# Data Analysis Superstore Project

# Content

- Introduction

- Data Cleaning and Normalization

- Questions and KPIs metrics

- Create charts and visualizations

- Creating Dashboard

Edit with WPS Office

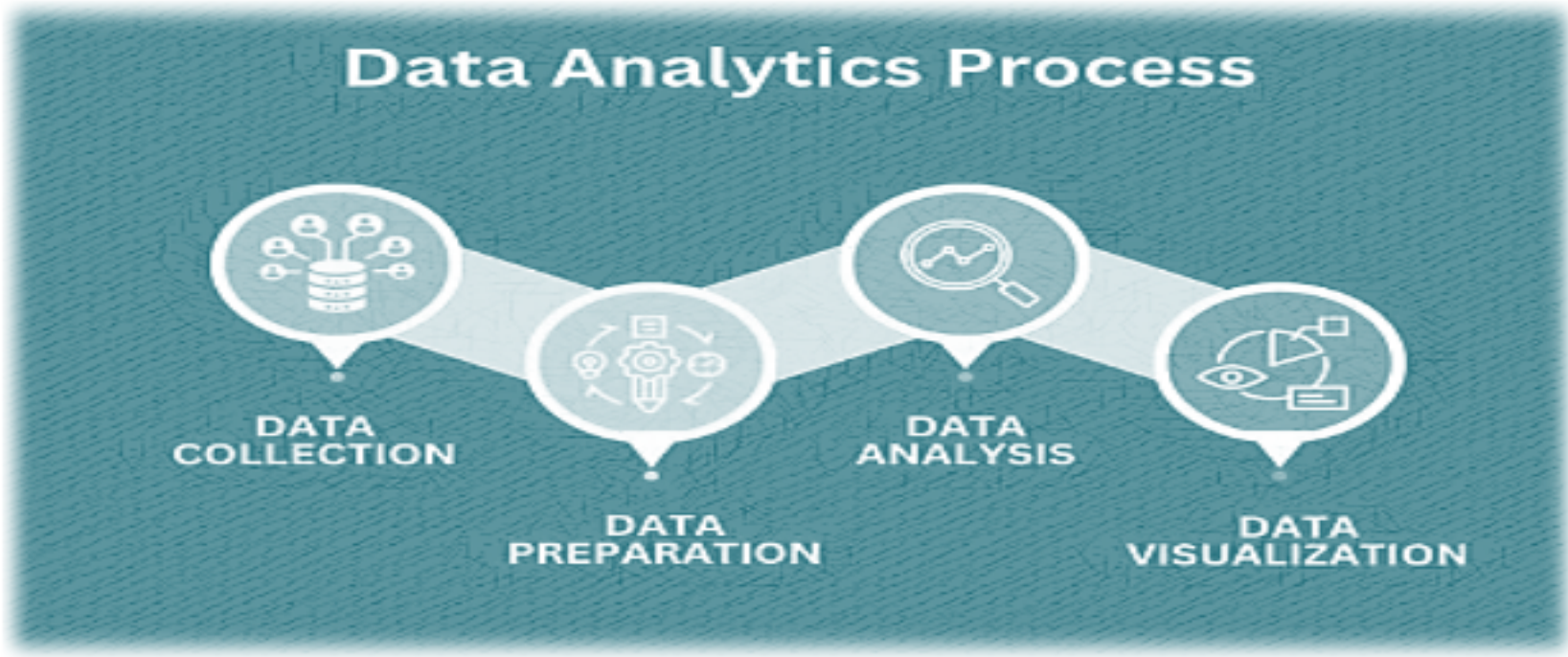Team members: Almoataz M. Gad | Abddelrahman Ali | Ahmed Adel |
Ahmed M. Elsafy

# ➤ Introduction

- In today's competitive retail environment, data-driven decision-making is essential for optimizing sales performance and enhancing customer satisfaction. This project focuses on analyzing the sales data of a superstore to uncover key insights and trends that can drive strategic improvements. By leveraging tools such as Excel, SQL Server, and Tableau, we have examined various factors like product categories, customer demographics and sales performance. The goal is to identify patterns and opportunities that can inform decisions on inventory management, pricing strategies, and customer targeting to increase profitability and market reach.

# ➤ Introduction

# ➤ Data Cleaning and Normalization

- Data Cleaning is the first step in the project. And in my team's opinion data cleaning process it the most important step.

- Cleaned data is important to build high quality insights.

- We used python programming language to achieve cleaned data.

- We used python libraries like : Pandas ,Numpy and matplotlip.

- Next screen show data cleaning process by python.

# ➤ Data Cleaning and Normalization

Data cleaning step can be done by many tools, our team used python programming language. By using libraries like pandas and numpy as shown.

# ➤ Data Cleaning and Normalization

# ➤ Data Cleaning and Normalization

Questions and KPIs metrics

# Questions and KPIs metrics

## 1. Customer Insights
• How many unique customers are there?

• What is the distribution of customers by segment (Customer, Corporate, Home Office)?

• Which segment generates the highest sales?

• How many orders does each customer place on average?

**Charts:**

• Bar chart: Number of customers by segment.

• Pie chart: Sales contribution by customer segment.

• Bar chart: Number of orders per customer.

## 2. Order Insights
• How many orders were placed, and what is the average order amount?

• What is the average time between the order date and the ship date?

• Which customers place the highest number of orders?

**Charts:**

• Line chart: Number of orders over time.

• Bar chart: Average ship time by customer.

• Line chart: Number of orders by month.

## 3. Product Insights
• Which products are the most and least sold by quantity?

• What are the best-selling product categories and sub-categories?

• How does product sales vary across different regions?

**Charts:**

• Bar chart: Sales by product category.

• Pie chart: Sales by sub-category.

• Heatmap: Product sales by region.

## 4. Sales and Revenue Insights
• What is the total sales amount, and how does it distribute across customer segments and regions?

• Which regions contribute the most to sales?

• What is the trend of sales over time?

**Charts:**

• Line chart: Total sales over time.

• Stacked bar chart: Sales by region.

• Pie chart: Sales distribution by region.

## 5. Geographical Insights
• How are sales distributed across different regions and cities?

• Which postal codes are the most active in terms of sales and orders?

• What is the regional contribution to overall revenue?

**Charts:**

• Map chart: Sales by postal code.

• Bar chart: Orders by region.

• Pie chart: Sales by city/state.

**Summary of Potential Questions:**

1. How many customers, orders, and products are there?

2. What is the sales distribution across customer segments?

3. What are the top-selling products?

4. Which customers contribute the most to revenue?

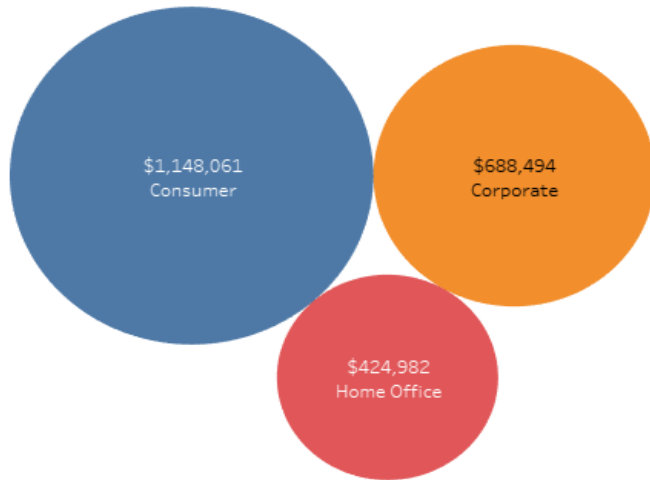5. What are the geographical patterns of sales?

Edit with WPS Office

# Create charts and visualizations

- In this Process our Team used **tableau desktop.**

- Next, a few examples of charts

- This chart shows amount of sales by Segment and it conclude that Consumer
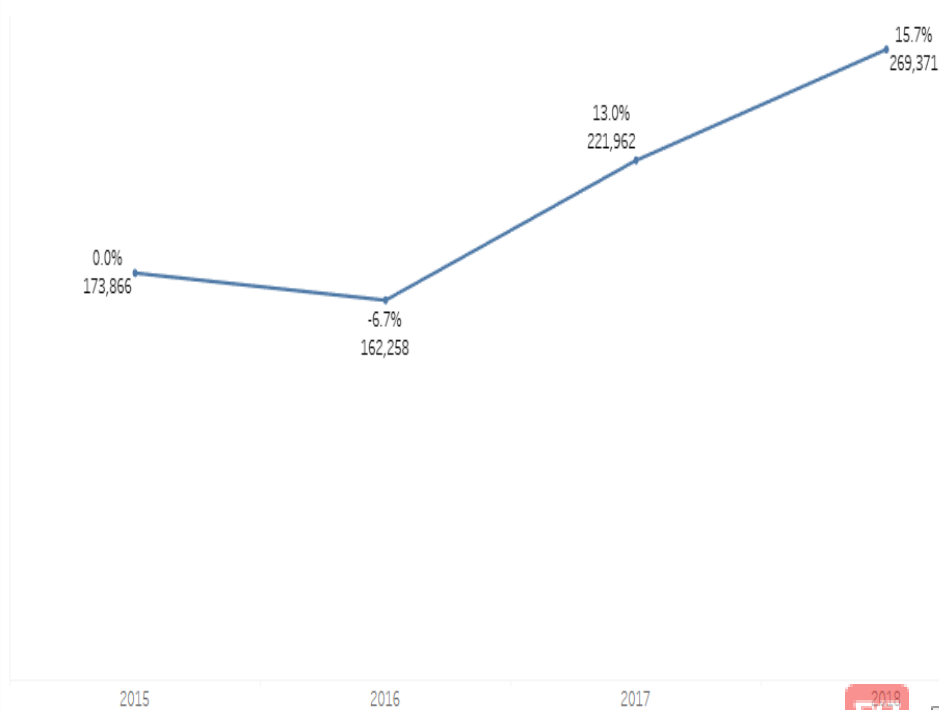
Total sales by Segment

$1,148,061
Consumer

$688,494
Corporate

$424,982
Home Office

# charts

- The growth rate is one of the most important KPI that give us an idea about the performance of the store by years.
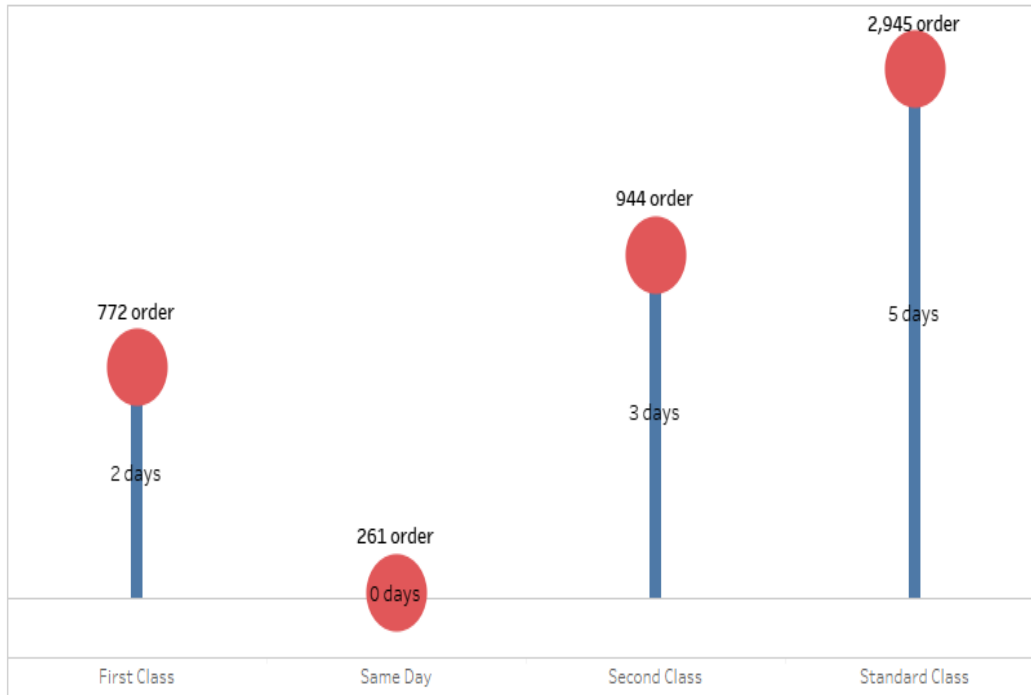
## Growth Rate by years



| | | | 15.7%<br>269,371 |
|---|---|---|---|
| | | 13.0%<br>221,962 | |
| 0.0%<br>173,866 | -6.7%<br>162,258 | | |

2015    2016    2017    2018

**charts**

- Lollipop chart is amazing figure to show the relation between number of orders and shipping mode and average days top ship.
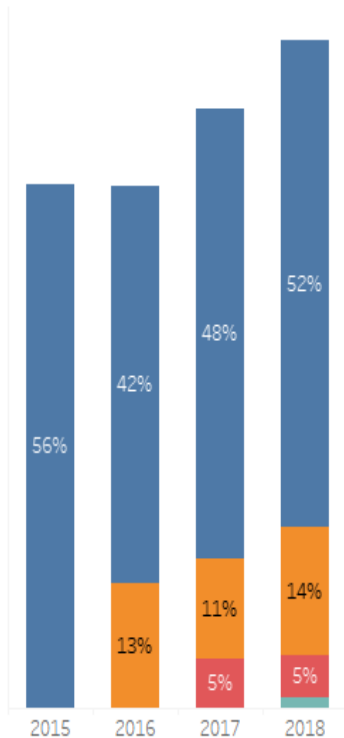

AVG Shipping days By Shipping Mode

charts

- This chart explains that old customers represented by blue color have the large scale in sales over years and represent the **loyalty** of the old customers

**Performance for old customers on sales**



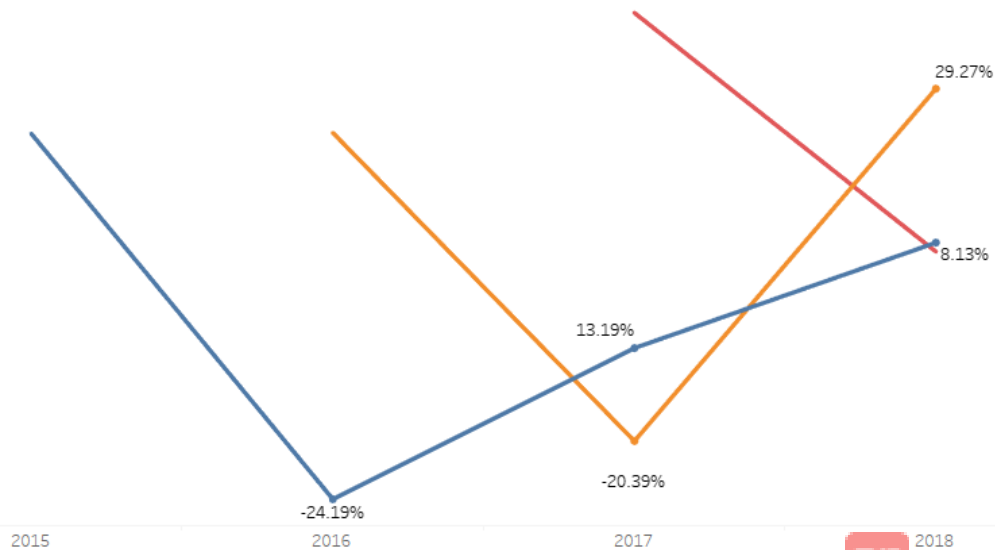| Year | Value |
|------|-------|
| 2015 | 56% |
| 2016 | 42%, 13% |
| 2017 | 48%, 11%, 5% |
| 2018 | 52%, 14%, 5% |

# charts

Digital Egypt Pioneers

- Retention Rate also considering important metric to measure how satisfied the customers

- If retention rate is high it means the same customers come again and make orders

Retention rate for customers



29.27%

8.13%

13.19%

-20.39%

-24.19%

2015    2016    2017    2018

**charts**

➤ The dashboard phase was divided into three reports.

Overview dashboard

Sales Report

Customer Performance

# Creating Dashboard

# Creating Dashboard

Creating Dashboard

Creating Dashboard