

SAMSUNG

# Superstore Sales Analysis

Team: 10



# Our Team



**Yousef Emad**



**Somaia Ahmed**



**Nour-Eldeen Anwar**



**Haneen El-Daly**

Together for Tomorrow!  
**Enabling People**

Education for Future Generations

# Overview

- Introduction
- Dataset Overview
- Business Goals
- EDA
- Data Preprocessing
- Modeling
- Deployment





# Introduction

Data analysis has become a critical tool for businesses seeking to drive profitability, optimize sales, and streamline operations.

Retailers generate vast amounts of data through every customer interaction, transaction, and product movement.

Harnessing this data allows businesses to gain deeper insights into customer behavior, product performance, and operational efficiencies.

# Dataset Overview:

- The "Superstore Sales" dataset captures transactional data from a retail store, covering various regions, product categories, and customer segments.
- It includes information such as sales, profit, discounts, shipping modes, and customer demographics.

# Business Goals:

- Understand how discounts impact profitability.
- Evaluate the impact of seasonality on sales trends.
- Identify the most profitable product combinations.

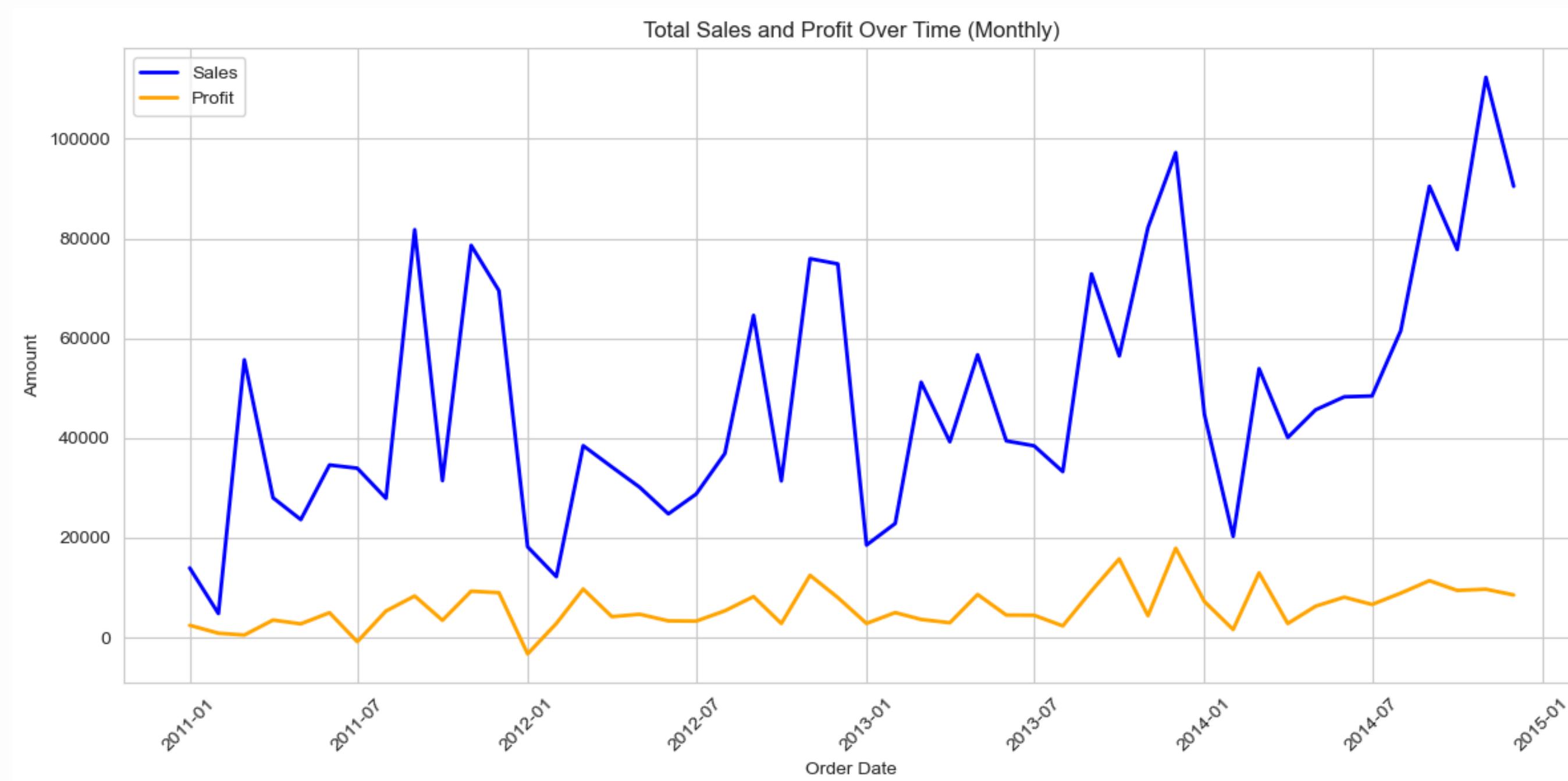


# SALES PERFORMANCE

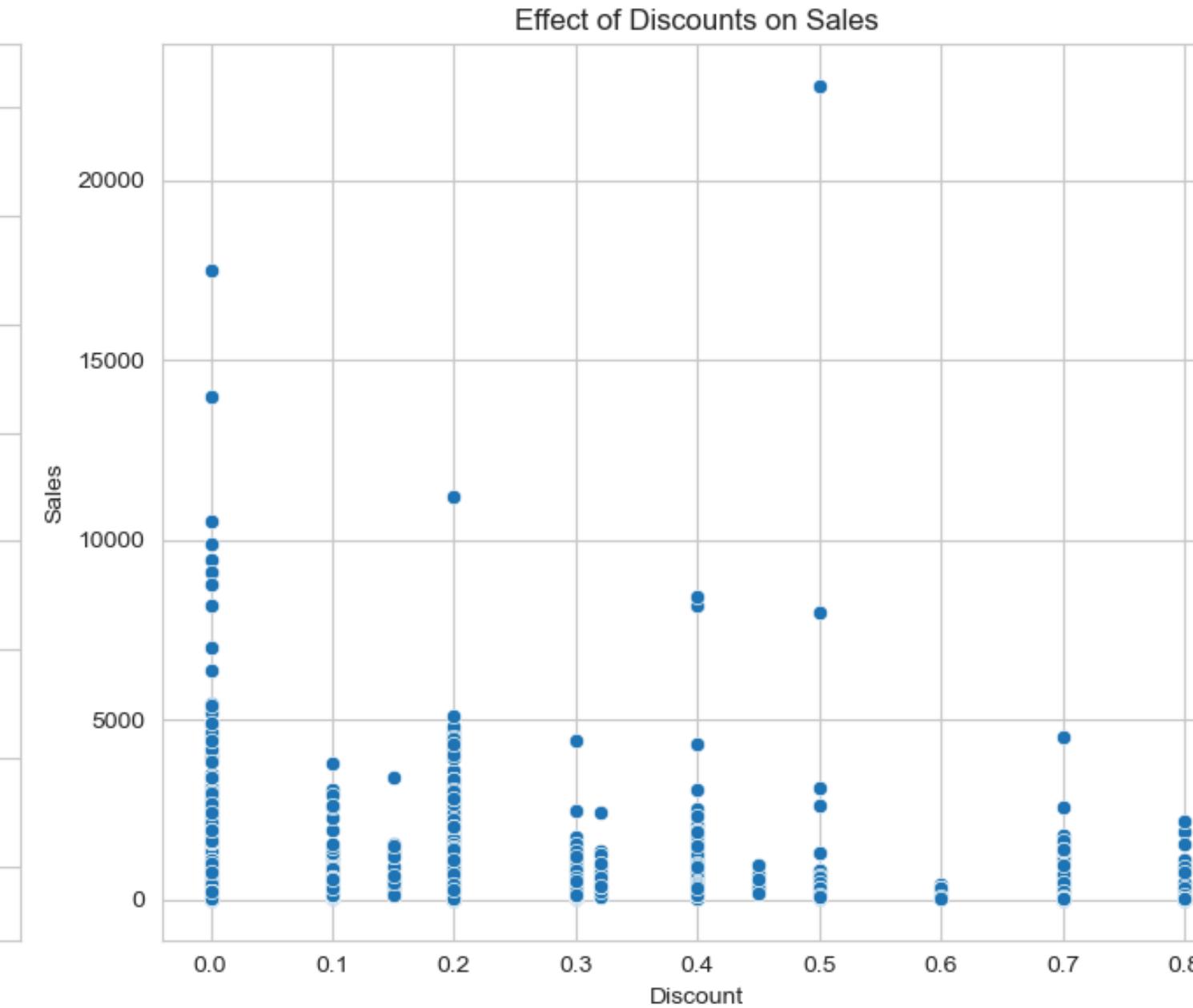
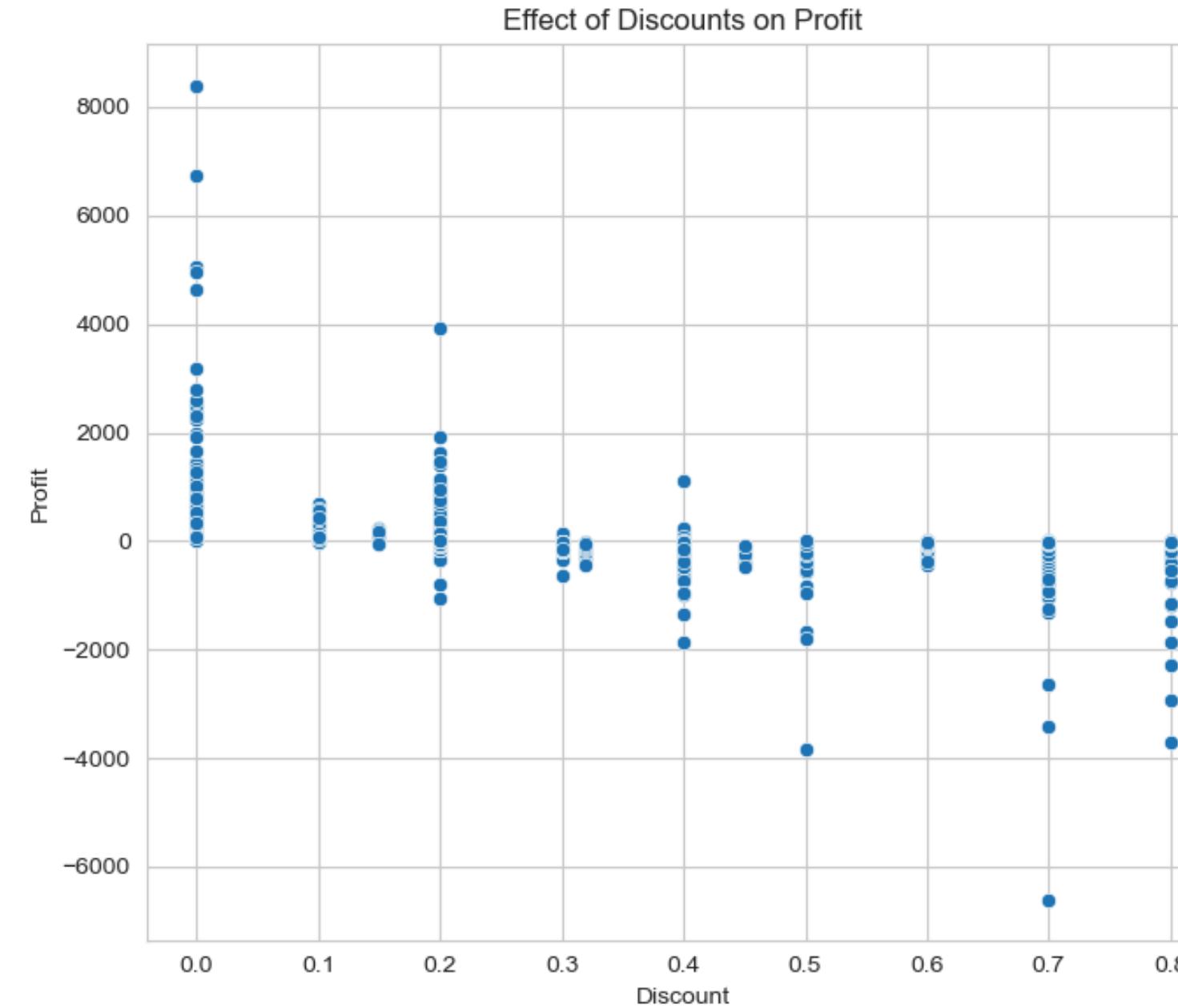
# Effect of Sales on Profit



# Total Sales and Profit Over Time



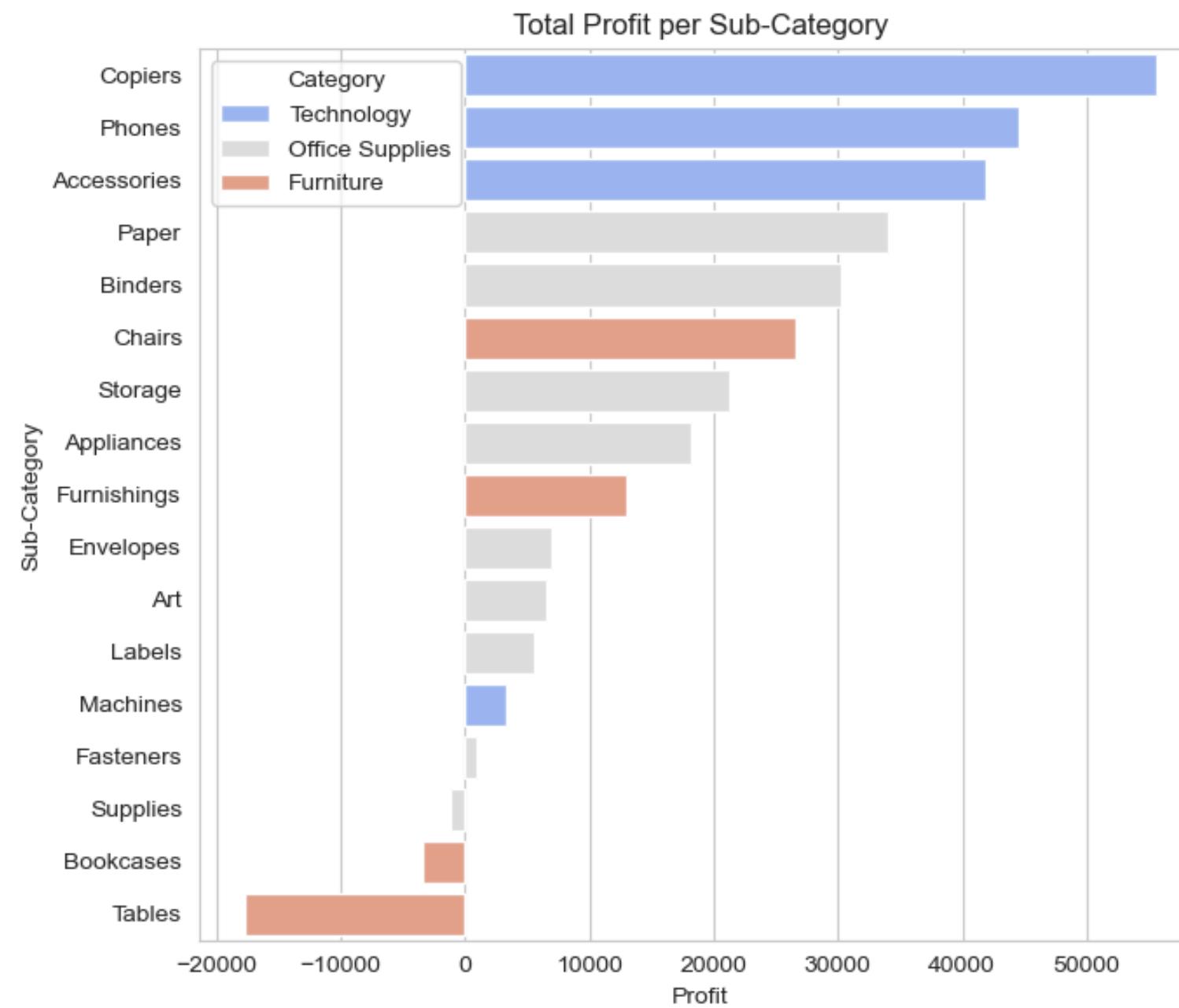
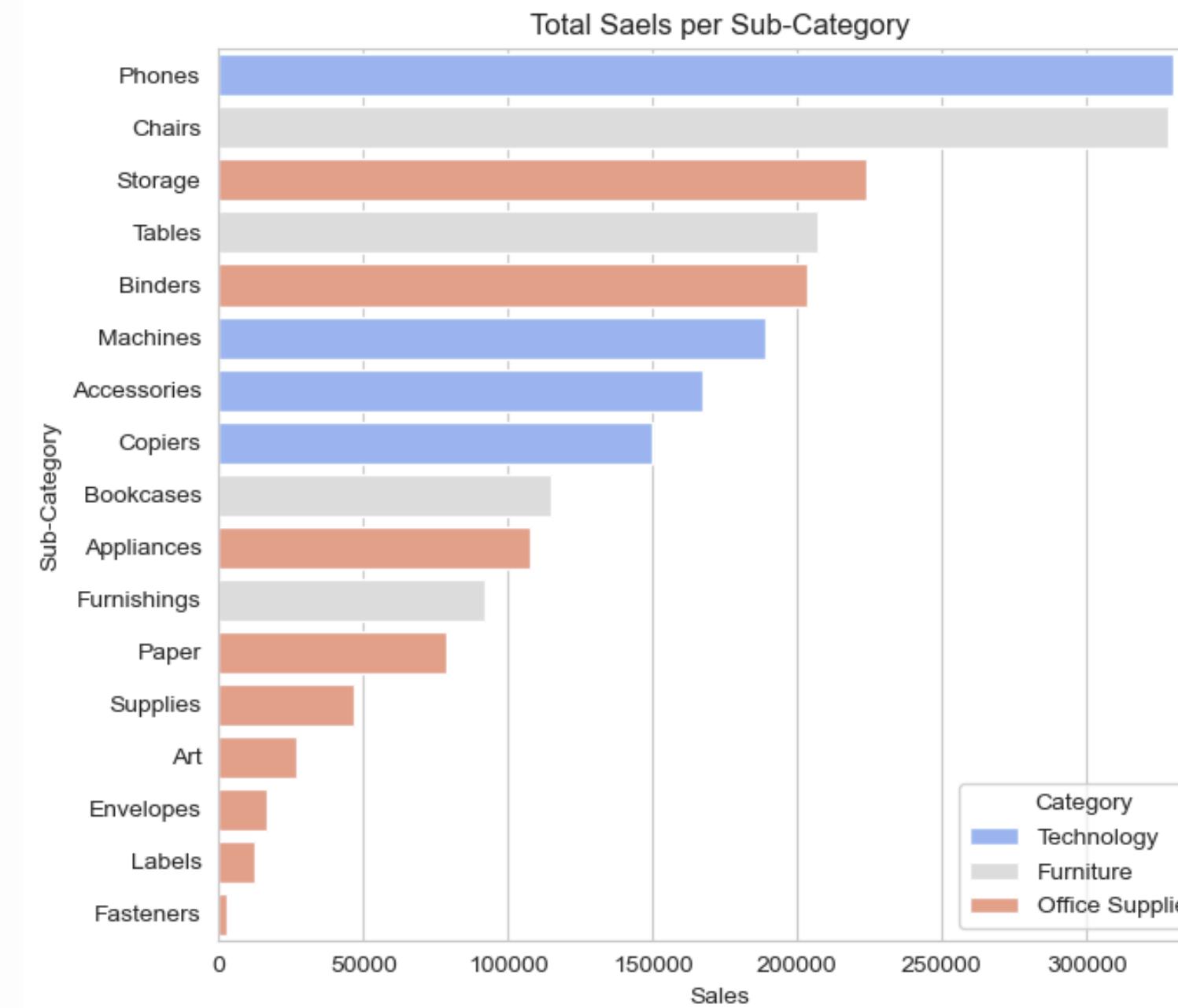
# Effect of Discounts on Profit and Sales



SAMSUNG

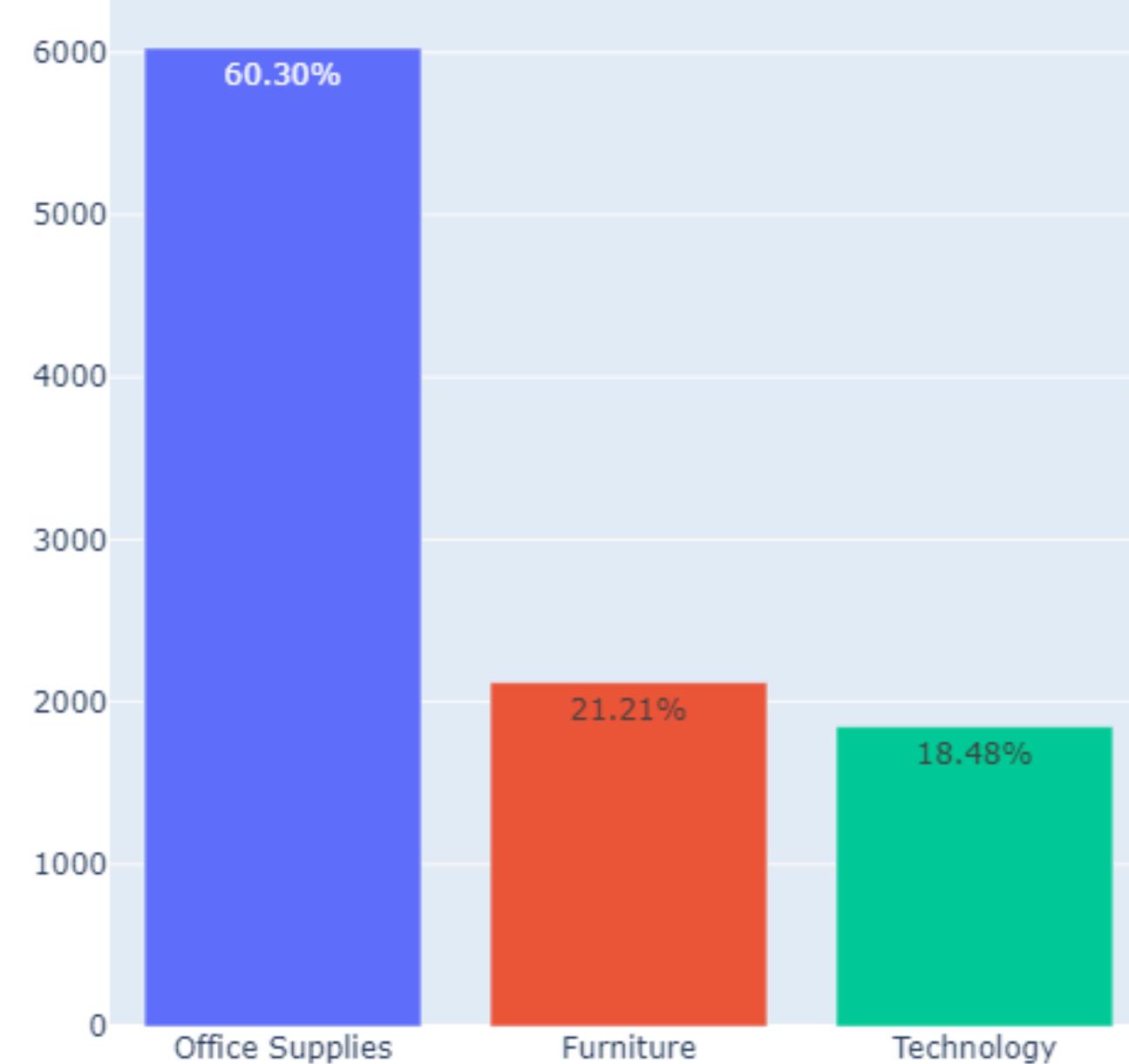
# PRODUCT CATEGORIES

# Product Categories

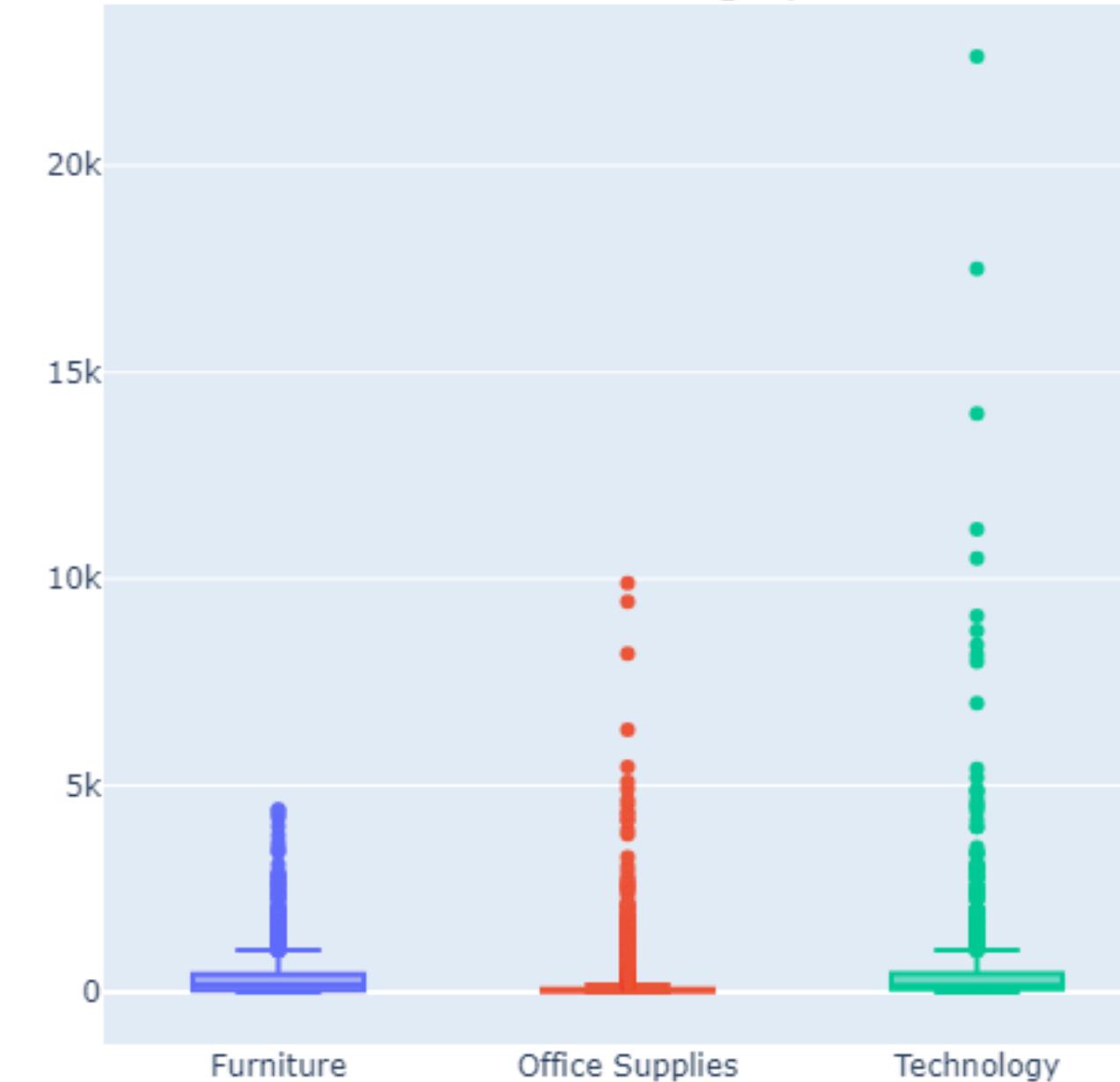


## Summary for Category

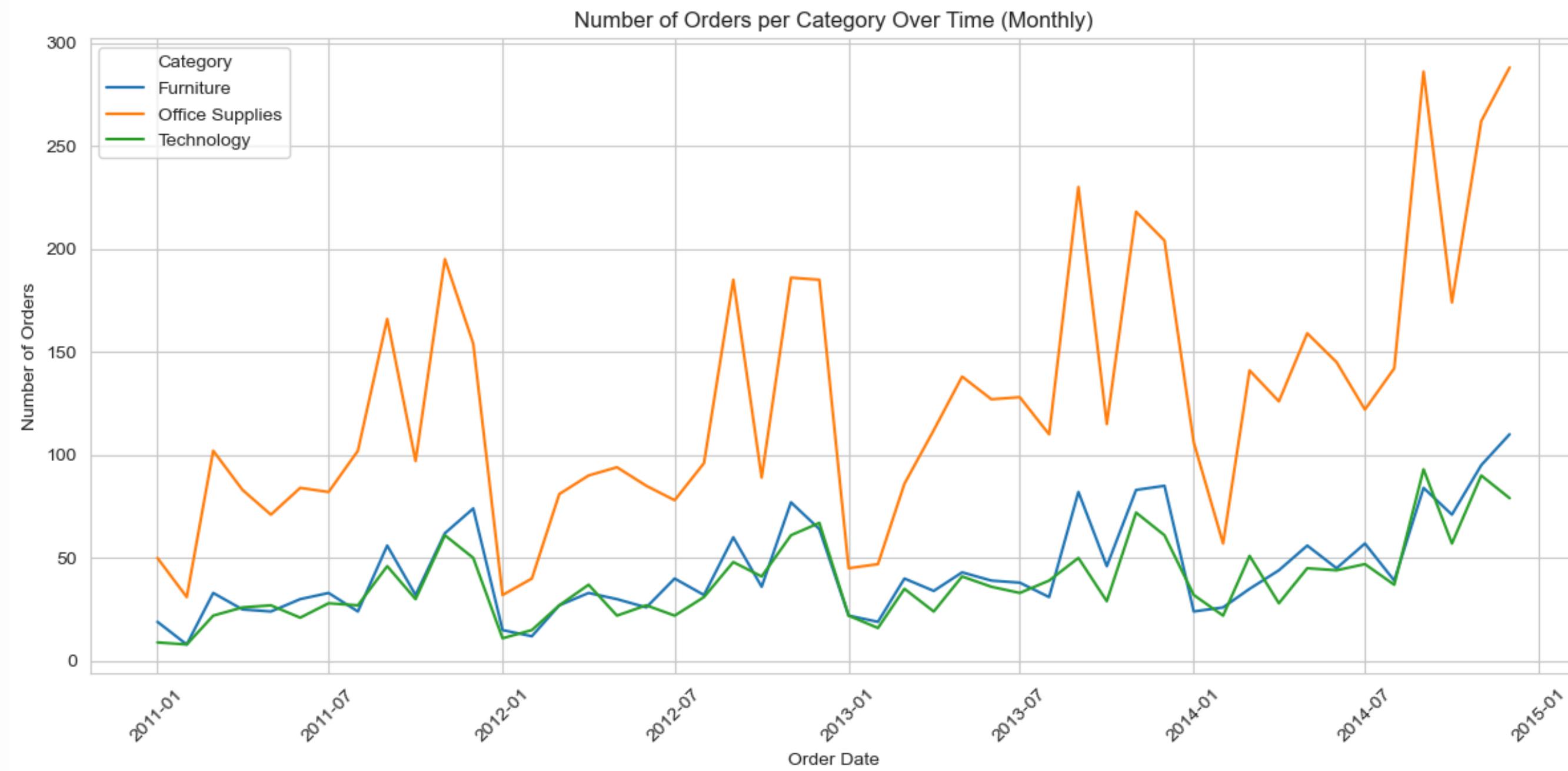
Category Frequency



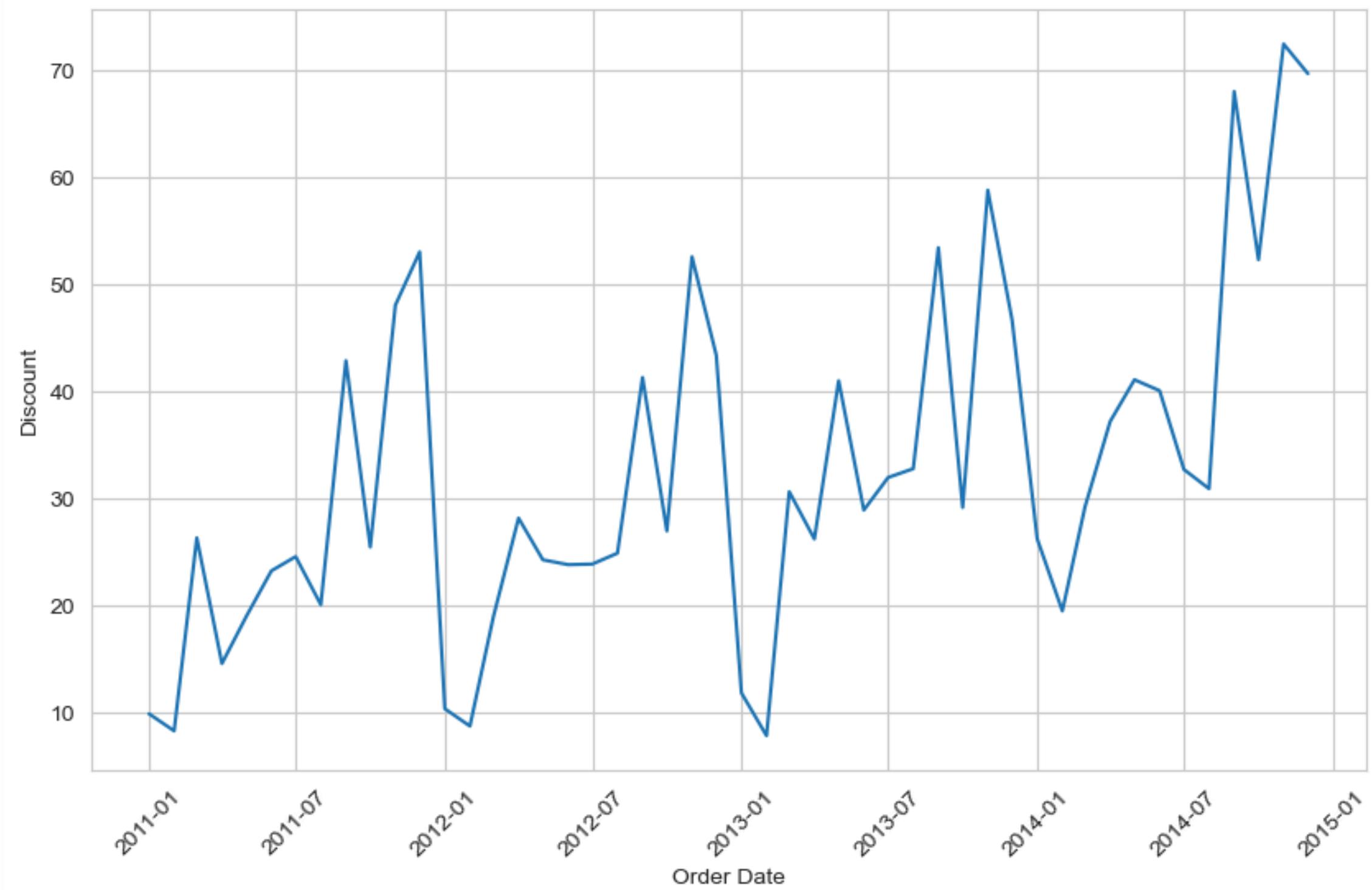
Sales vs Category



# Number of Orders per Category Over Time



Total Discount Over Time

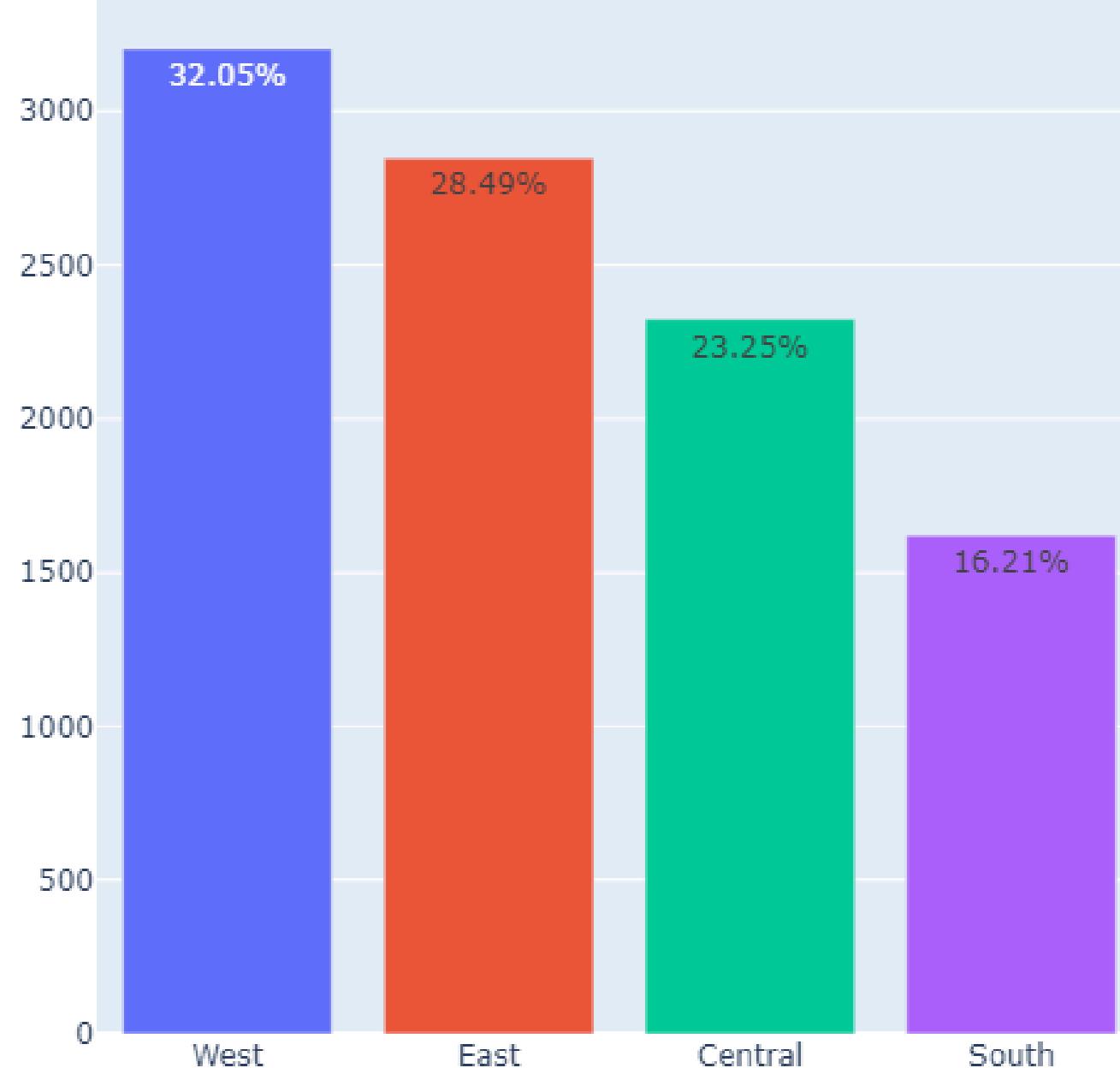


SAMSUNG

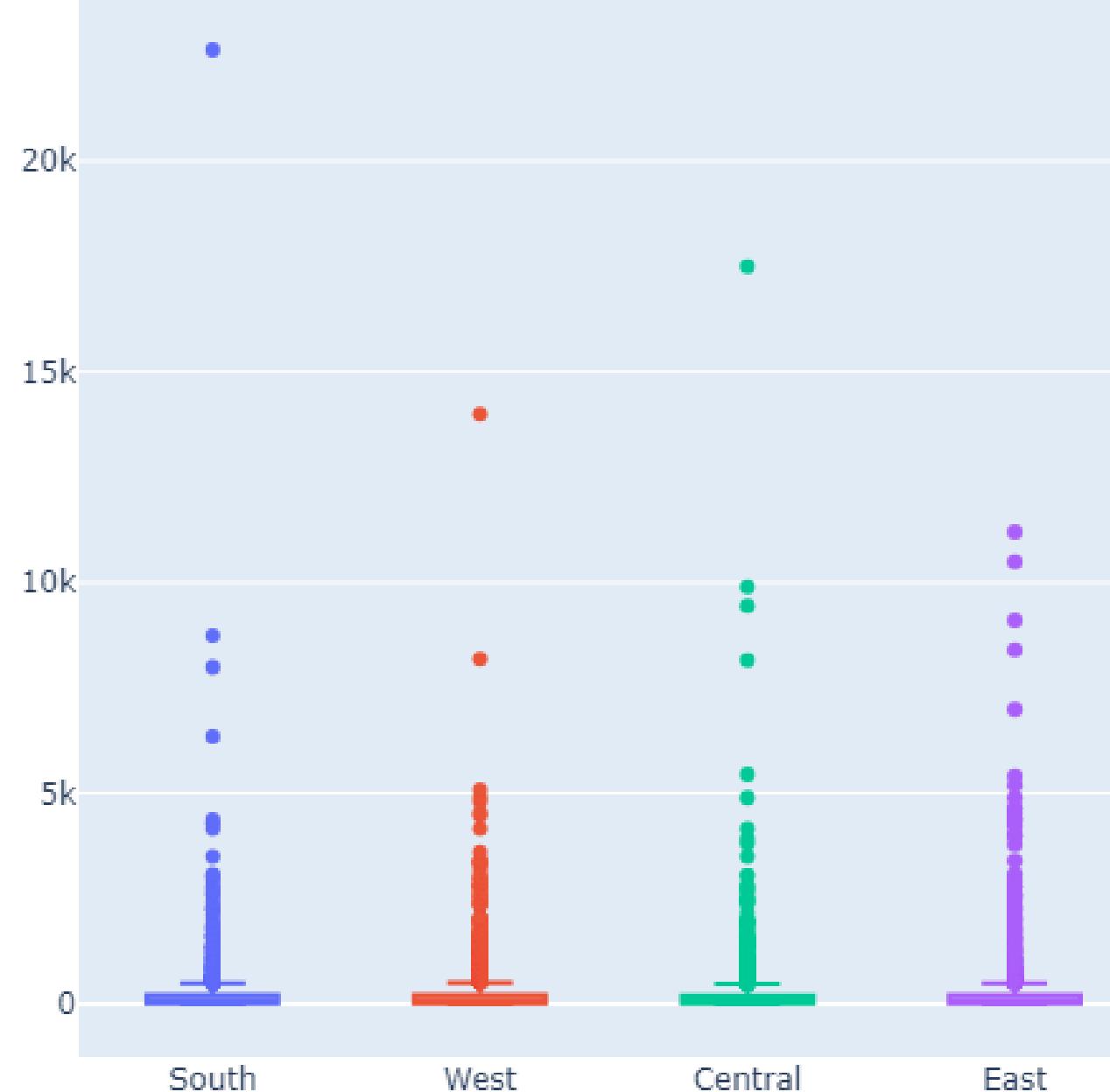
# GEOGRAPHICAL INSIGHTS

## Summary for Region

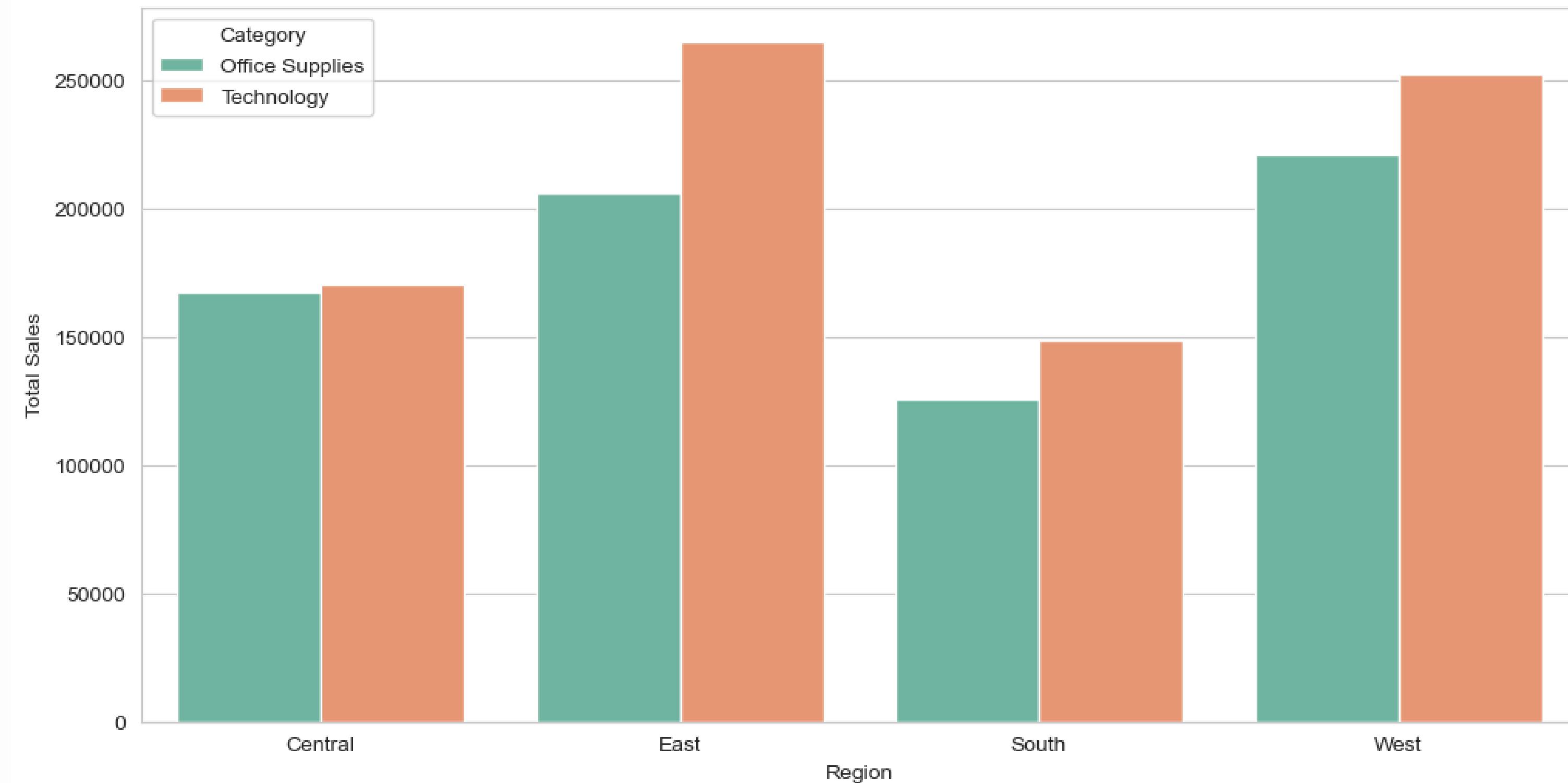
Region Frequency



Sales vs Region

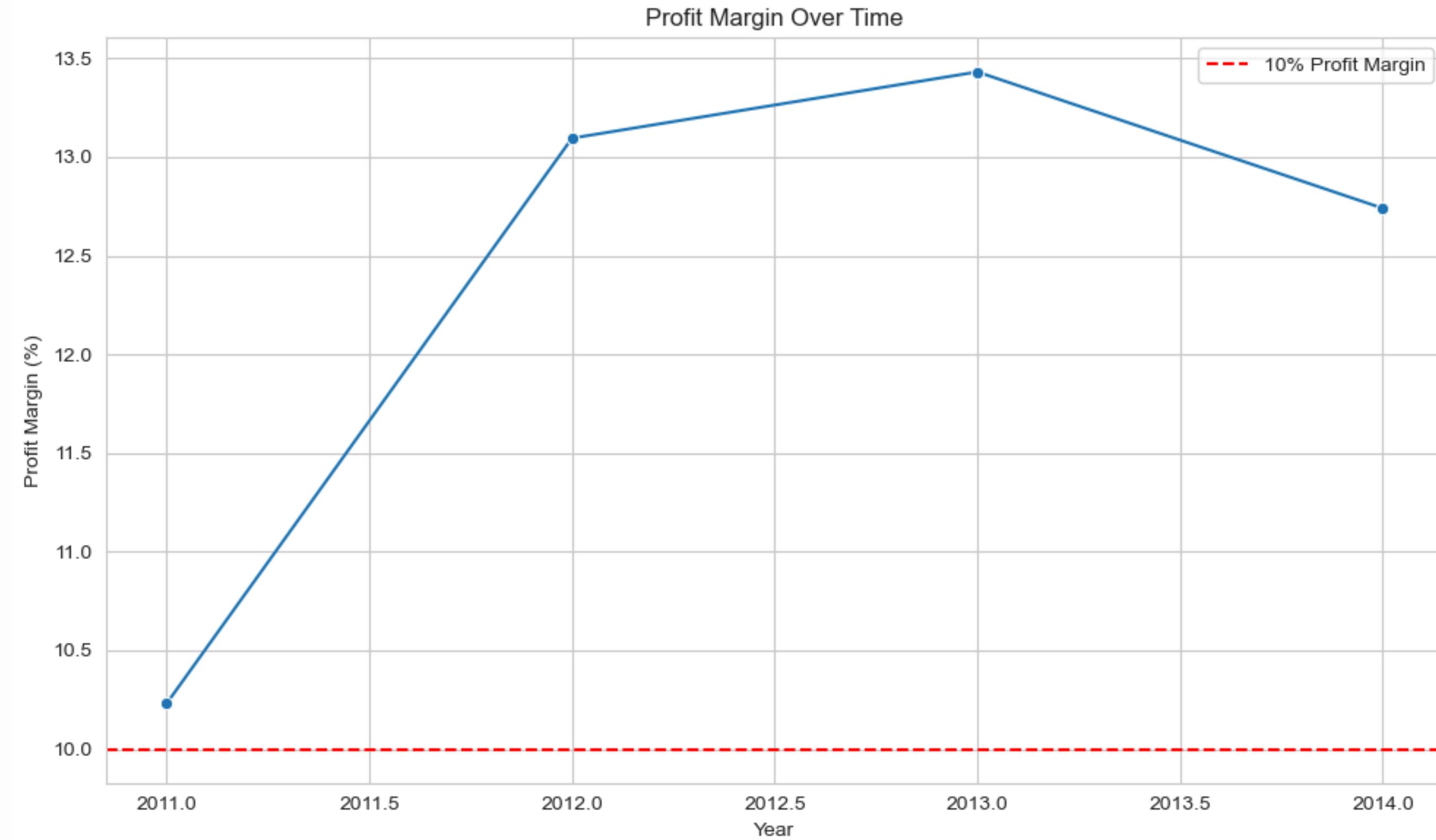


## Sales of Office Supplies and Technology in Different Regions

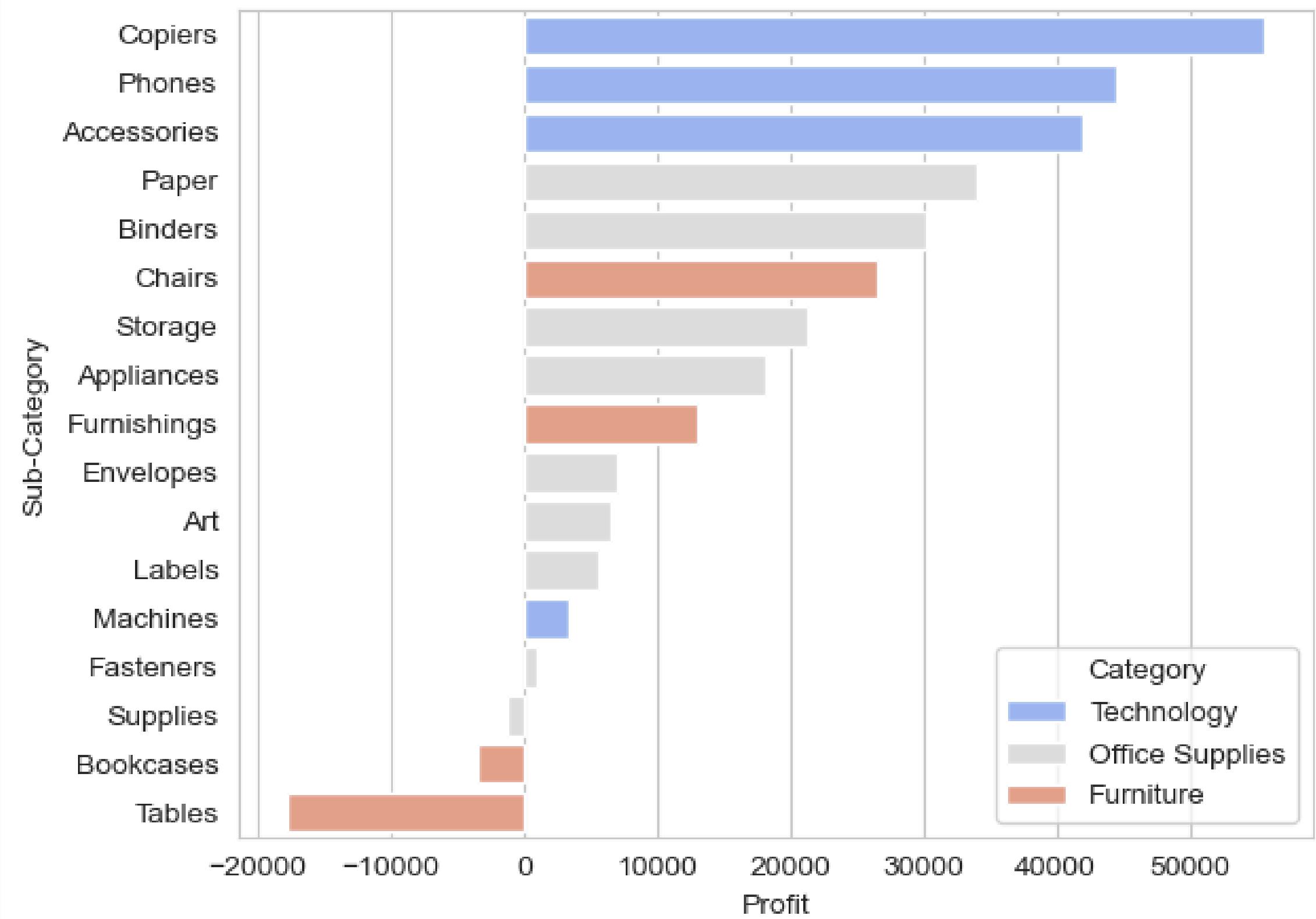


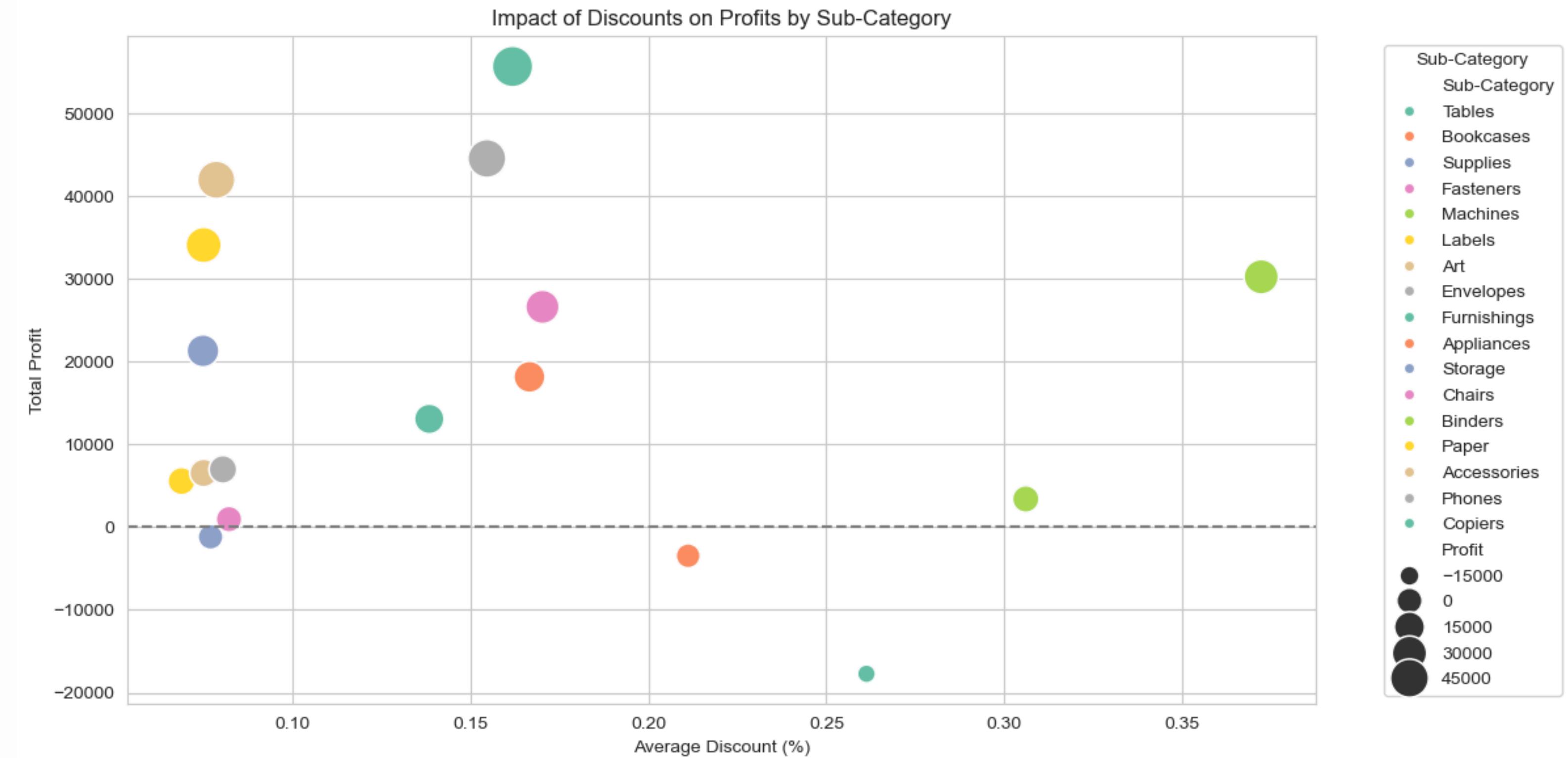
SAMSUNG

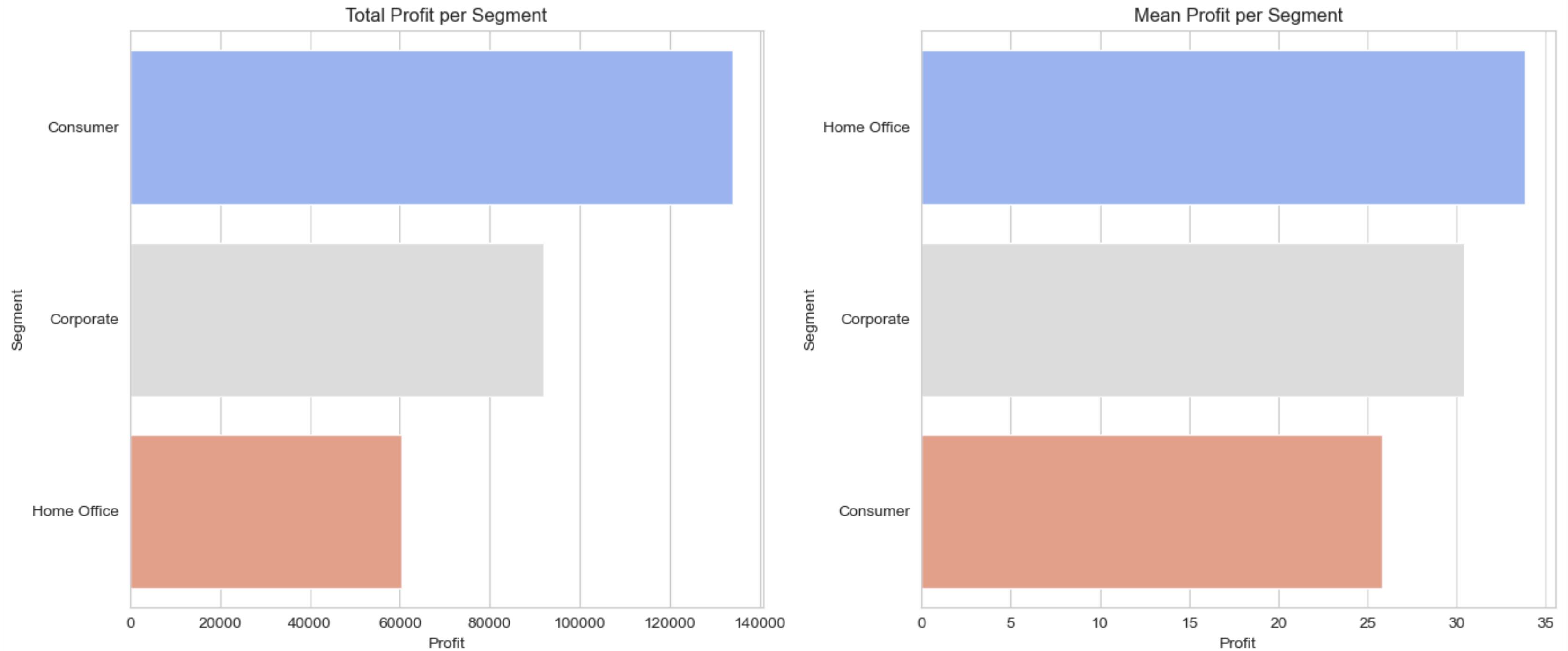
# PROFITABILITY



### Total Profit per Sub-Category







SAMSUNG

# BUSINESS INSGHITS



## Problem

- High Discount Dependency Leading to Low Profit Margins



## Solution

- Redefine Discount Strategies

Action: Implement data-driven discount strategies that vary by product category and customer segment.



## Problem

- Inconsistent Sales Across Product Categories



## Solution

- Boost Sales in Underperforming Product Categories

Action: Launch targeted marketing campaigns promoting underperforming categories like Furniture and Technology.



## Problem

- Uneven Regional Sales Distribution



## Solution

- Region-Specific Sales Strategies:

Action: Use sales data to identify regional demand patterns and preferences.



## Problem

- Limited Insight into Customer Behavior



## Solution

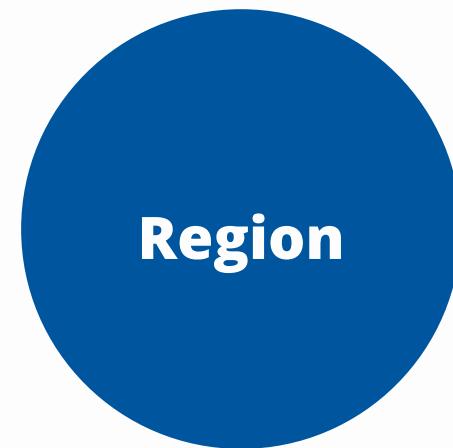
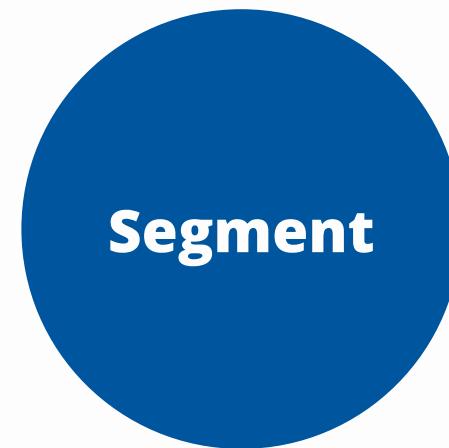
- Develop Customer Loyalty and Retention Programs:  
Action: Use purchase history to personalize offers and encourage long-term customer retention.



# Data Preprocessing

# Encoding

## Label encoder



# Feature Engineering

01

**Profit Margin**

02

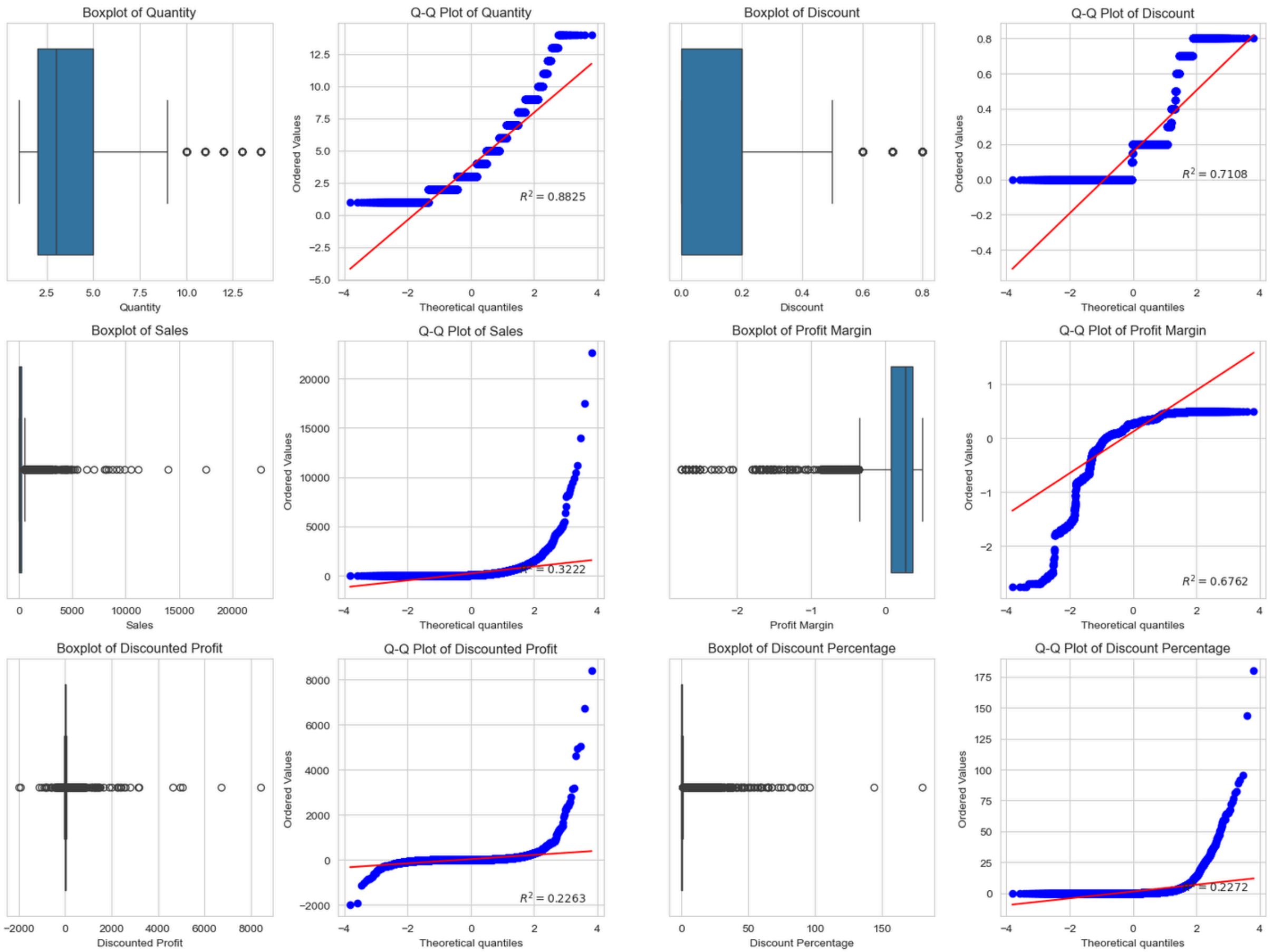
**Discounted Profit**

03

**Discount Percentage**

SAMSUNG

# SKEWNESS HANDLING



# Feature Transformation

01

**log\_transform**

Quantity  
Discounted Profit

02

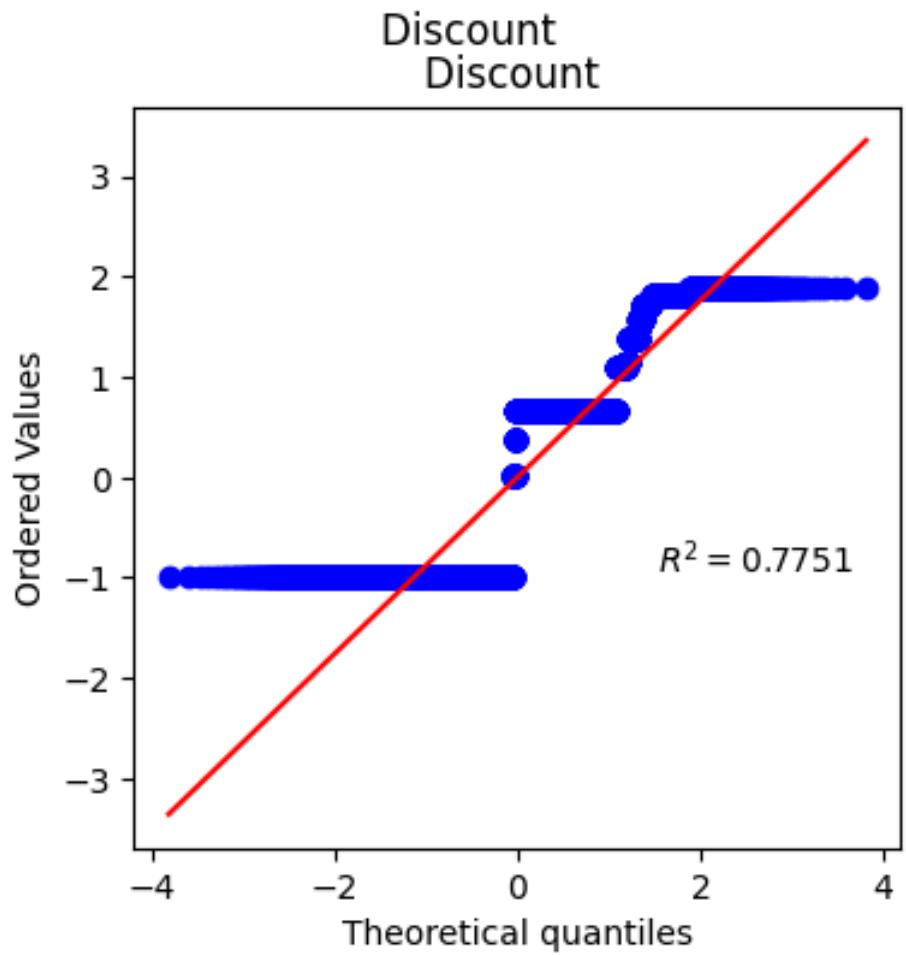
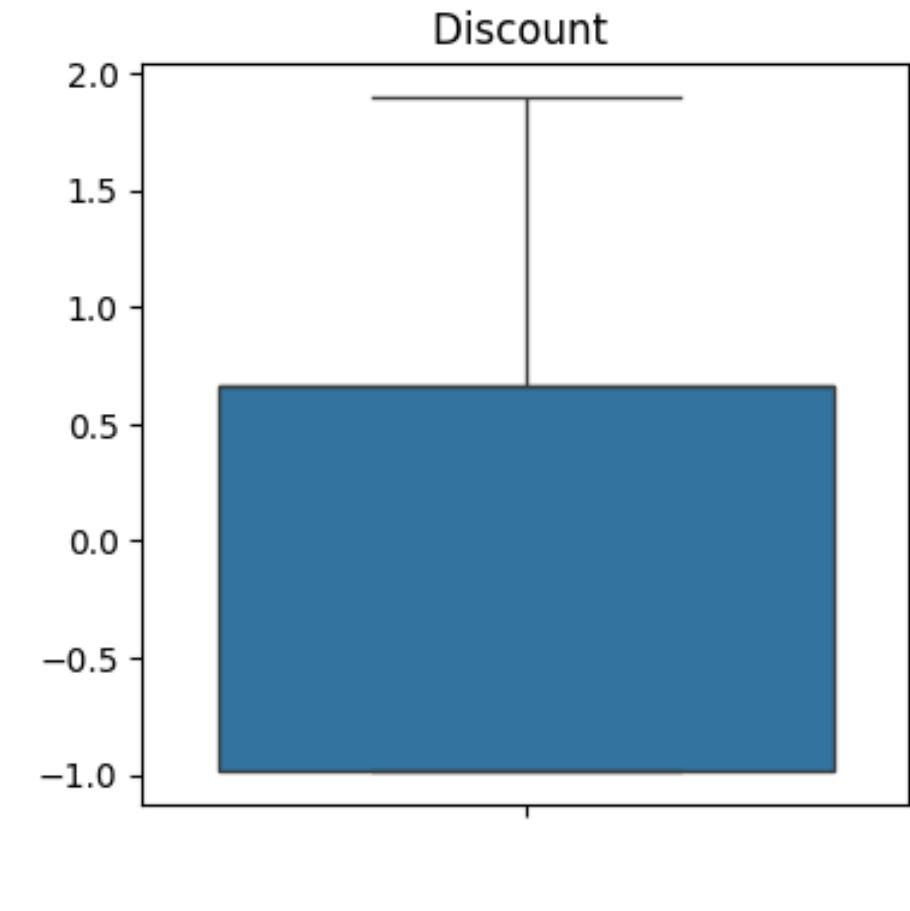
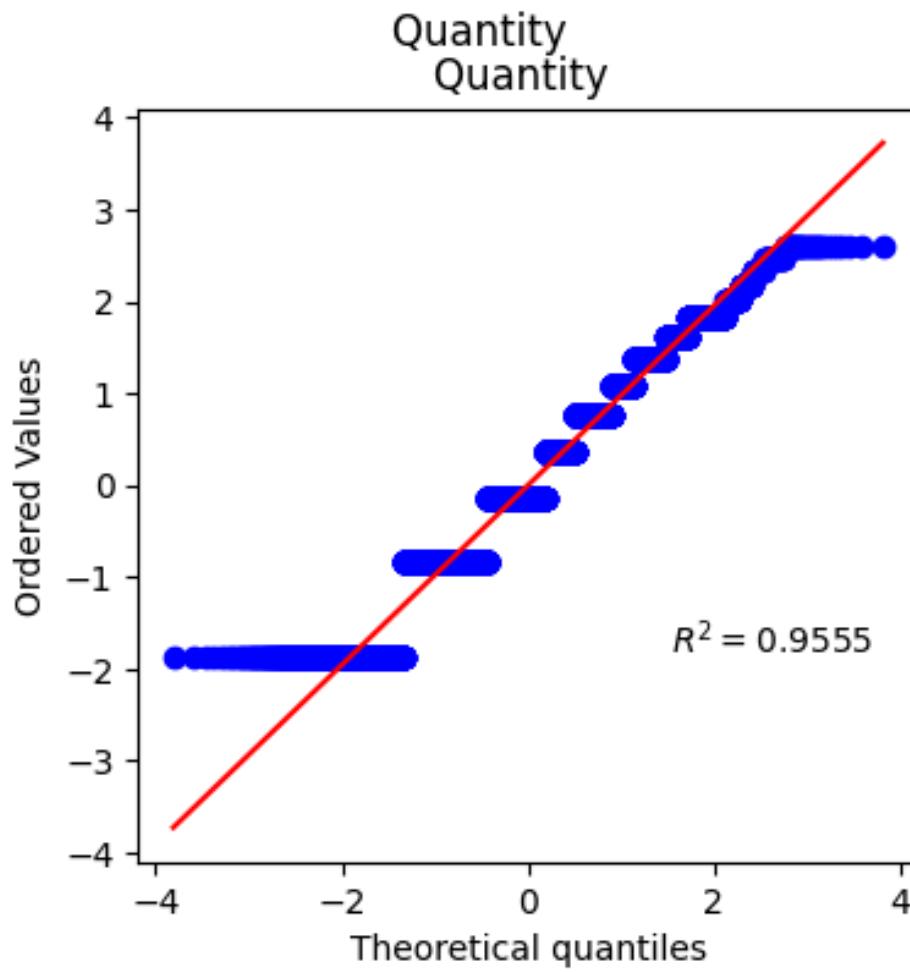
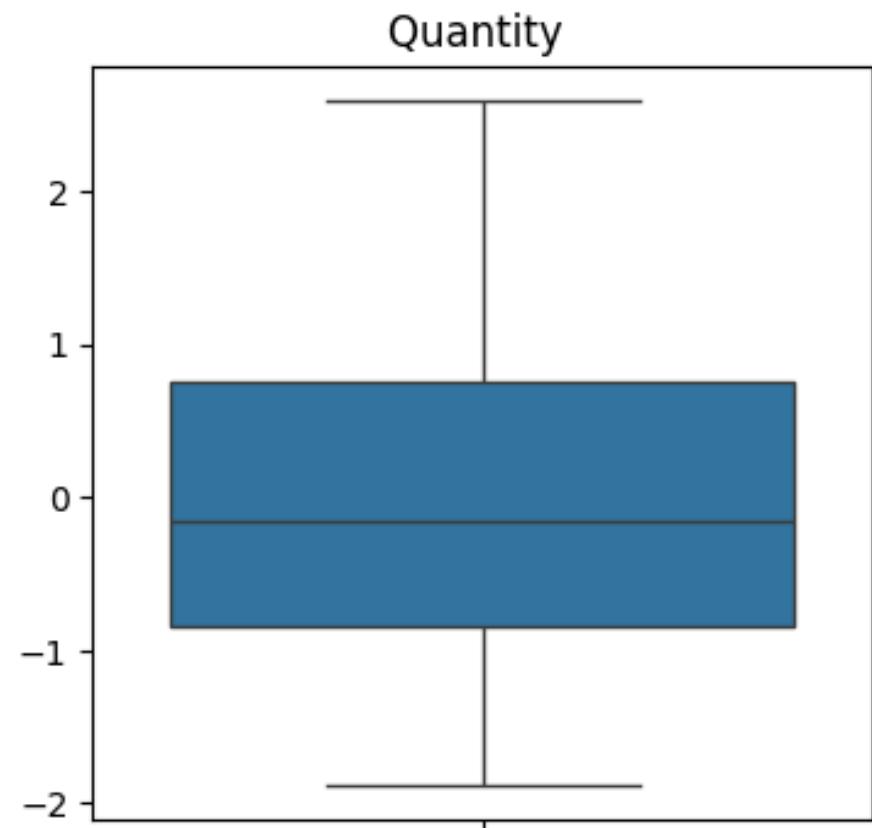
**power\_transform**

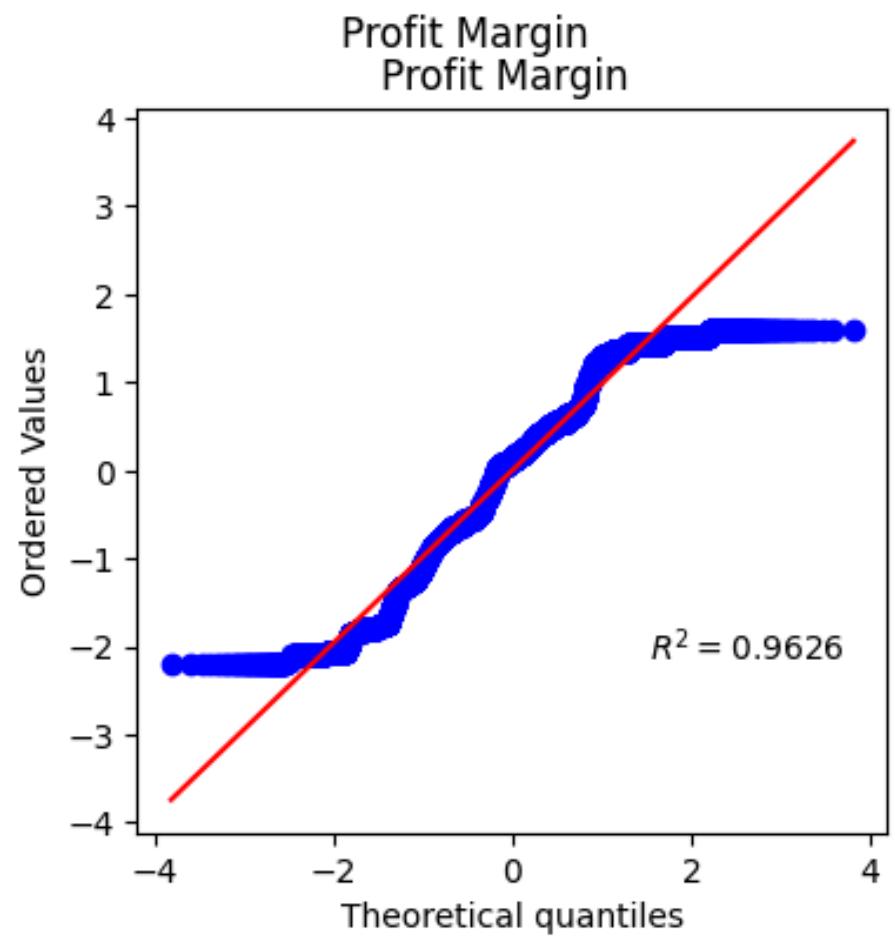
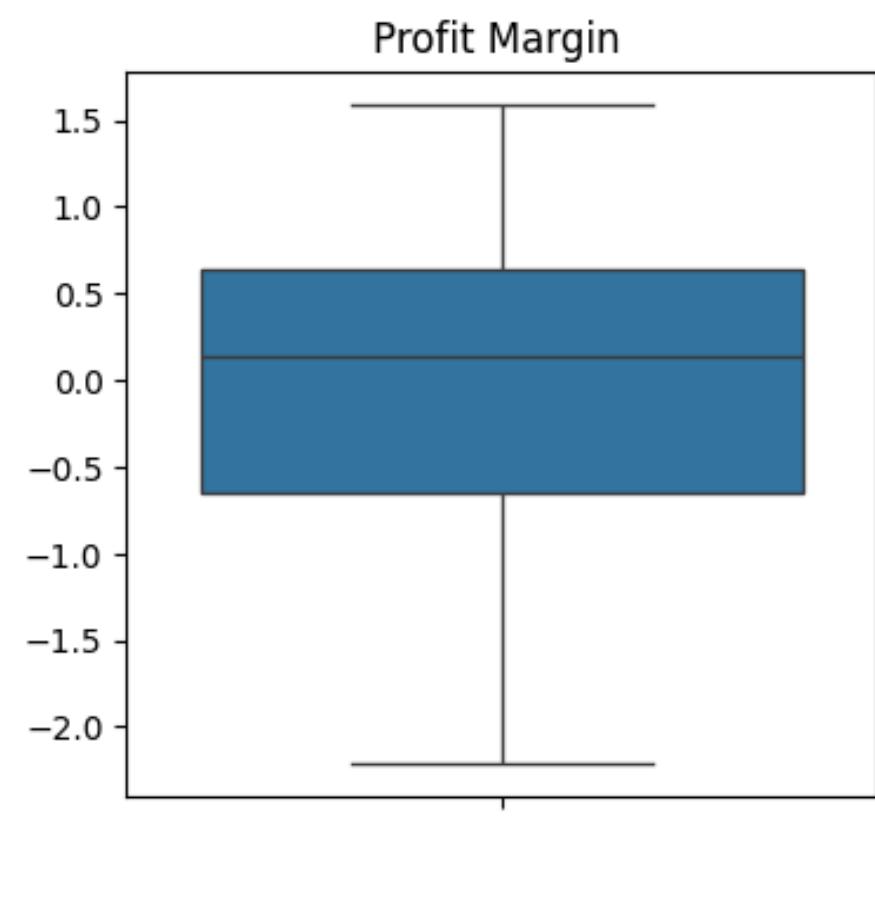
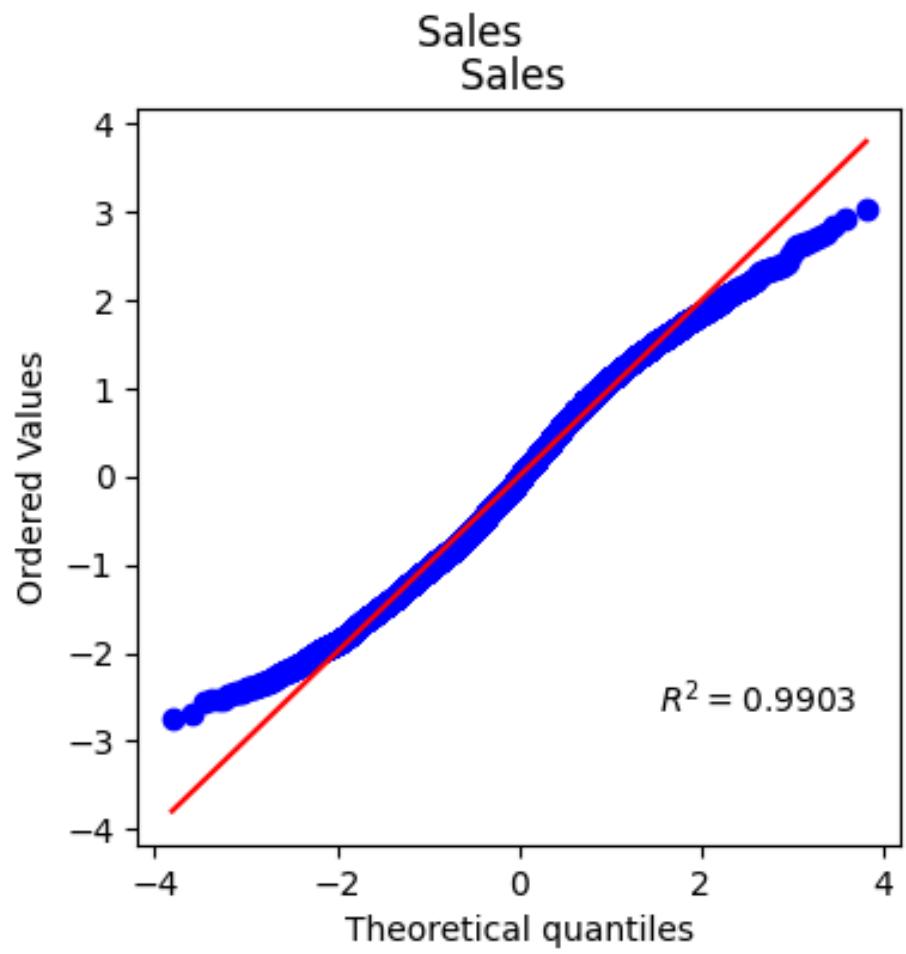
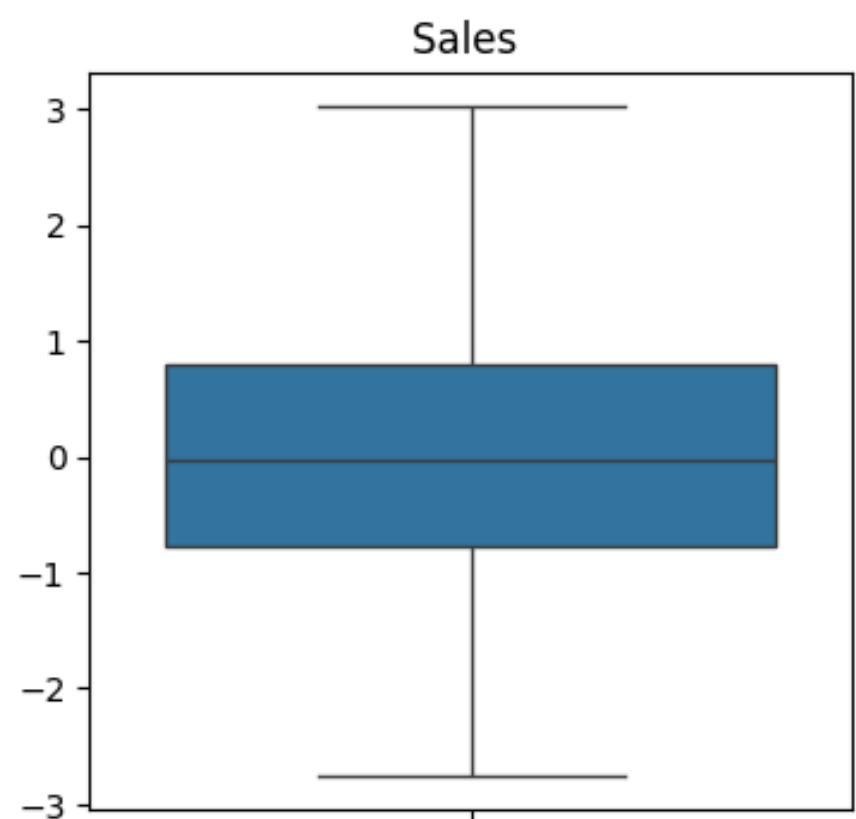
Profit Margin  
Sales  
Discount Percentage

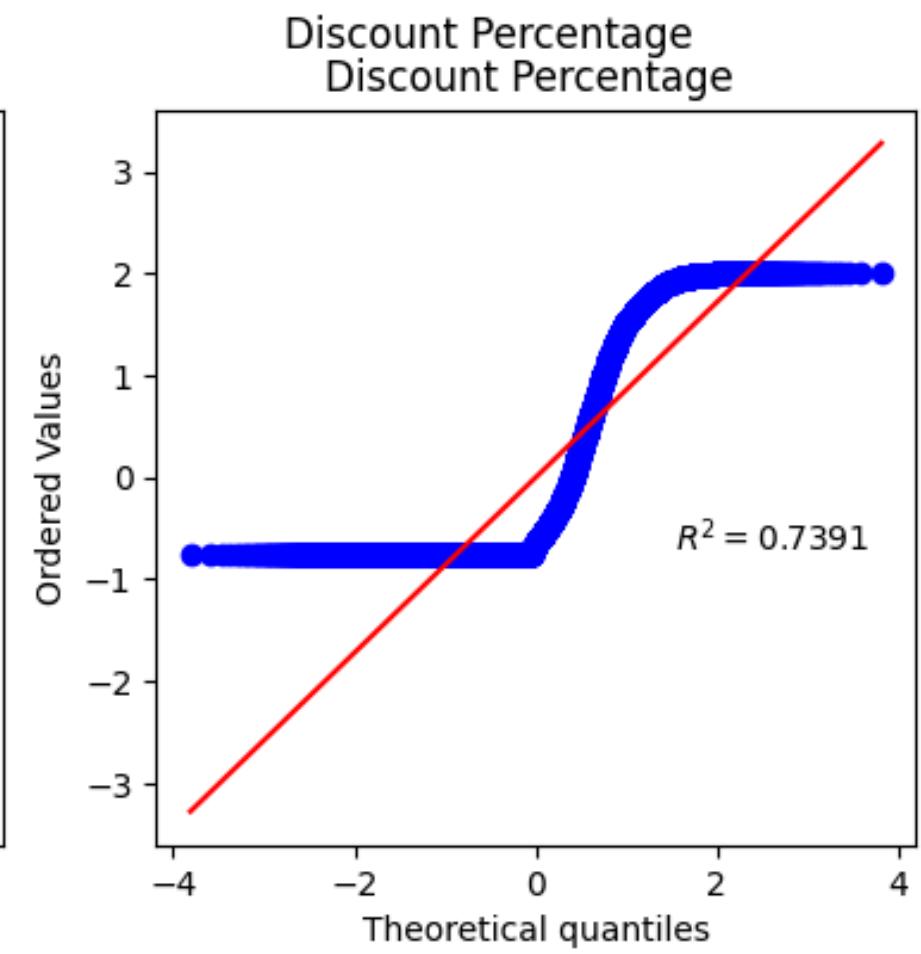
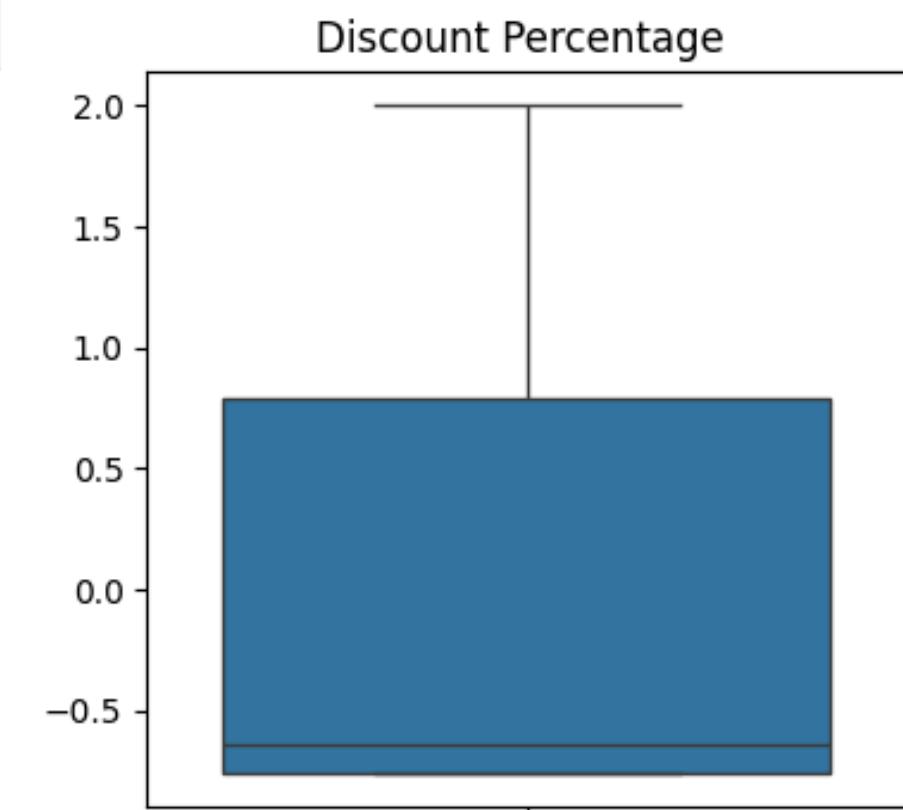
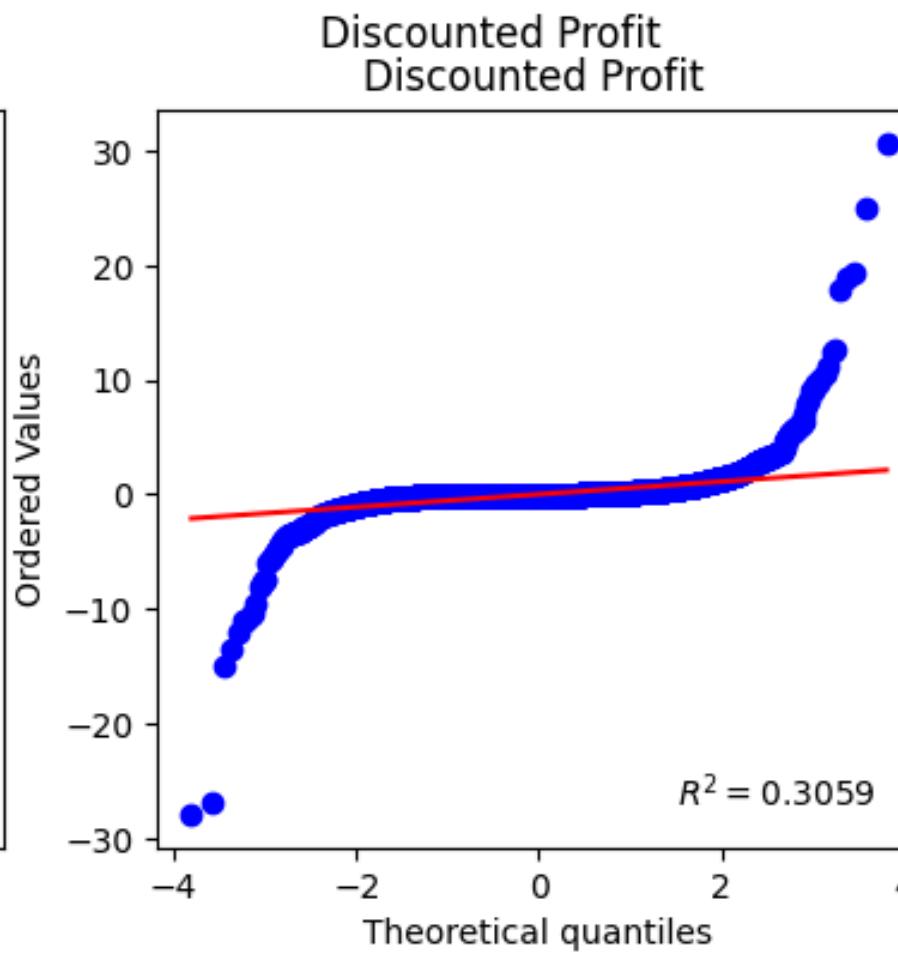
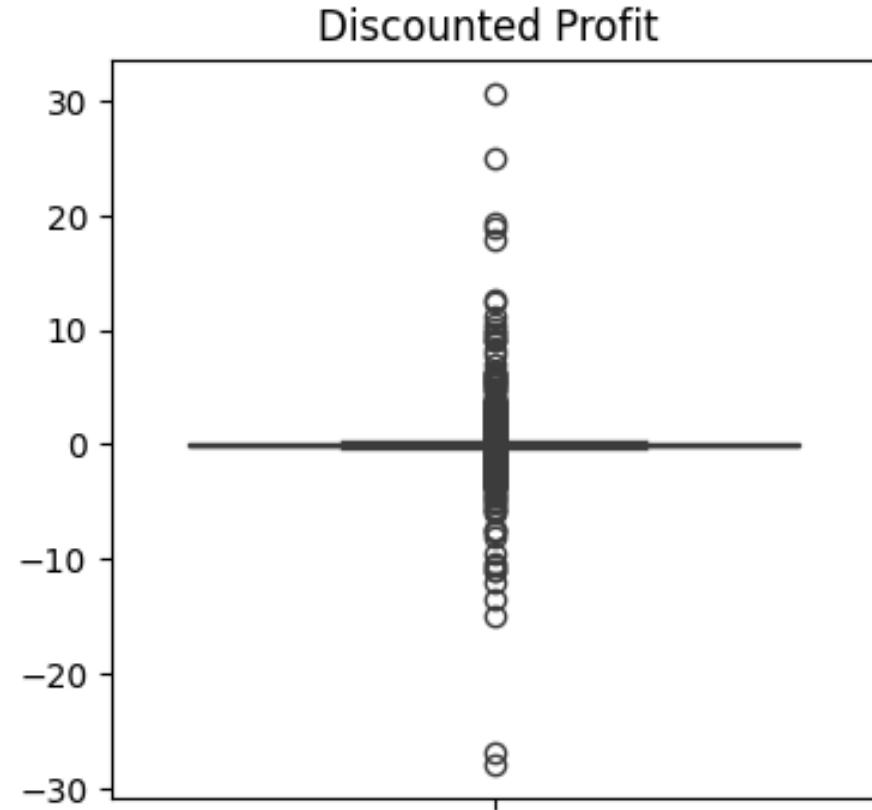
03

**sqrt\_transform**

Discount







# Model Overview



# Model Selection

We built multiple regression models to predict key metrics in sales performance.

The models tested include:

- K-Nearest Neighbors (KNN)
- Random Forest
- XGBoost
- AdaBoost

To ensure the accuracy and relevance of predictions, we used hyperparameter tuning and scaling techniques to optimize each model's performance.

# Scaling Techniques

Scaling techniques are crucial in machine learning, especially when models rely on distance measurements or when features vary in magnitude.

We tested three scaling methods:

- RobustScaler: Handles outliers by using the interquartile range.
- MinMaxScaler: Scales values between 0 and 1.
- StandardScaler: Scales data to have zero mean and unit variance.

The performance of models varies based on how the input features are scaled.

# Hyperparameter Tuning

Hyperparameter tuning was done using GridSearchCV, which tests different parameter combinations to find the most effective configuration.

Key hyperparameters:

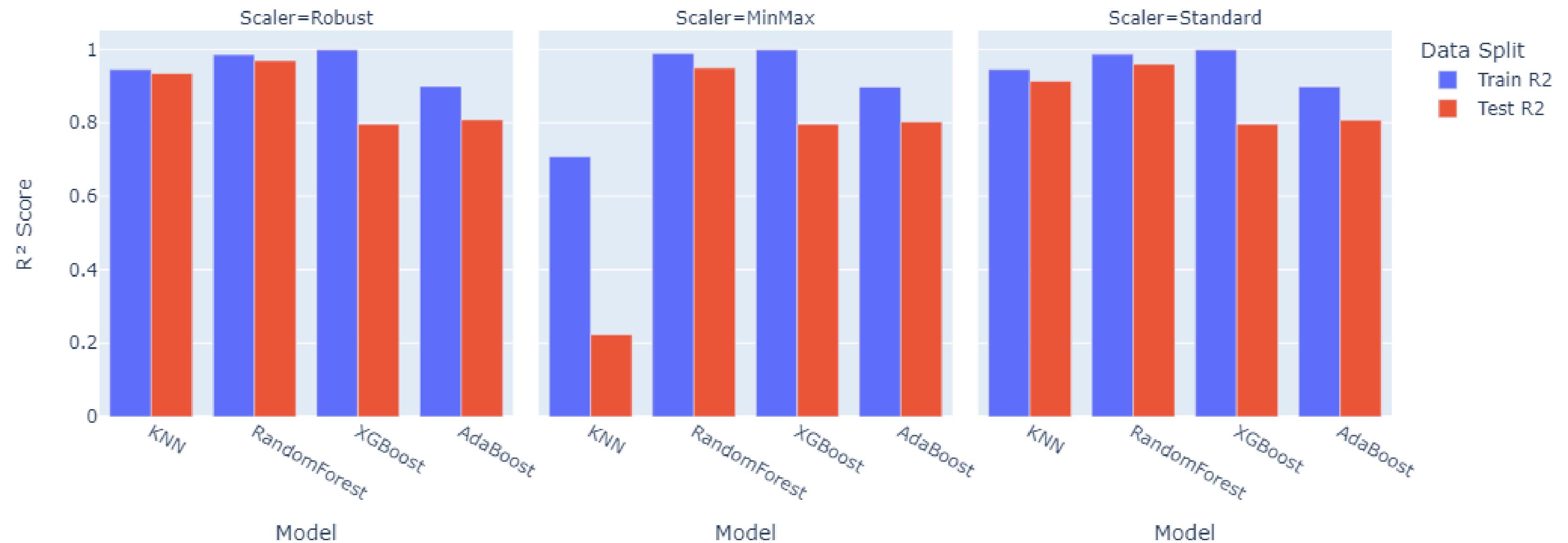
- KNN: n\_neighbors (3, 5, 7)
- RandomForest: n\_estimators, max\_depth
- XGBoost: n\_estimators, max\_depth, learning\_rate
- AdaBoost: n\_estimators, learning\_rate

Tuning improves model accuracy, making predictions more reliable for decision-making.

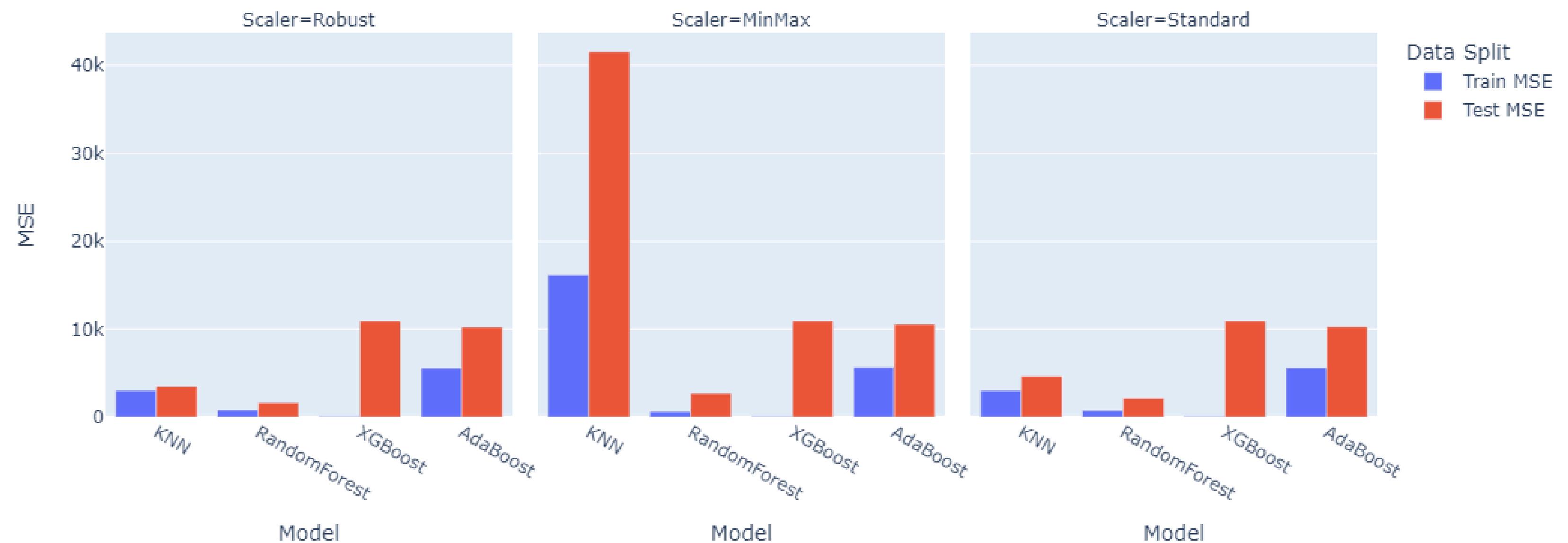
# MODEL EVALUATION

	<b>Scaler</b>	<b>Model</b>	<b>Train R2</b>	<b>Test R2</b>	<b>Train MSE</b>	<b>Test MSE</b>	<b>Train RMSE</b>	<b>Test RMSE</b>	<b>Train MAE</b>	<b>Test MAE</b>
0	Robust	KNN	0.944956	0.934284	3052.301997	3516.813055	55.247642	59.302724	6.040234	9.144991
1	MinMax	KNN	0.707839	0.223556	16200.985060	41551.509818	127.283090	203.841875	34.181396	57.830215
2	Standard	KNN	0.945054	0.912716	3046.896187	4671.023586	55.198697	68.344887	15.868940	24.689630
3	Robust	RandomForest	0.984758	0.968853	845.192388	1666.824503	29.072193	40.826762	1.559644	3.945505
4	MinMax	RandomForest	0.987976	0.948960	666.759451	2731.406999	25.821686	52.262864	1.427619	4.407531
5	Standard	RandomForest	0.985799	0.959706	787.455433	2156.334119	28.061636	46.436345	2.806085	5.368107
6	Robust	XGBoost	0.998386	0.795704	89.478879	10932.943725	9.459328	104.560718	3.469071	11.984223
7	MinMax	XGBoost	0.998386	0.795704	89.478879	10932.943725	9.459328	104.560718	3.469071	11.984223
8	Standard	XGBoost	0.998386	0.795704	89.478879	10932.943725	9.459328	104.560718	3.469071	11.984223
9	Robust	AdaBoost	0.898779	0.808420	5612.952096	10252.450796	74.919638	101.254387	41.215511	44.825566
10	MinMax	AdaBoost	0.897264	0.802070	5696.924893	10592.269689	75.477976	102.918753	41.451512	45.405819
11	Standard	AdaBoost	0.898097	0.807508	5650.727679	10301.236380	75.171322	101.495007	39.797622	43.257175

## Train vs Test R<sup>2</sup> Comparison for Models with Different Scalers



## Train vs Test MSE Comparison for Models with Different Scalers



# Best Models



## Random Forest with Robust Scaler

**Train R2: 98.5%**

**Test R2: 97%**

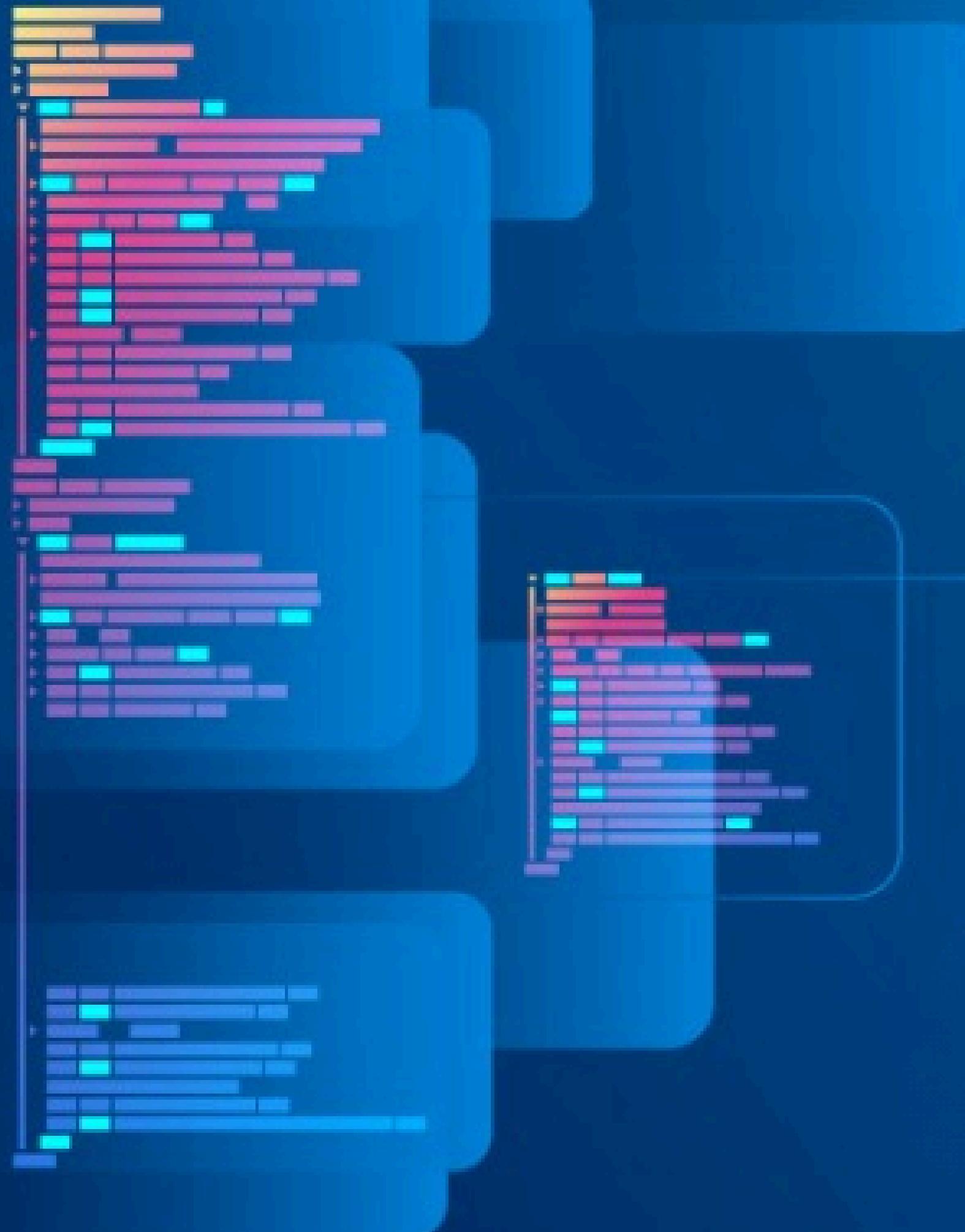


## KNN with Robust Scaler

**Train R2: 94.5%**

**Test R2: 93.5%**

# Deployment



SAMSUNG



# Website



Scan me!

The image shows a screenshot of a mobile application interface. At the top, the word "Upload" is displayed in green. Below it is a QR code with the word "ATT" and a paperclip icon above it. A blue button labeled "scan" is located at the bottom of the QR code area. At the very bottom of the screen, a message states "The model results accuracy is 96.88%." The background of the app is light blue.



# Dashboard



**SAMSUNG**

**THANK YOU!**



Together for Tomorrow!  
**Enabling People**

Education for Future Generations

©2020 SAMSUNG. All rights reserved.

Samsung Electronics Corporate Citizenship Office holds the copyright of book.

This book is a literary property protected by copyright law so reprint and reproduction without permission are prohibited.

To use this book other than the curriculum of Samsung innovation Campus or to use the entire or part of this book, you must receive written consent from copyright holder.