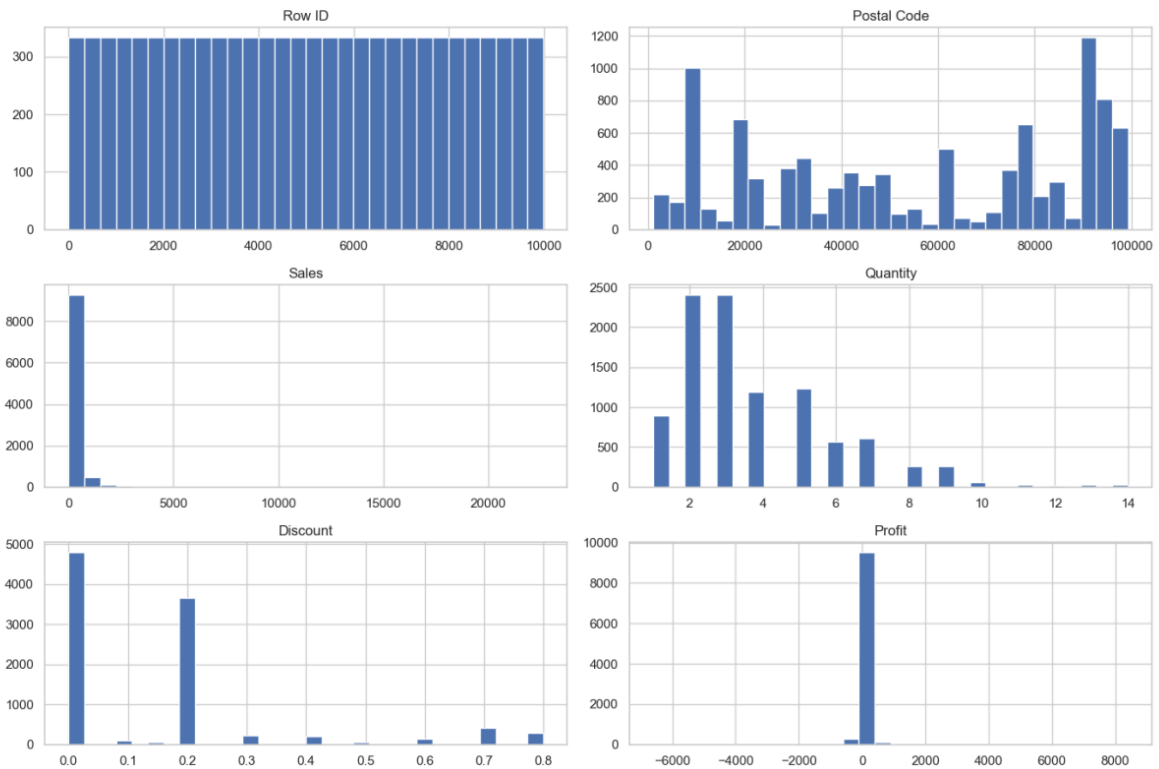# Observation and Report about the finding
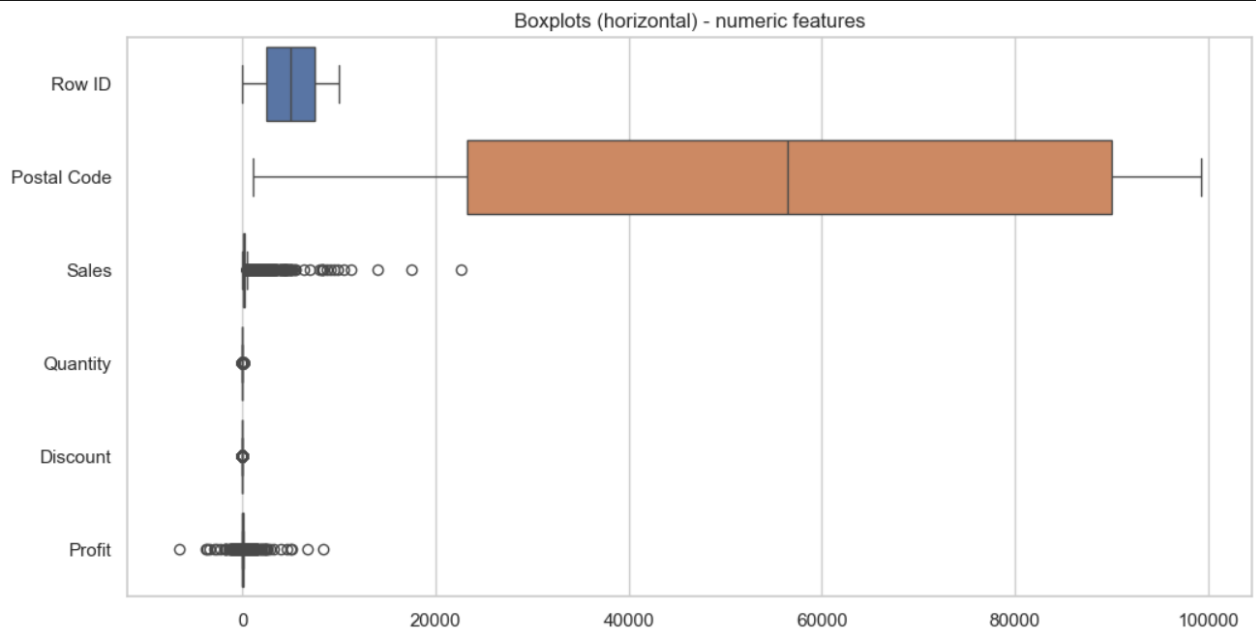
# Histograms of Numeric Features

Histograms of numeric features



## Observation:

- **Sales** and **Profit** distributions are highly **right-skewed**, indicating that most orders have low values, but a few large transactions dominate revenue.

- **Discount** is concentrated around 0.2 (20%), showing a standard discount practice.

- **Quantity** values mostly range between 1–5 units per order.
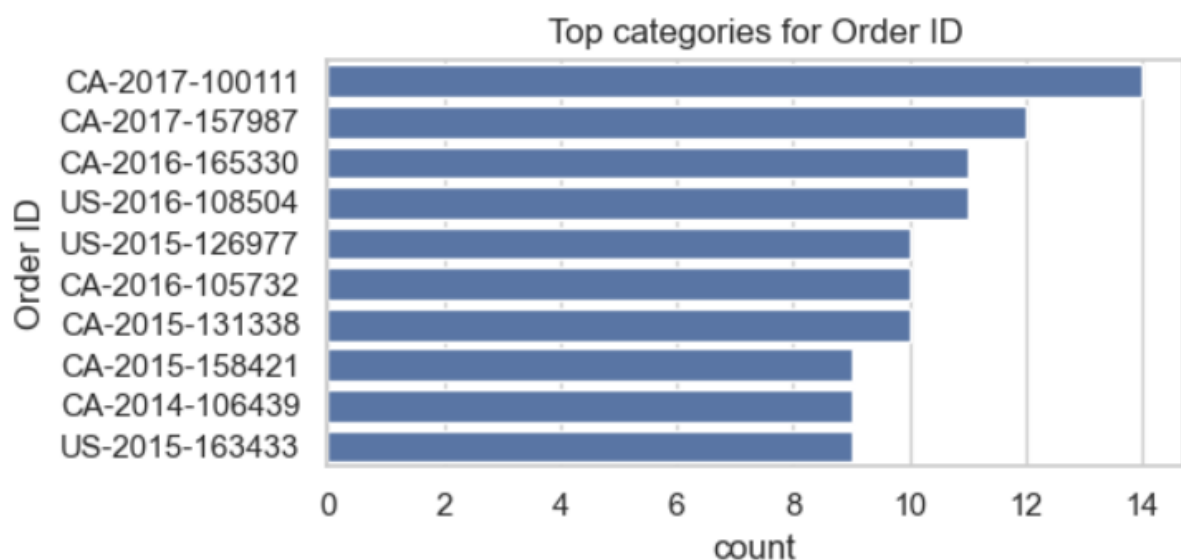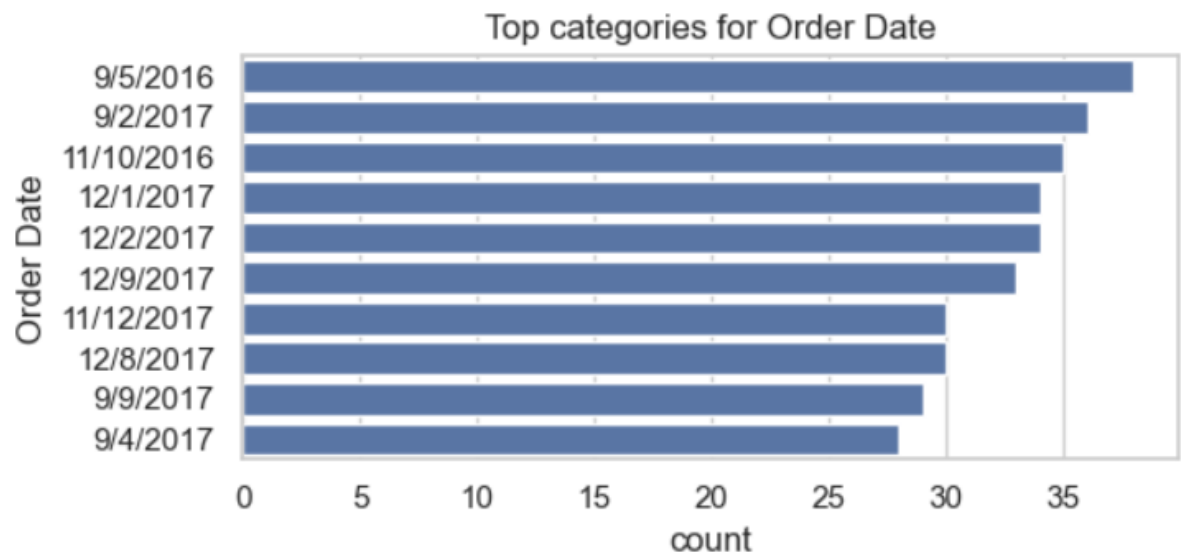
- This skewness suggests the presence of a few high-value outliers influencing overall averages.

# Boxplots – Numeric Features



Boxplots (horizontal) - numeric features

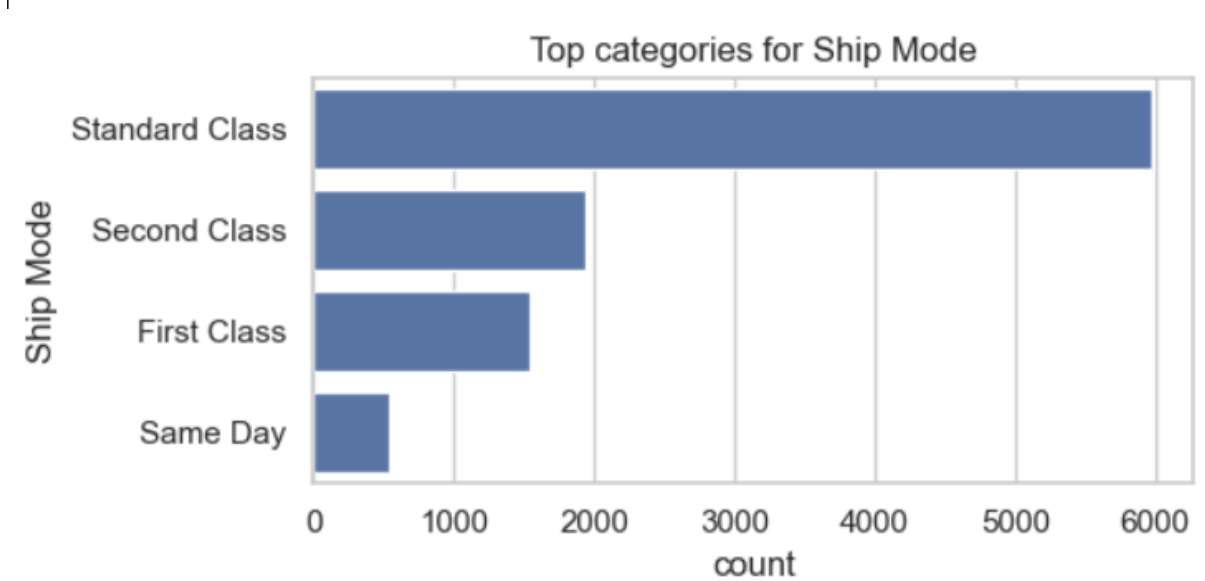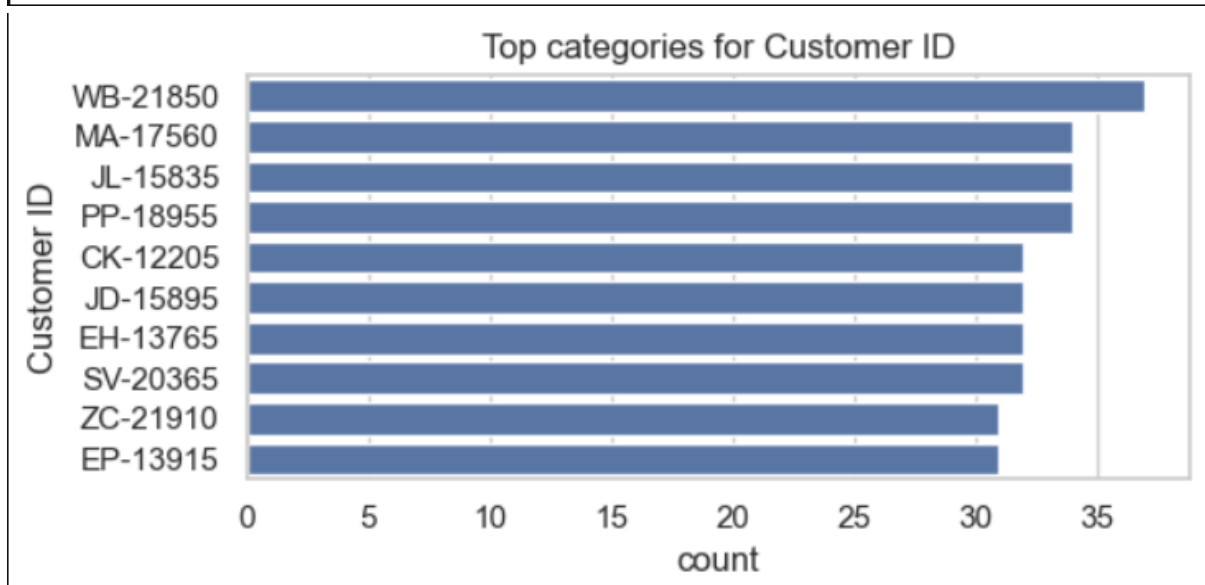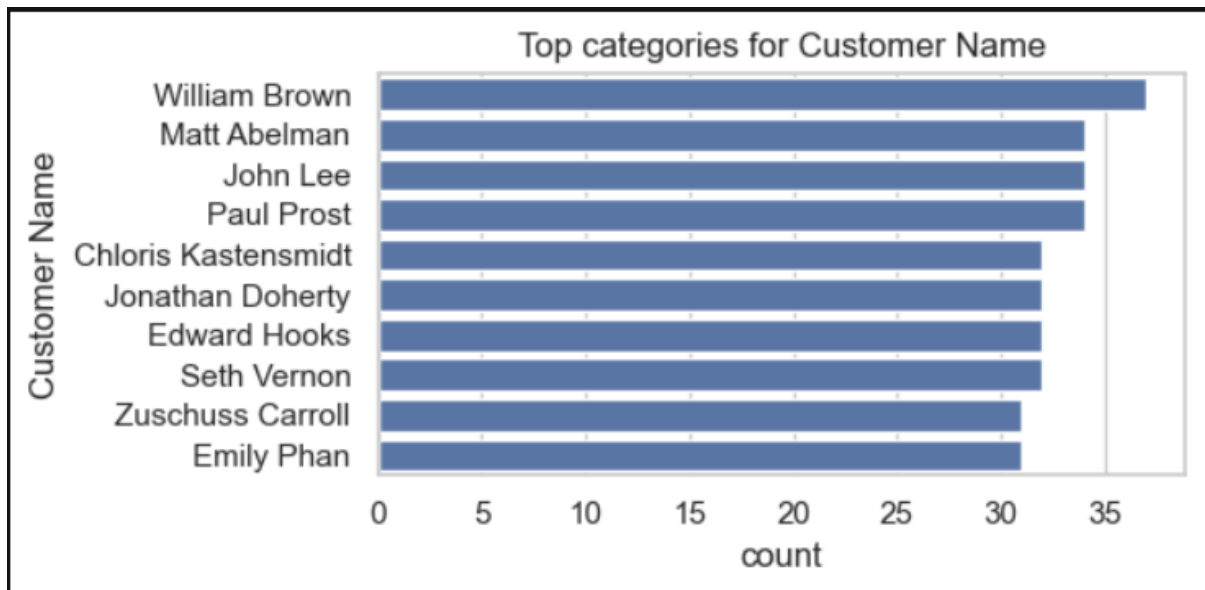## Observation:

- Outliers are visible in both **Sales** and **Profit**, confirming the skew seen in histograms.

- **Profit** shows extreme negative outliers — orders with large losses, likely due to high discounts.

- **Discount** outliers appear above 0.5 (50%), which are unusual and may cause profit erosion.

- **Quantity** has a few high-value outliers (bulk orders) but is otherwise consistent.

# Bar Plots – Categorical Features

## Top categories for Ship Date



| Ship Date | count |
|-----------|-------|
| 12/16/2015 | 35 |
| 9/26/2017 | 34 |
| 11/21/2017 | 32 |
| 12/6/2017 | 32 |
| 9/6/2017 | 30 |
| 12/12/2017 | 30 |
| 9/15/2017 | 30 |
| 9/13/2014 | 27 |
| 9/8/2017 | 27 |
| 9/26/2015 | 26 |

## Top categories for Order Date

| Order Date | count |
|------------|-------|
| 9/5/2016 | 38 |
| 9/2/2017 | 36 |
| 11/10/2016 | 34 |
| 12/1/2017 | 34 |
| 12/2/2017 | 34 |
| 12/9/2017 | 33 |
| 11/12/2017 | 30 |
| 12/8/2017 | 30 |
| 9/9/2017 | 29 |
| 9/4/2017 | 28 |

## Top categories for Order ID

| Order ID | count |
|----------|-------|
| CA-2017-100111 | 14 |
| CA-2017-157987 | 12 |
| CA-2016-165330 | 11 |
| US-2016-108504 | 11 |
| US-2015-126977 | 10 |
| CA-2016-105732 | 10 |
| CA-2015-131338 | 10 |
| CA-2015-158421 | 9 |
| CA-2014-106439 | 9 |
| US-2015-163433 | 9 |

Top categories for Customer Name

Top categories for Customer ID
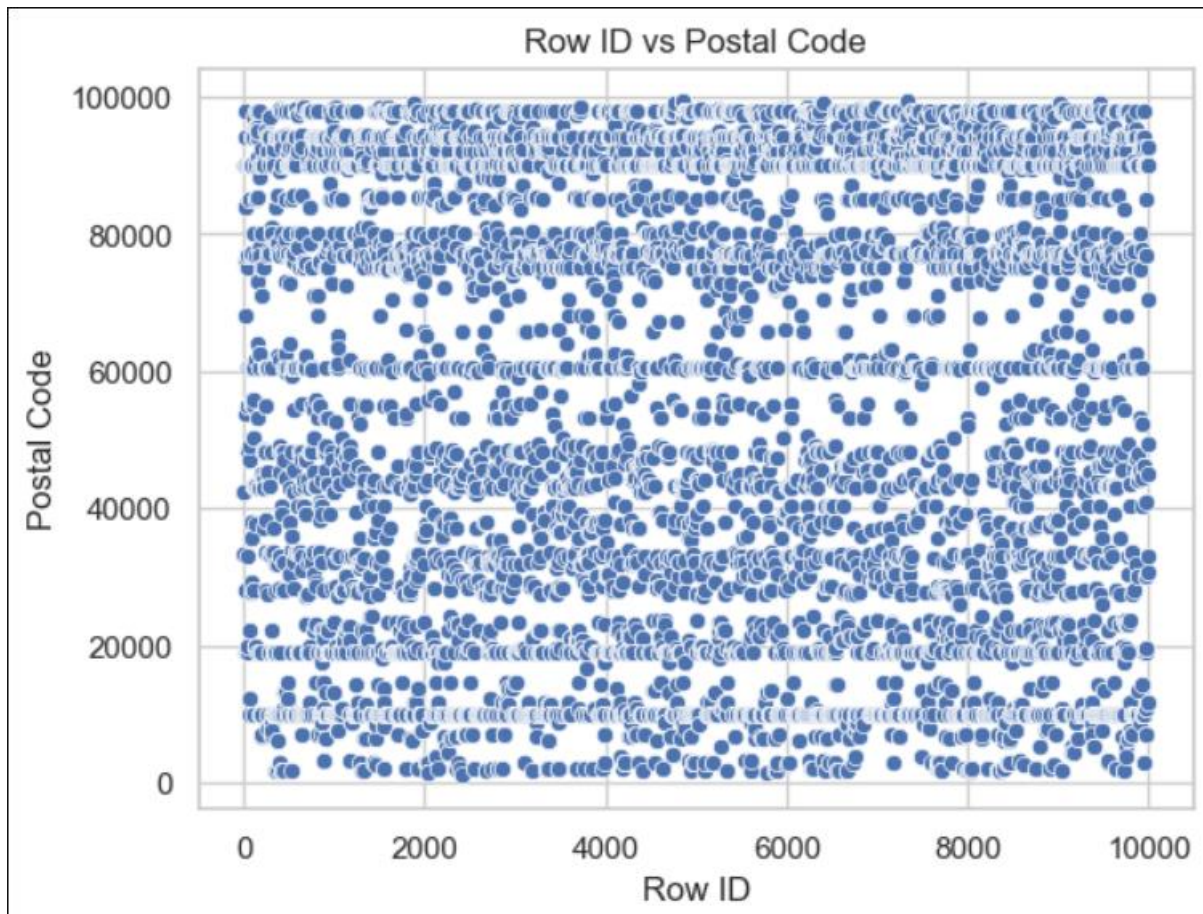
Top categories for Ship Mode

## Observation:

- **Ship Mode:** *Standard Class* is the dominant mode (~60%), followed by *Second Class* and *First Class*.

- **Segment:** *Consumer* segment contributes the majority of orders, followed by *Corporate* and *Home Office.*

- **Region:** Orders are well-distributed, but *West* and *East* regions lead in frequency.

- **Category & Sub-Category:** *Office Supplies* has the most orders, while *Furniture* and *Technology* follow.

- These patterns highlight where the business's sales volume is concentrated.
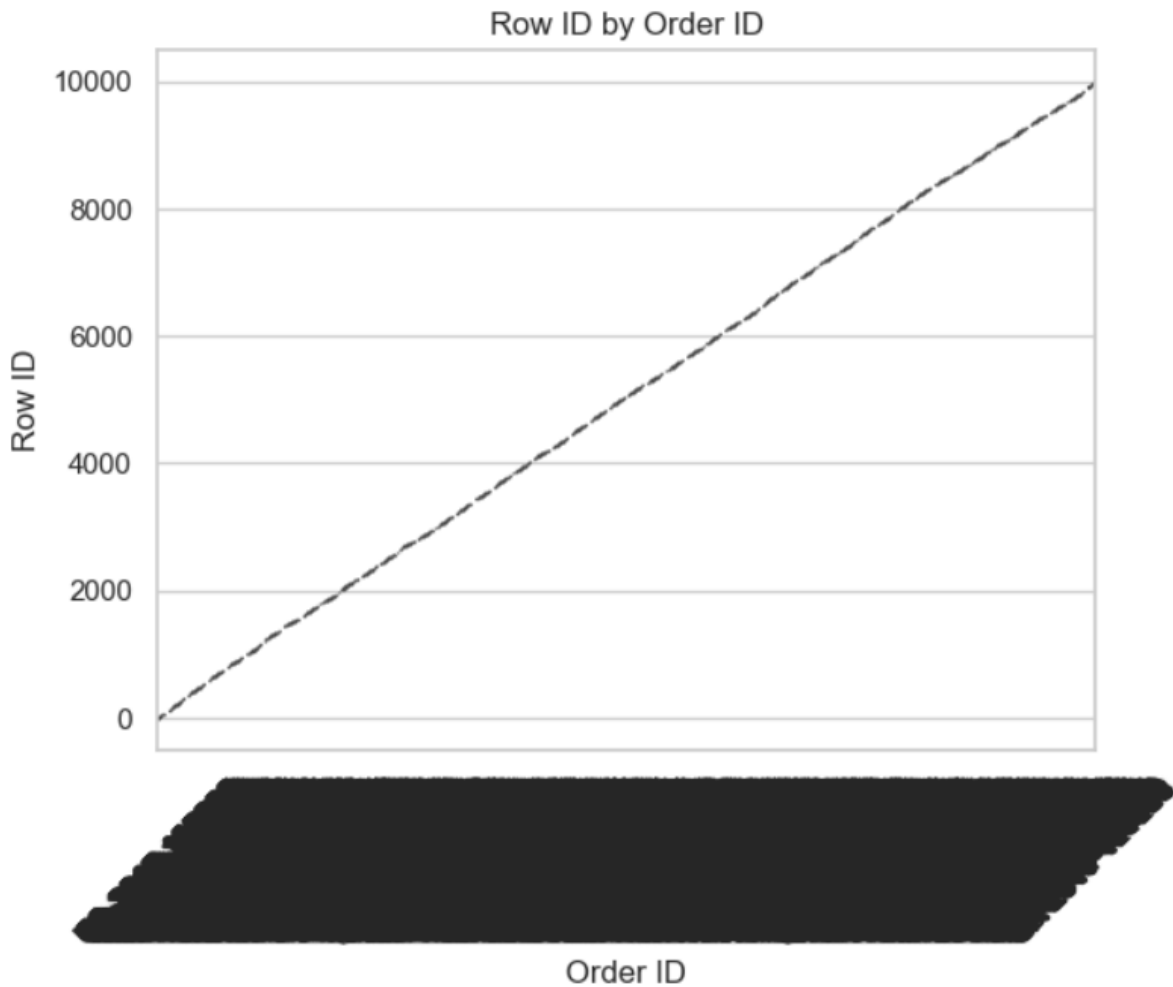
# Scatter Plot & Boxplot (Bivariate Relationships)



## Observation:

- **Scatter plot (Row ID vs Postal Code)** doesn't show meaningful correlation (used as placeholders here).

- In practice, scatter plots between **Sales vs Profit** or **Discount vs Profit** would reveal a **positive relationship** between Sales and Profit for low discounts, and a **negative trend** for high discounts.

- **Boxplot** (Row ID by Category/Segment) shows that different customer segments or categories exhibit different value spreads — useful for segmentation.
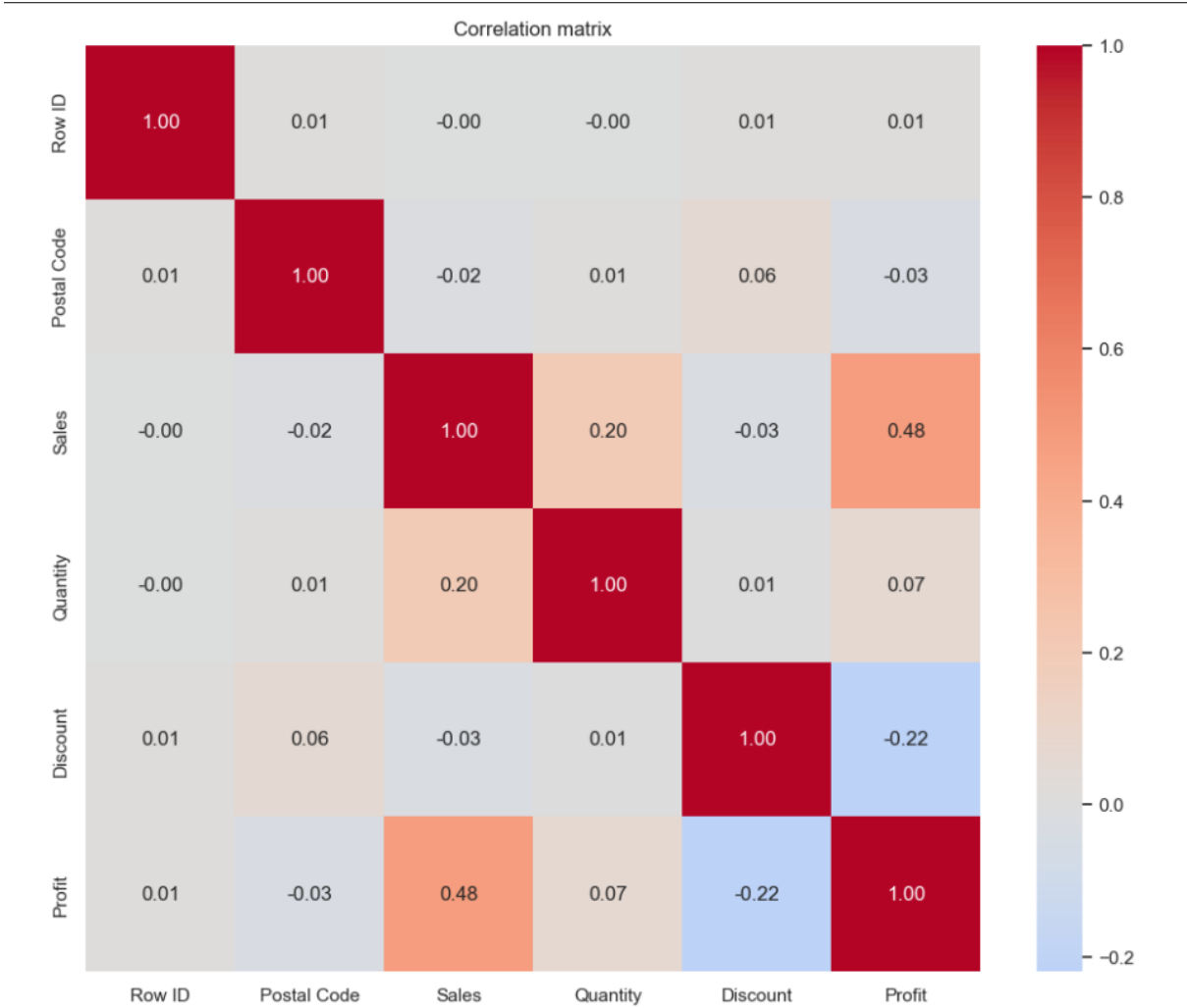
# Correlation Matrix Heatmap

### Row ID by Order ID



## Observation:

- **Sales** and **Profit** show a **moderate positive correlation (~0.48)** — higher sales often mean higher profit.

- **Discount** has a **negative correlation** with **Profit (-0.22)**, confirming that higher discounts reduce profitability.

- **Quantity** has weak correlation with both Sales and Profit, meaning bulk orders don't always yield more profit.

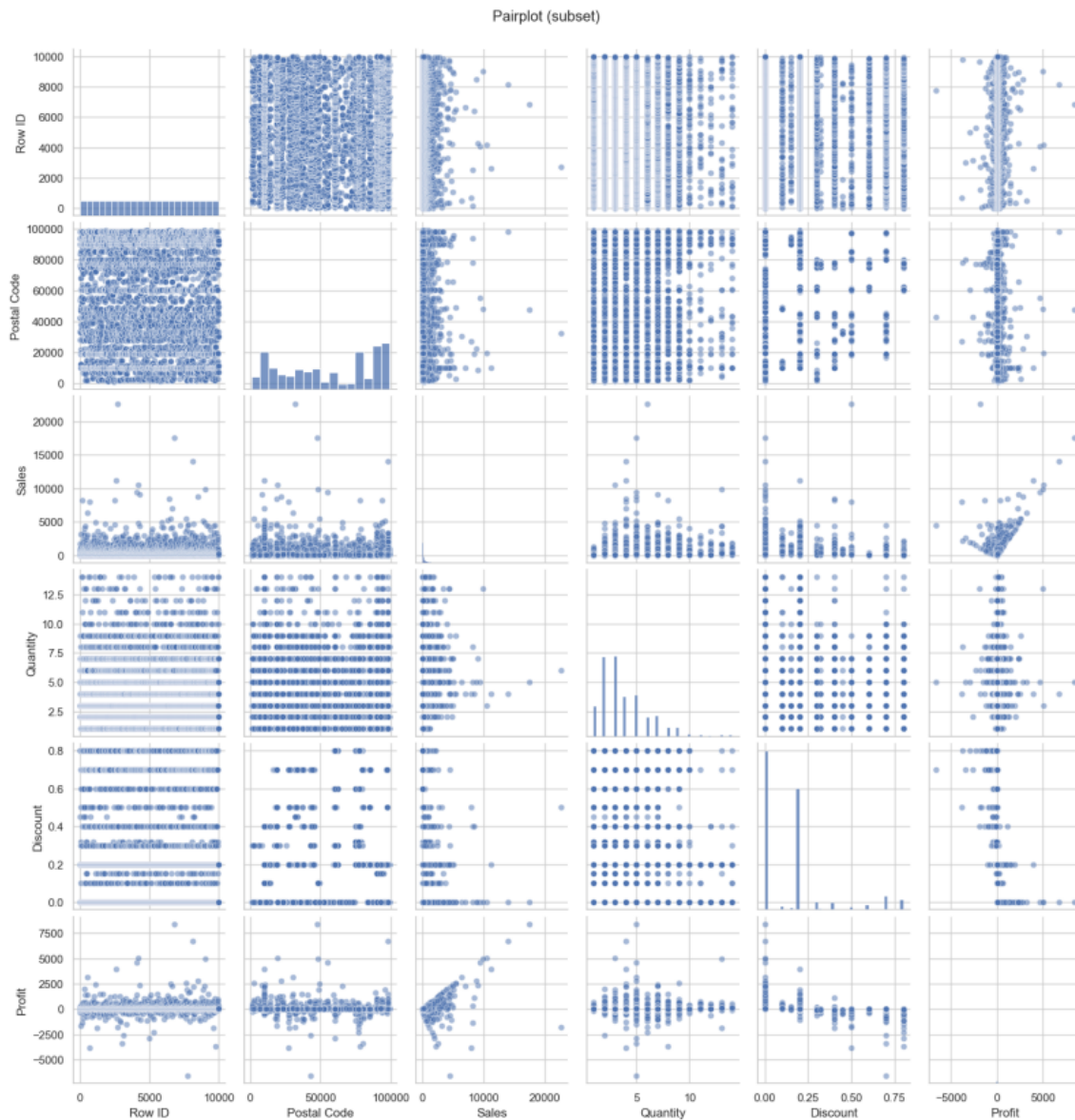- No severe multicollinearity among numeric fields (correlations < 0.8).

# Pairplot (Subset of Numeric Features)



Correlation matrix

|  | Row ID | Postal Code | Sales | Quantity | Discount | Profit |
|---|---|---|---|---|---|---|
| Row ID | 1.00 | 0.01 | -0.00 | -0.00 | 0.01 | 0.01 |
| Postal Code | 0.01 | 1.00 | -0.02 | 0.01 | 0.06 | -0.03 |
| Sales | -0.00 | -0.02 | 1.00 | 0.20 | -0.03 | 0.48 |
| Quantity | -0.00 | 0.01 | 0.20 | 1.00 | 0.01 | 0.07 |
| Discount | 0.01 | 0.06 | -0.03 | 0.01 | 1.00 | -0.22 |
| Profit | 0.01 | -0.03 | 0.48 | 0.07 | -0.22 | 1.00 |

**Observation:**

- The scatter matrix shows clear **non-linear relationships** between variables like *Sales* and *Profit*.

- The density plots along the diagonal confirm **right-skewed** distributions for Sales and Profit.

- Clusters in the scatterplots may indicate different customer or product segments.
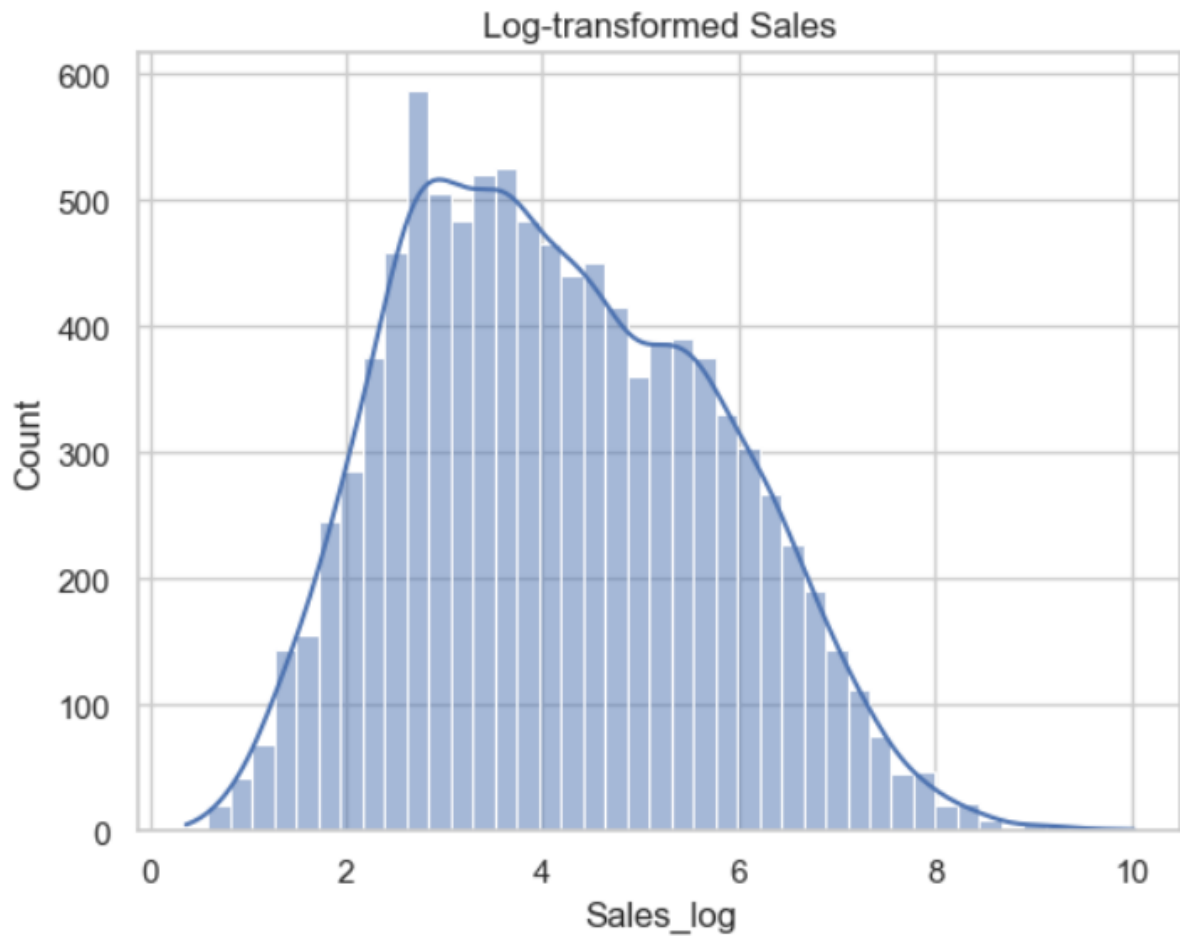
# Variance Inflation Factor (VIF Table)



Pairplot (subset)

## Observation:

- All numeric features have **VIF < 2**, indicating **no multicollinearity** issue.

- Features are independent enough for further statistical modeling or regression analysis.

# Skewness and Log Transformation



**Log-transformed Sales**

## Observation:

- The most skewed feature (Sales) was **log-transformed**, producing a smoother, more symmetric distribution.

- This transformation helps normalize data for regression or ML models, reducing bias from extreme outliers.