer_id
55 )
56     select customers.customer_city, round(avg(count_per_order.oc),2) average_orders
57     from customers join count_per_order
58     on customers.customer_id = count_per_order.customer_id
59     group by customers.customer_city order by average_orders desc;
60
```

```
60
61 -- Calculate the percentage of total revenue contributed by each product category.
62
63 • select * from order_items;
64 • select * from products;
65
66 • select * from payments;
67
68 • select upper(products.product_category) category,
69 round((sum(payments.payment_value)/(select sum(payment_value) from payments))*100,2) sales_percentage
70 from products join order_items
71 on products.product_id = order_items.product_id
72 join payments
73 on payments.order_id = order_items.order_id
74 group by category order by sales_percentage desc;
75
76
```

```
-- Identify the correlation between product price and the number of times a product has been purchased.
```

- ```
select products.product_category,
count(order_items.product_id) as total_count,
round(avg(order_items.price),2) as avg_price
from products join order_items
on products.product_id = order_items.product_id
group by products.product_category;
```

```
-- Calculate the total revenue generated by each seller, and rank them by revenue.
```

- ```
select *, dense_rank() over(order by revenue desc) as rn from  
(select order_items.seller_id, sum(payments.payment_value)  
revenue from order_items join payments  
on order_items.order_id = payments.order_id  
group by order_items.seller_id) as a;
```

```
4
5 -- 1.Calculate the moving average of order values for each customer over their order history.
6 • WITH cte AS (
7     SELECT
8         orders.customer_id,
9         orders.order_purchase_timestamp,
10        payments.payment_value AS payment,
11        AVG(payments.payment_value) OVER (
12            PARTITION BY orders.customer_id
13            ORDER BY orders.order_purchase_timestamp
14            ROWS BETWEEN 2 PRECEDING AND CURRENT ROW
15        ) AS mov_avg
16    FROM
17        payments
18    JOIN
19        orders ON payments.order_id = orders.order_id
20 )
21 SELECT
22     customer_id,
23     order_purchase_timestamp,
24     payment,
25     mov_avg
26 FROM
27     cte
28 ORDER BY
29     customer_id,
30     order_purchase_timestamp;
```

```
122
123 -- 2. Calculate the cumulative sales per month for each year.
124 • select * from orders;
125 • select * from payments;
126 • WITH sales_data AS (
127     SELECT
128         YEAR(orders.order_purchase_timestamp) AS years,
129         MONTH(orders.order_purchase_timestamp) AS months,
130         ROUND(SUM(payments.payment_value), 2) AS payment
131     FROM orders
132     JOIN payments ON orders.order_id = payments.order_id
133     GROUP BY
134         YEAR(orders.order_purchase_timestamp),
135         MONTH(orders.order_purchase_timestamp)
136     ORDER BY
137         YEAR(orders.order_purchase_timestamp),
138         MONTH(orders.order_purchase_timestamp)),
139 • cumulative_sales_data AS (
140     SELECT years, months, payment,
141         SUM(payment) OVER (ORDER BY years, months) AS cumulative_sales
142     FROM sales_data)
143 SELECT years, months, payment, cumulative_sales
144 FROM cumulative_sales_data
145 ORDER BY years, months;
146
```



```
148 -- 3. Calculate the year-over-year growth rate of total sales.
149
150 • WITH a AS (
151     SELECT
152         YEAR(orders.order_purchase_timestamp) AS years,
153         ROUND(SUM(payments.payment_value), 2) AS payment
154     FROM orders
155     JOIN payments ON orders.order_id = payments.order_id
156     GROUP BY
157         YEAR(orders.order_purchase_timestamp)
158     ORDER BY
159         YEAR(orders.order_purchase_timestamp)
160 ),
161 percentage_change AS (
162     SELECT years, payment,
163         ((payment - LAG(payment, 1) OVER (ORDER BY years)) /
164         LAG(payment, 1) OVER (ORDER BY years)) * 100 AS percentage_change
165     FROM a)
166 SELECT years, payment, percentage_change
167 FROM percentage_change
168 ORDER BY years;
169
```

```
169
170
171 -- 4. Identify the top 3 customers who spent the most money in each year.
172
173 WITH customer_payments AS (
174     SELECT
175         YEAR(orders.order_purchase_timestamp) AS years,
176         orders.customer_id,
177         SUM(payments.payment_value) AS payment,
178         DENSE_RANK() OVER (
179             PARTITION BY YEAR(orders.order_purchase_timestamp)
180             ORDER BY SUM(payments.payment_value) DESC
181         ) AS d_rank
182     FROM orders
183     JOIN payments ON payments.order_id = orders.order_id
184     GROUP BY YEAR(orders.order_purchase_timestamp), orders.customer_id
185 )
186 SELECT years, customer_id, payment, d_rank
187 FROM customer_payments
188 WHERE d_rank <= 3
189 ORDER BY years, d_rank;
190
191
```

Key Findings

Sales Trends:

Monthly/Yearly Sales Growth: Identified periods with significant increases or decreases in sales.

Seasonal Patterns: Observed peaks in sales during certain months, indicating seasonal trends.

Customer Behavior:

Top Customers: Identified customers who make the most frequent or highest value purchases.

Average Order Value: Calculated the average value of customer orders to understand purchasing power.

Customer Segmentation: Grouped customers into segments based on purchasing behavior and demographics.

Product Performance:

Best-Selling Products: Determined which products have the highest sales volume.

Product Categories: Analyzed which categories perform best in terms of sales and customer preference.

Inventory Turnover: Assessed how quickly products move through inventory.

Revenue Analysis:

Total Revenue: Calculated total revenue generated over different periods.

Revenue by Category: Analyzed revenue contributions from different product categories.

Order Analysis:

Order Volume: Monitored the number of orders placed over time.

Order Status: Evaluated the distribution of order statuses (e.g., completed, canceled, returned).

Customer Retention:

Repeat Purchase Rate: Measured the percentage of customers who make repeat purchases.

Customer Lifetime Value (CLV): Estimated the long-term value of customers based on purchasing behavior.

Geographical Insights:

Top Regions: Identified regions or locations with the highest sales.

Regional Preferences: Analyzed product preferences across different regions

Conclusion:

These key findings provide a comprehensive view of the e-commerce platform's performance, customer behavior, and product dynamics. By leveraging this data, stakeholders can make informed decisions to optimize operations, enhance customer satisfaction, and drive business growth.