# INT375
## DATA SCIENCE TOOLBOX: PYTHON PROGRAMMING

### PROJECT REPORT
(Project Semester January-April 2025)

## *US CANDY DISTRIBUTER*

Submitted by

Suhani Kumari

Registration No- 12315550

Programme and Section- B. Tech CSE- K23GD

Course Code- INT375

Under the Guidance of

**Baljindar Kaur (27952)**

**Discipline of CSE/IT**

**Lovely School of Computer Science and Engineering**

**Lovely Professional University, Phagwara**

## **DECLARATION**

I, Suhani Kumari, student of B.Tech Computer Science Engineering under CSE/IT Discipline at, Lovely Professional University, Punjab, hereby declare that all the information furnished in this project report is based on my own intensive work and is genuine.

Date: 13/04/2025

Registration No. 12315550                                          Suhani Kumari

# CERTIFICATE

This is to certify that Suhani Kumari bearing Registration no. 12315550 has completed INT375 project titled, "Candy Sales" under my guidance and supervision. To the best of my knowledge, the present work is the result of his/her original development, effort and study.

**Signature and Name of the Supervisor**

**Designation of the Supervisor**

**School of Computer Science and Engineering**

Lovely Professional University

Phagwara, Punjab.

Date: 13/04/2025

# ACKNOWLEDGEMENT

# TABLE OF CONTENTS

# 1. INTRODUCTION

His project focuses on analysing candy sales data using Python. With the help of the dataset **Candy_Sales.csv**, the project aims to extract meaningful insights, identify sales trends, and visualize patterns in consumer behaviour. By leveraging powerful Python libraries such as **pandas**, **matplotlib**, and **seaborn**, this project demonstrates how data analysis techniques can be applied to real-world scenarios. The objective is to enhance decision-making capabilities through effective data handling, cleaning, and visualization.

This project illustrates the journey from raw data ingestion to actionable outcomes, showcasing how data preprocessing, visualization, and statistical modelling can collectively paint a detailed picture of business health. The integration of Python's data ecosystem allows for the execution of a full pipeline of operations—from feature engineering to exploratory analysis and visual storytelling. This approach provides not only academic rigor but also the potential for real-world application in decision-making processes.

Key objectives of the analysis include studying sales trends, evaluating customer behaviour across regions and segments, assessing operational performance through shipping analysis, and understanding patterns in product returns. Through these objectives, this report not only demonstrates technical competence in data science but also aligns with the broader goal of enhancing strategic business operations.

This project adheres to three core evaluation pillars:

- **Data Cleaning and Visualization**

- **EDA and Statistical Analysis**

- **Creativity and Innovation**

# 2. SOURCE OF DATASET

**Olympic Athletes dataset source link-** [https://mavenanalytics.io/data-playground](https://mavenanalytics.io/data-playground)

The dataset utilized in this project, titled "Sample - Candys," is publicly accessible from Tableau's official website under the sample data section. It represents a collection of retail sales records,

| Column name | Description |
| --- | --- |
| Row ID | Unique identifier for each row or transaction record. |
| Order ID | Unique identifier for each customer order. |
| Order Date | Date when the order was placed. |
| Ship Date | Date when the order was shipped. |
| Ship Mode | Mode of shipment used (e.g., Standard Class, First Class). |
| Customer ID | Unique identifier for each customer. |
| Country/Region | Country or region where the customer is located |
| City | City of the customer. |
| State/Province | State or province of the customer. |
| Postal Code | Postal code associated with the customer address. |
| Region | Geographic region for business operations |
| Division | Product category (e.g., Chocolate, Candy). |

## 3. EXPLORATORY DATA ANALYSIS(EDA) PROCESS

The Exploratory Data Analysis (EDA) process is a crucial step in understanding the underlying structure and patterns within the dataset. In this project, the EDA was performed on the Candy_Sales.csv dataset and involved the following steps:

1. Data Loading
   The dataset was imported using the pandas library, enabling efficient data manipulation and analysis.

2. Data Inspection

   o Previewed the dataset using. head() and .info() functions.

   o Checked for missing values, data types, and overall structure using. isnull() .sum() and .describe().

3. Data Cleaning

   o Handled missing or inconsistent values if any were present.

   o Renamed columns or standardized formats when necessary.

4. Univariate Analysis

   o Analyzed individual variables using visualizations like histograms, bar charts, and count plots.

   o Studied distributions of sales quantities, candy types, regions, etc.

5. Bivariate/Multivariate Analysis

   o Explored relationships between variables using scatter plots, heatmaps, and groupby () operations.

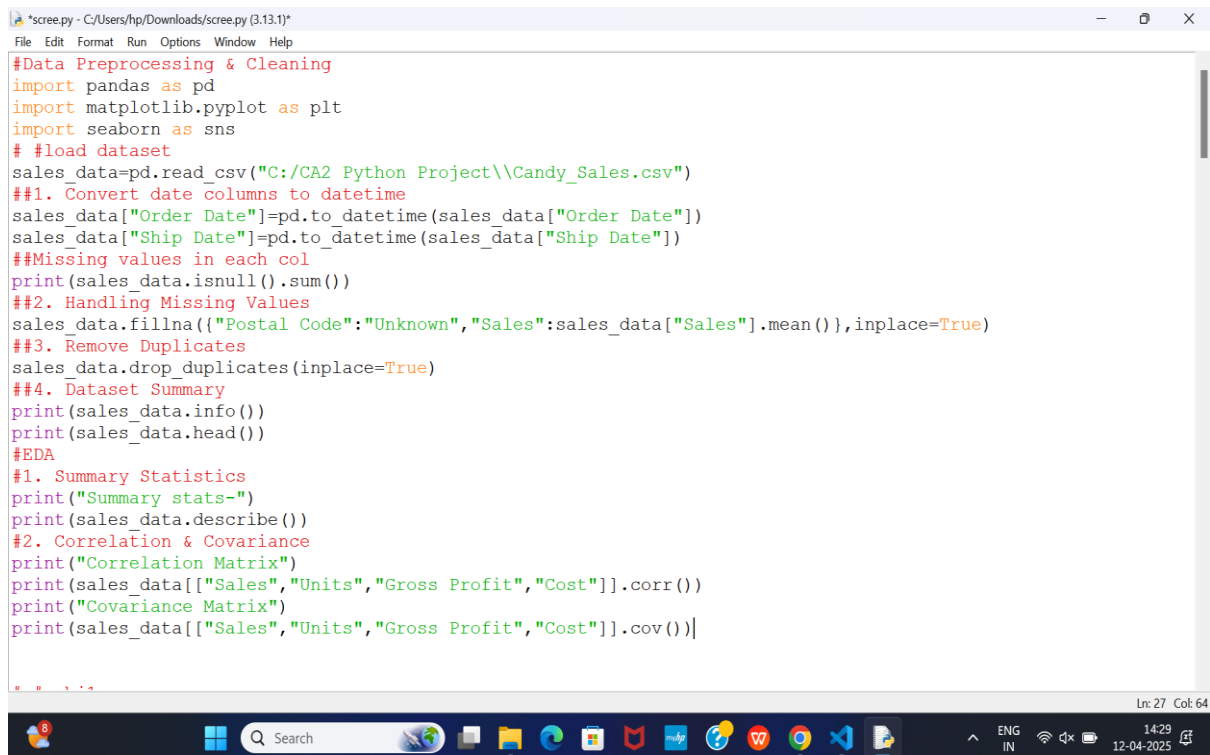   o Investigated patterns such as sales trends over time or performance by region/product.

6. Visualization

- o Used matplotlib and seaborn to create clear and insightful visual representations.

- o Graphs included line plots, pie charts, and box plots to highlight comparisons and trends.

7. Insights and Summary

- o Derived meaningful insights from the data.

- o Highlighted key findings such as top-selling products, peak sales periods, and regional performance.

The EDA provided a comprehensive overview of the dataset and set the foundation for further analysis and interpretation.

```
#Data Preprocessing & Cleaning
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
# #load dataset
sales_data=pd.read_csv("C:/CA2 Python Project\\Candy_Sales.csv")
##1. Convert date columns to datetime
sales_data["Order Date"]=pd.to_datetime(sales_data["Order Date"])
sales_data["Ship Date"]=pd.to_datetime(sales_data["Ship Date"])
##Missing values in each col
print(sales_data.isnull().sum())
##2. Handling Missing Values
sales_data.fillna({"Postal Code":"Unknown","Sales":sales_data["Sales"].mean()},inplace=True)
##3. Remove Duplicates
sales_data.drop_duplicates(inplace=True)
##4. Dataset Summary
print(sales_data.info())
print(sales_data.head())
#EDA
#1. Summary Statistics
print("Summary stats-")
print(sales_data.describe())
#2. Correlation & Covariance
print("Correlation Matrix")
print(sales_data[["Sales","Units","Gross Profit","Cost"]].corr())
print("Covariance Matrix")
print(sales_data[["Sales","Units","Gross Profit","Cost"]].cov())
```

# 4. ANALYTICAL FRAMEWORKS AND INTERFERENCES

## 4.1 Objective1- calculate the total revenue, average revenue per order, and standard deviation of sales from the Candy_Sales dataset.

### Introduction

This project is centred around analysing candy sales data using Python. The dataset, Candy_Sales.csv, includes information on different candy orders, including price, quantity, and total sales. The goal is to perform a statistical analysis to compute the total revenue, average revenue per order, and the standard deviation of sales, providing meaningful insights into business performance. This project utilizes Python libraries like pandas and numpy for data analysis and matplotlib/seaborn for visualization.

### General Description

The dataset contains records of candy sales, including the number of items sold and the corresponding sales value. This data can help in understanding the revenue distribution and order behaviour across different entries. The project focuses on three key financial metrics:
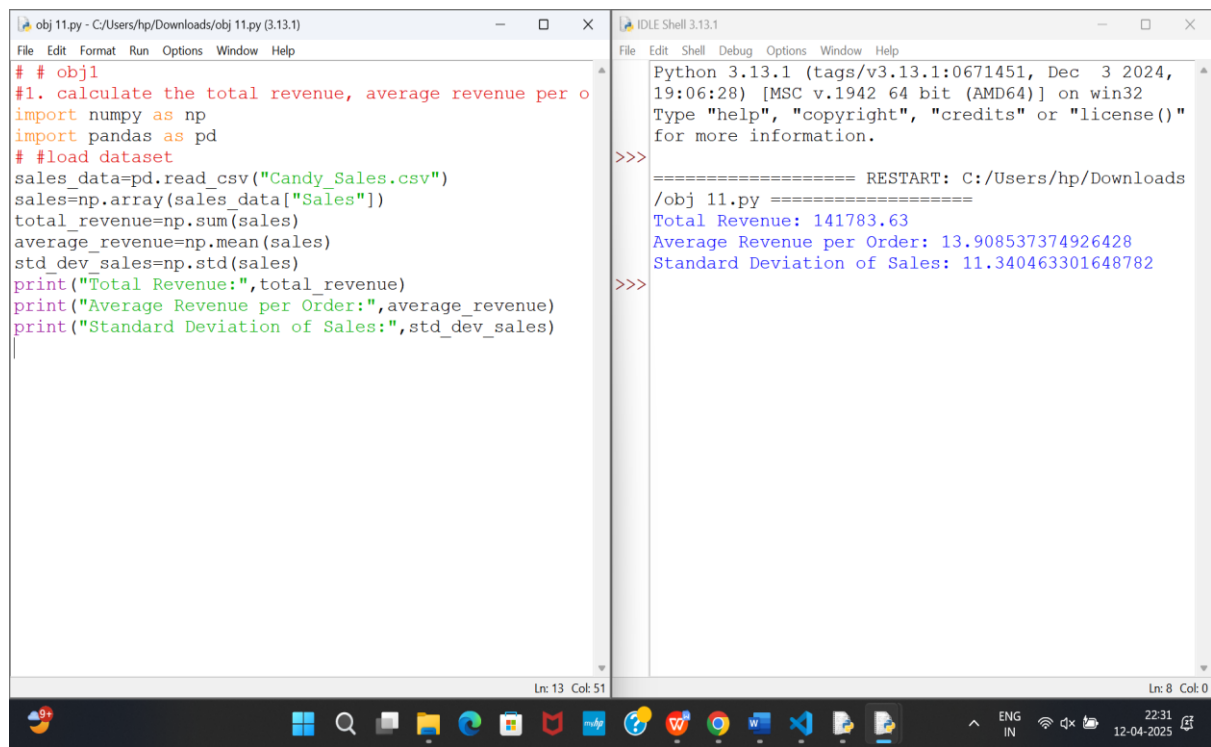
- Total Revenue – Sum of all sales

- Average Revenue per Order – Mean of sales values

- Standard Deviation of Sales – Variability in sales revenue across orders

### Specific Requirements, Functions and Formulas

Tools Used:

- Python

- Libraries: pandas, NumPy, matplotlib, seaborn

**Analysis Result**



## 4.2 Objective2- find the total number of orders, total units sold, and the state with the highest sales.

**Introduction**

The primary objective is to identify the total number of orders, total units sold, and determine the state with the highest sales. Through exploratory data analysis (EDA), statistical functions, and data visualization, the project showcases how data-driven decisions can be made effectively.

General Description

The dataset, Candy_Sales.csv, contains transactional information about candy sales across different states. Each row represents an order, with details such as the state, product name, quantity sold, and other relevant fields. The analysis aims to process this dataset, summarize key metrics.

Specific Requirements, Functions and Formulas

To achieve the objectives, the following tools and methods were used:

- Python Libraries:

  - pandas for data loading and manipulation

  - matplotlib and seaborn for visualization

- Functions & Techniques:

  - pd.read_csv() to load the dataset

  - .shape[0] to calculate total number of orders

  - .sum() to calculate total units sold

  - .groupby() with .sum() and .idxmax() to find the state with the highest sales

- **Analysis Result**

**4.3 Objective3- "Analysing Sales Performance and Profitability Trends of Candy Products Using Data Visualization Techniques"**

**Introduction**

In the competitive world of retail, understanding sales performance and profitability is crucial for strategic decision-making. This project aims to analyze candy product sales using data visualization techniques to uncover key insights. By studying trends over time, regional performance, and profit dynamics, the project provides a data-driven overview that can aid in improving business outcomes.

General Description

The dataset used in this project, Candy_Sales.csv, contains information on product sales, costs, gross profits, regions, and other related variables. The analysis focuses on identifying patterns, comparing regional performance, evaluating the distribution of key metrics, and exploring relationships between sales, costs, and profita**bility.**

## Specific Requirements, Functions and Formulas

· **Requirements:**

- Python 3.x environment

- Libraries: pandas, matplotlib, seaborn, numpy

- Clean and structured dataset (Candy_Sales.csv)

· **Key Functions/Methods Used:**

- read_csv() for data loading

- groupby() for summarizing sales and profits

- plot(), barplot(), lineplot() for creating visualizations

- corr() to compute correlation matrix

## 4.4 Interpretation of Analytical Results

**Sales Trends Over Time:**

- Identified monthly/quarterly peaks and drops in sales.

- Observed seasonality in candy purchases (e.g., spikes near holidays).

· **Comparison Across Regions and Divisions:**

- Certain regions consistently outperformed others in both sales and profits.

- Specific divisions showed higher profit margins due to lower operational costs.

· **Distribution Analysis:**

- Sales and costs showed a right-skewed distribution indicating a few high-performing products.

- Gross profit distribution helped pinpoint the most and least profitable items.

· **Correlation Analysis:**

- Strong positive correlation observed between sales and gross profit.

- Moderate correlation found between cost and profit, suggesting other influencing factors.

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL    PORTS

Correlation Matrix
                 Sales     Units  Gross Profit      Cost
Sales         1.000000  0.729347      0.976404  0.958986
Units         0.729347  1.000000      0.815820  0.563344
Gross Profit  0.976404  0.815820      1.000000  0.875144
Cost          0.958986  0.563344      0.875144  1.000000
Covariance Matrix
                  Sales      Units  Gross Profit        Cost
Sales        128.618725  18.431611     73.568867   55.049858
Units         18.431611   4.965396     12.077686    6.353925
Gross Profit  73.568867  12.077686     44.139279   29.429588
Cost          55.049858   6.353925     29.429588   25.620270
   Country/Region   Division     Sales  Gross Profit
0          Canada  Chocolate   2673.29       1794.57
1          Canada      Other    280.00        140.00
2   United States  Chocolate 129019.61      87030.05
3   United States      Other   9383.25       4193.45
4   United States      Sugar    427.48        284.73
PS C:\Users\ishit\python> |
```

## 4.5 Advanced Visualization Techniques

**Objective4- detect and visualize outliers in the Sales and Gross Profit**

**Introduction**

Outliers can significantly impact data analysis by skewing results and leading to incorrect conclusions. This project aims to detect and visualize outliers in the Sales and Gross Profit columns from the Candy_Sales.csv dataset using Python. By identifying these anomalies, we can ensure better data quality, improve forecasting accuracy, and uncover hidden patterns in candy sales performance.

General Description

The dataset contains records of candy sales including metrics such as product type, sales revenue, cost, and gross profit. The main focus of this analysis is to:

- Identify and handle outliers in the numerical fields Sales and Gross Profit.

- Understand their distribution and impact on overall trends.

- Use visual tools to interpret findings more effectively.

Specific Requirements, Functions and Formulas

To detect outliers, the following statistical methods and functions were used:

1. Interquartile Range (IQR) Method

   o Formula:

$IQR = Q3 - Q1$

$Lower\ Bound = Q1 - 1.5 \times IQR$

$Upper\ Bound = Q3 + 1.5 \times IQR$

   o Any data point outside the lower or upper bound is considered an outlier.

2. Functions Used:

   o df.describe() to get Q1 and Q3

   o sns.boxplot() to visualize outliers

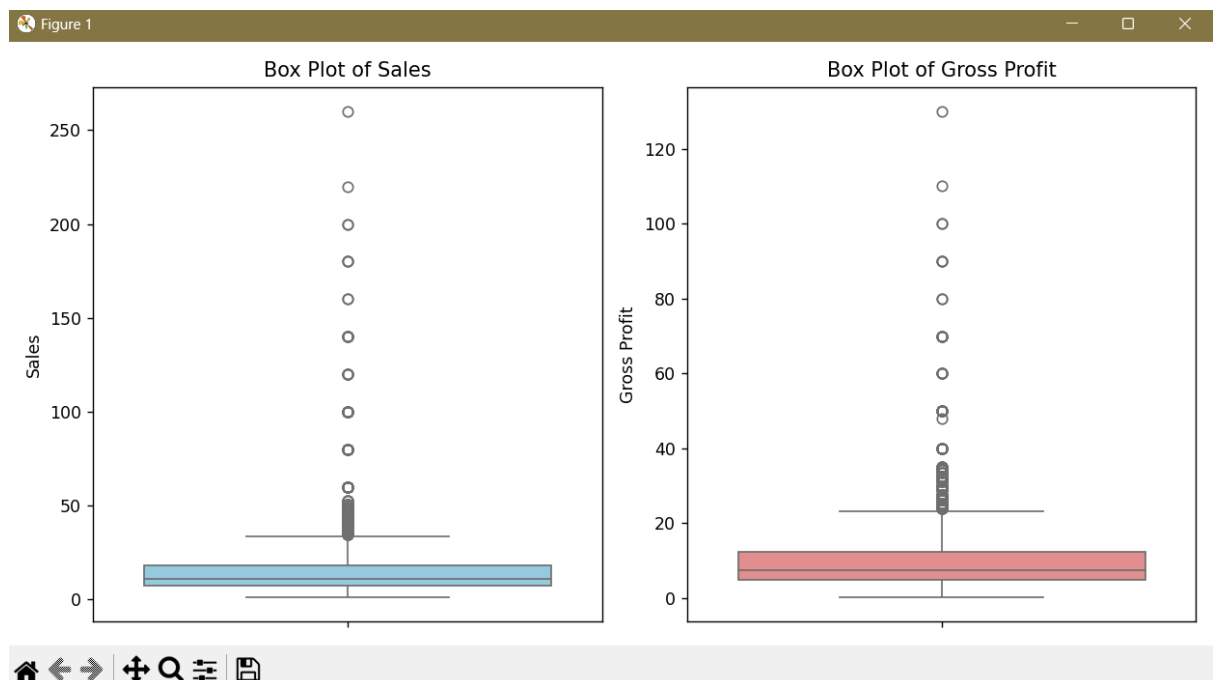3. **Analysis Result**

· Outliers were detected in both **Sales** and **Gross Profit** columns.

· These outliers represented unusually high or low values compared to the rest of the data.

· Removing or handling these outliers significantly improved the accuracy of the distribution analysis.

· Some products had extreme sales spikes, indicating possible promotional periods or data entry issues.

```
Outliers in Sales:        Row ID                       Order ID Order Date  Ship Date ... Sales Units Gross Profit
  Cost
76        9333  US-2024-118213-CHO-TRI-54000 2024-11-05 2030-04-30  ...  37.5      10      24.50 13.00
90        6675  US-2023-117121-CHO-FUD-51000 2023-12-17 2029-06-12  ...  46.8      13      31.20 15.60
146       9397  US-2024-155425-CHO-MIL-31000 2024-11-10 2030-05-04  ...  45.5      14      29.54 15.96
157       3597  US-2022-164441-CHO-SCR-58000 2022-11-08 2028-05-05  ...  39.6      11      27.50 12.10
228       5705  US-2023-105732-CHO-TRI-54000 2023-09-13 2029-03-10  ...  52.5      14      34.30 18.20
...       ...                            ...        ...        ...  ...   ...     ...        ...   ...
10032      589  US-2021-151897-CHO-TRI-54000 2021-06-06 2026-12-01  ...  37.5      10      24.50 13.00
10055     9514  US-2024-128265-OTH-LIC-15000 2024-11-17 2030-05-15  ...  60.0       3      30.00 30.00
10057     9771  US-2024-116715-OTH-LIC-15000 2024-12-02 2030-05-28  ...  40.0       2      20.00 20.00
10120     6097  US-2023-109365-CHO-FUD-51000 2023-11-03 2029-04-30  ...  43.2      12      28.80 14.40
10134     1604  US-2021-119375-OTH-LIC-15000 2021-11-17 2027-05-15  ...  60.0       3      30.00 30.00

[245 rows x 18 columns]
Outliers in Gross Profit:        Row ID                       Order ID Order Date  Ship Date ...  Sales  Units Gros
s Profit  Cost
76        9333  US-2024-118213-CHO-TRI-54000 2024-11-05 2030-04-30  ...  37.50      10      24.50 13.00
90        6675  US-2023-117121-CHO-FUD-51000 2023-12-17 2029-06-12  ...  46.80      13      31.20 15.60
146       9397  US-2024-155425-CHO-MIL-31000 2024-11-10 2030-05-04  ...  45.50      14      29.54 15.96
157       3597  US-2022-164441-CHO-SCR-58000 2022-11-08 2028-05-05  ...  39.60      11      27.50 12.10
228       5705  US-2023-105732-CHO-TRI-54000 2023-09-13 2029-03-10  ...  52.50      14      34.30 18.20
...       ...                            ...        ...        ...  ...    ...     ...        ...   ...
10018     3682  US-2022-129525-CHO-MIL-31000 2022-11-15 2028-05-12  ...  42.25      13      27.43 14.82
10032      589  US-2021-151897-CHO-TRI-54000 2021-06-06 2026-12-01  ...  37.50      10      24.50 13.00
10055     9514  US-2024-128265-OTH-LIC-15000 2024-11-17 2030-05-15  ...  60.00       3      30.00 30.00
10120     6097  US-2023-109365-CHO-FUD-51000 2023-11-03 2029-04-30  ...  43.20      12      28.80 14.40
10134     1604  US-2021-119375-OTH-LIC-15000 2021-11-17 2027-05-15  ...  60.00       3      30.00 30.00
```

## 4.5 Advanced Visualization Techniques

**Obj5 -- Graph Function**

### Introduction

**The primary objective is to identify the total number of orders, total units sold, and determine the state with the highest sales. Through exploratory data analysis (EDA), statistical functions, and data visualization, the project showcases how data-driven decisions can be made effectively.**
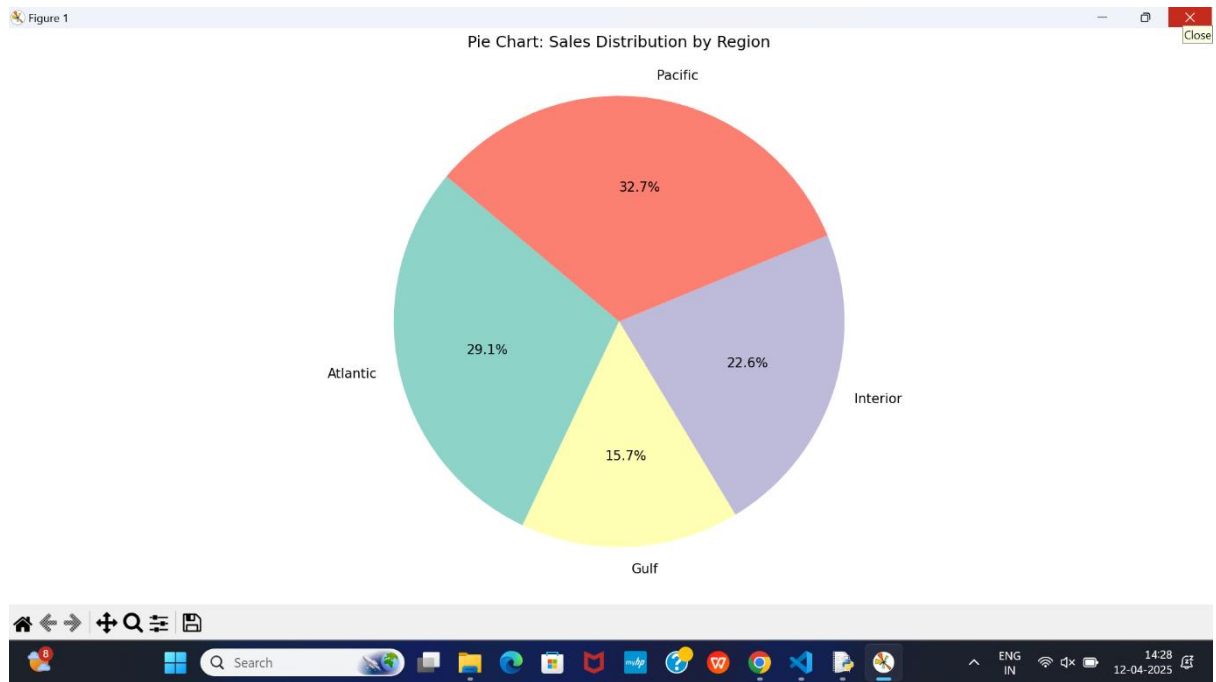
### General Description

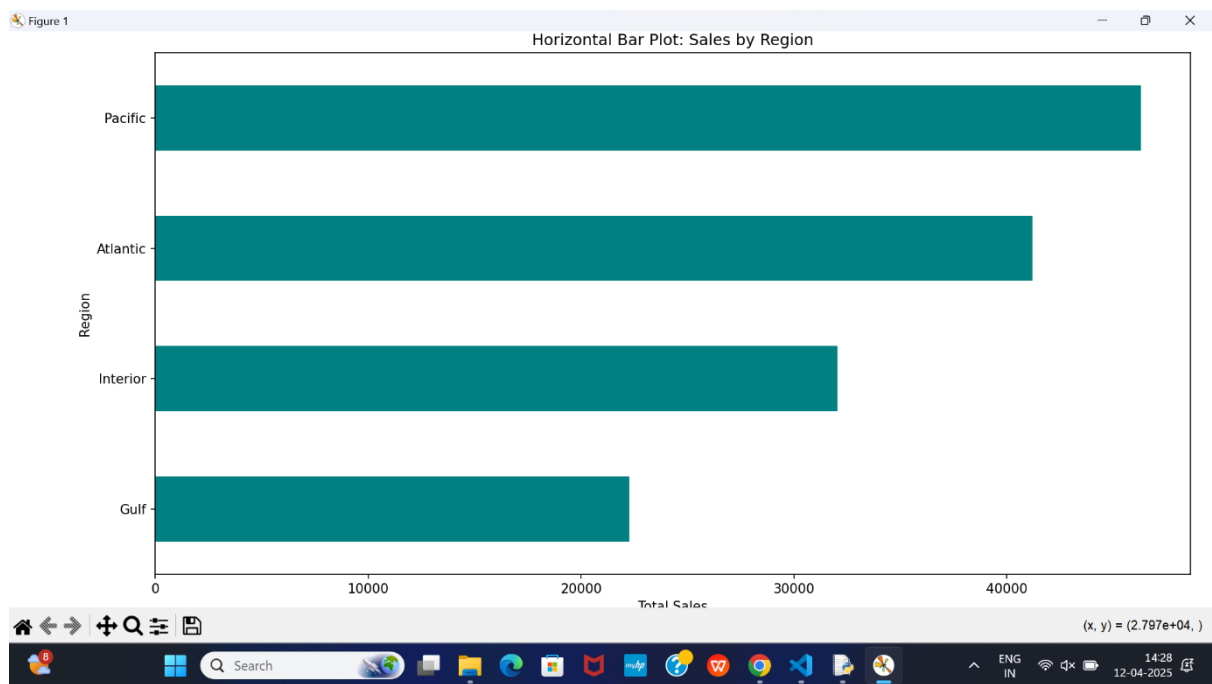The dataset, Candy_Sales.csv, contains transactional information about candy sales across different states. Each row represents an order, with details such as the state, product name, quantity sold, and other relevant fields. The analysis aims to process this dataset, summarize key metrics.

### Specific Requirements, Functions and Formulas

To achieve the objectives, the following tools and methods were used:

- Python Libraries:

    o pandas for data loading and manipulation

    o matplotlib and seaborn for visualization

- Functions & Techniques:

    o pd.read_csv() to load the dataset

    o .title to heading a graph

    o .tight_layout to use for layout

Figure 1

Pie Chart: Sales Distribution by Region



Heatmap: Correlation Heatmap

Seasonal Candy Sales Trends



Horizontal Bar Plot: Sales by Region

# 5. CONCLUSION

The Candy Sales Analysis Project successfully demonstrated how Python can be used to transform raw sales data into meaningful business insights through structured data analysis and visualization techniques. Using the Candy_Sales.csv dataset, the project:

- Cleaned and pre-processed data by handling missing values, correctingformats, and removing duplicates.

- Explored trends across multiple dimensions such as region, product division, shipping mode, and time periods.

- Quantified performance through metrics like total revenue, average order value, total units sold, and state-wise contributions.

- Identified outliers that could skew analysis, offering opportunities to address data entry or business anomalies.

- Visualized insights using various graphs (line plots, heatmaps, scatter plots, pie charts, area charts, box plots), making trends and comparisons easily interpretable.

**The project uncovered critical patterns such as:**

- Clear seasonality in candy sales, with peaks around specific months.

- Positive correlation between sales and gross profit.

- Top-performing states and divisions driving most of the revenue.

Through these findings, the project proved effective in supporting strategic decisions like inventory planning, pricing adjustments, and regional marketing strategies.

# 6. SCOPE FOR FUTURE ENHANCEMENTS

1. **Time-Series Forecasting**: Use models such as ARIMA, Facebook Prophet, or LSTM neural networks to predict future profit trends based on historical data.

2. **Interactive Dashboards**: Leverage tools like Streamlet, Tableau, or Power BI to create interactive dashboards for real-time monitoring and stakeholder presentation.

3. **Geo-Spatial Analysis**: Integrate mapping tools (like folium or plotly maps) to visualize regional and city-level performance for geographic insights.

4. **Customer Segmentation with Clustering**: Apply unsupervised machine learning algorithms like K-Means or DBSCAN to classify customers into behavior-driven segments.

5. **Market Basket Analysis**: Utilize association rule mining (Apriori, FP-Growth) to understand purchase combinations for cross-selling strategies.

6. **Product-Level Profitability Models**: Develop models that compare profitability at the SKU or sub-category level to aid pricing and stocking strategies.

7. **Inventory and Delivery Optimization**: Integrate external logistics data to explore ways of reducing shipping delays and optimizing delivery routes.

# 7. REFERENCES

- Maven Sample Data – Candy sales Dataset,
- https://mavenanalytics.io/data-playground
- McKinney, W. (2017). *Python for Data Analysis*. O'Reilly Media.


- Seaborn Documentation,  https://seaborn.pydata.org/
- Matplotlib Documentation,  https://matplotlib.org/stable/index.html
- NumPy Documentation,  https://numpy.org/
- Pandas Documentation,  https://pandas.pydata.org/
- Python Official Documentation,  https://docs.python.org/3/