SUPERSTORE DATA ANALYSIS USING DATA ANALYTICS AND DATA SCIENCE

SUPERSTORE DATA ANALYSIS USING AI WITH DATA SCIENCE

Presented By

STUDENT NAME: S.SHRRIVATHSAN

COLLEGE NAME: PANIMALAR ENGINEERING COLLEGE

DEPARTMENT: COMPUTER SCIENCE DEPARTMENT

EMAIL ID: SHRRIVATHSANY2K@GMAIL.COM

AICTE STUDENT ID: STU680F6B03BD1FE1745840899

INTERNSHIP ID:

INTERNSHIP_1748937226683EAA0A58ABC



OUTLINE

- PROBLEM STATEMENT (SHOULD NOT INCLUDE SOLUTION)
- PROPOSED SYSTEM/SOLUTION
- SYSTEM DEVELOPMENT APPROACH (TECHNOLOGY USED)
- ALGORITHM & DEPLOYMENT
- RESULT (OUTPUT IMAGE)
- Conclusion
- FUTURE SCOPE
- REFERENCES



PROBLEM STATEMENT

PROBLEM STATEMENT:

THE OBJECTIVE IS TO ANALYZE THE HISTORICAL SALES AND PROFIT DATA OF A SUPERSTORE TO IDENTIFY KEY TRENDS, PATTERNS, AND ANOMALIES, AND SUBSEQUENTLY LEVERAGE THESE INSIGHTS TO OPTIMIZE BUSINESS OPERATIONS, ENHANCE CUSTOMER SATISFACTION, AND MAXIMIZE PROFITABILITY. THIS INVOLVES USING DATA SCIENCE TECHNIQUES, INCLUDING DESCRIPTIVE ANALYTICS, PREDICTIVE MODELING, AND POTENTIALLY PRESCRIPTIVE ANALYTICS, TO ADDRESS SPECIFIC BUSINESS QUESTIONS AND INFORM STRATEGIC DECISION-MAKING.

KEY OBJECTIVES

- IDENTIFY KEY DRIVERS OF SALES AND PROFITABILITY: UNDERSTANDING WHICH FACTORS CONTRIBUTE MOST SIGNIFICANTLY TO SALES AND PROFIT
 ACROSS DIFFERENT PRODUCTS, REGIONS, AND CUSTOMER SEGMENTS.
- FORECAST FUTURE SALES AND DEMAND: ACCURATELY PREDICT SALES AND DEMAND FOR VARIOUS PRODUCTS TO OPTIMIZE INVENTORY LEVELS AND
 PREVENT STOCKOUTS OR OVERSTOCKING, WHICH CAN REDUCE WASTE AND SAVE COSTS.
- SEGMENT CUSTOMERS AND PERSONALIZE RECOMMENDATIONS: IDENTIFY DISTINCT CUSTOMER SEGMENTS BASED ON THEIR PREFERENCES AND BUYING BEHAVIOR TO TAILOR MARKETING CAMPAIGNS AND PRODUCT RECOMMENDATIONS, ENHANCING CUSTOMER SATISFACTION AND LOYALTY.
- IDENTIFY UNDERPERFORMING PRODUCTS AND CATEGORIES: PINPOINT PRODUCTS OR CATEGORIES EXPERIENCING LOW SALES OR PROFITABILITY TO DEVELOP STRATEGIES FOR IMPROVEMENT, SUCH AS ADJUSTING PRICING, PROMOTIONAL ACTIVITIES, OR EXPLORING ALTERNATIVE SOURCING.
- IMPROVE SHIPPING AND LOGISTICS EFFICIENCY: ANALYZE SHIPPING MODE PERFORMANCE AND OPTIMIZE LOGISTICS TO REDUCE COSTS AND DELIVERY TIMES, ULTIMATELY IMPROVING CUSTOMER SATISFACTION AND RETENTION.



DATA SCIENCE APPROACH AND AI TECHNIQUES

THE PROJECT WILL INVOLVE THE FOLLOWING STEPS AND TECHNIQUES:

- DATA COLLECTION AND CLEANING:
 - GATHER HISTORICAL SALES DATA, INCLUDING PRODUCT INFORMATION, CUSTOMER DETAILS, ORDER HISTORY, AND SHIPPING INFORMATION.
 - CLEAN AND PREPROCESS THE DATA TO HANDLE MISSING VALUES, INCONSISTENCIES, AND OUTLIERS, ENSURING DATA QUALITY FOR ANALYSIS AND MODEL BUILDING.
- EXPLORATORY DATA ANALYSIS (EDA):
 - CONDUCT EDA TO UNDERSTAND THE DATA STRUCTURE, DISTRIBUTIONS, AND INITIAL TRENDS IN SALES, PROFIT, AND OTHER KEY METRICS ACROSS VARIOUS DIMENSIONS (CATEGORIES, REGIONS, SEGMENTS, ETC.).
 - UTILIZE VISUALIZATIONS LIKE BAR CHARTS, SCATTER PLOTS, AND GEOGRAPHICAL MAPS TO UNCOVER PATTERNS, DISPARITIES, AND ANOMALIES IN THE SALES DATA.
- FEATURE ENGINEERING:
 - CREATE NEW FEATURES FROM EXISTING DATA, SUCH AS SEASONALITY INDICATORS, DISCOUNT RATIOS, OR CUSTOMER LIFETIME VALUE METRICS, TO ENHANCE THE PREDICTIVE POWER OF MODELS.



Model Building and Training:

- Sales Forecasting: Employ Al and machine learning models like time series forecasting (e.g., ARIMA or Prophet) to predict future sales based on historical data and seasonal patterns.
- Customer Segmentation: Apply unsupervised learning algorithms (e.g., K-Means clustering) to segment customers based on purchase history, demographics, and behavior.
- Product Recommendation: Build personalized recommendation systems using collaborative filtering or content-based filtering techniques to suggest products to customers based on their preferences and browsing history.
- Profitability Analysis: Utilize regression models or machine learning algorithms to identify factors influencing profitability and predict potential losses or areas of improvement.

Model Evaluation:

- Evaluate the performance of the developed models using appropriate metrics, such as Mean Absolute Percentage Error (MAPE) for forecasting or accuracy/precision/recall for classification tasks.
- Deployment and Monitoring:
 - Deploy the models into a production environment for continuous analysis and prediction.
 - Establish a monitoring system to track model performance and ensure the generated insights remain accurate and relevant over time.



PROPOSED SOLUTION

The proposed system aims to address the challenge of predicting the required bike count at each hour to ensure a stable supply of rental bikes. This involves leveraging data analytics and machine learning techniques to forecast demand patterns accurately. The solution will consist of the following components:

DATA COLLECTION:

- GATHER HISTORICAL DATA ON BIKE RENTALS, INCLUDING TIME, DATE, LOCATION, AND OTHER RELEVANT FACTORS.
- UTILIZE REAL-TIME DATA SOURCES, SUCH AS WEATHER CONDITIONS, EVENTS, AND HOLIDAYS, TO ENHANCE PREDICTION ACCURACY.

DATA PREPROCESSING:

- CLEAN AND PREPROCESS THE COLLECTED DATA TO HANDLE MISSING VALUES, OUTLIERS, AND INCONSISTENCIES.
- FEATURE ENGINEERING TO EXTRACT RELEVANT FEATURES FROM THE DATA THAT MIGHT IMPACT BIKE DEMAND.

MACHINE LEARNING ALGORITHM:

- IMPLEMENT A MACHINE LEARNING ALGORITHM, SUCH AS A TIME-SERIES FORECASTING MODEL (E.G., ARIMA, SARIMA, LSTM, XGBOOST), TO PREDICT BIKE COUNTS BASED ON HISTORICAL PATTERNS.
- CONSIDER INCORPORATING OTHER FACTORS LIKE WEATHER CONDITIONS, DAY OF THE WEEK, AND SPECIAL EVENTS TO IMPROVE PREDICTION ACCURACY.

DEPLOYMENT:

- Develop a user-friendly interface or application that provides real-time predictions for bike counts at different hours.
- DEPLOY THE SOLUTION ON A SCALABLE AND RELIABLE PLATFORM, CONSIDERING FACTORS LIKE SERVER INFRASTRUCTURE, RESPONSE TIME, AND USER ACCESSIBILITY.

EVALUATION:

- ASSESS THE MODEL'S PERFORMANCE USING APPROPRIATE METRICS SUCH AS MEAN ABSOLUTE ERROR (MAE), ROOT MEAN SQUARED ERROR (RMSE), OR OTHER RELEVANT METRICS.
- Fine-tune the model based on feedback and continuous monitoring of prediction accuracy.



SYSTEM APPROACH

The "System Approach" section outlines the overall strategy and methodology for developing and implementing the rental bike prediction system. Here's a suggested structure for this section:

- System requirement's
- Library required to build the model



SYSTEM REQUIREMENTS

1. HARDWARE REQUIREMENTS:

- PROCESSOR: INTEL 15/17 OR EQUIVALENT (MULTI-CORE RECOMMENDED)
- RAM: MINIMUM 8 GB (16 GB RECOMMENDED FOR FASTER PROCESSING)
- STORAGE: MINIMUM 1 GB FREE SPACE (FOR DATASET AND LIBRARIES)

2. SOFTWARE REQUIREMENTS:

- OS: Windows 10/11, MacOS, or any Linux distribution
- PYTHON VERSION: 3.8 OR ABOVE
- IDE: Jupyter Notebook, VS Code, or Google Colab



LIBRARY REQUIRED TO BUILD THE MODEL

1. Data Handling & Manipulation

- •PANDAS For reading and manipulating tabular data
- •NUMPY FOR NUMERICAL OPERATIONS

2. DATA VISUALIZATION

- •MATPLOTLIB For basic plotting and visualization
- •SEABORN FOR ADVANCED STATISTICAL PLOTS

3. DATA PREPROCESSING

•SKLEARN.PREPROCESSING - For encoding, scaling, imputing data

4. MODEL BUILDING

- •SKLEARN.LINEAR_MODEL For linear regression models
- •SKLEARN.ENSEMBLE For models like Random Forest, Gradient Boosting
- •XGBOOST For extreme gradient boosting (if needed)
- •LIGHTGBM For faster and efficient boosting (optional)



LIBRARY REQUIRED TO BUILD THE MODEL

5. MODEL EVALUATION

•SKLEARN.METRICS - For performance evaluation metrics like RMSE, R2, MAE

6. OPTIONAL (ADVANCED VISUALIZATION & TUNING)

- •PLOTLY For interactive charts
- •SCIPY For statistical operations
- •SKLEARN.MODEL_SELECTION For cross-validation and hyperparameter tuning



ALGORITHM & DEPLOYMENT

1. PROJECT OBJECTIVE

- ANALYZE THE SUPERSTORE DATASET TO:
- **UNDERSTAND SALES PERFORMANCE**
- IDENTIFY TOP-PERFORMING SEGMENTS, PRODUCTS, AND REGIONS
- PREDICT FUTURE SALES
- DETECT ANOMALIES (E.G., LOSS-GENERATING ITEMS)
- PROVIDE DECISION SUPPORT VIA A DEPLOYED WEB APPLICATION

2. DATA SCIENCE WORKFLOW

- A. EXPLORATORY DATA ANALYSIS (EDA)

 •VISUALIZE SALES, PROFIT BY REGION, CATEGORY, AND CUSTOMER SEGMENT
- HEATMAPS FOR CORRELATION
- TIME SERIES OF SALES OVER TIME
- DENTIFY LOSS-MAKING SUB-CATEGORIES

B. FEATURE ENGINEERING

- •EXTRACT YEAR, MONTH FROM ORDER DATE
- •COMPUTE PROFIT MARGIN
- •AGGREGATE METRICS BY CUSTOMER/PRODUCT/REGION

C. MACHINE LEARNING ALGORITHMS

CHOOSE ML ALGORITHMS BASED ON YOUR GOALS:

1. SALES PREDICTION

- •ALGORITHMS: LINEAR REGRESSION, RANDOM FOREST, XGBOOST
 •INPUT: PRODUCT, REGION, TIME, DISCOUNT, QUANTITY
- •OUTPUT: EXPECTED SALES/PROFIT



ALGORITHM & DEPLOYMENT

2. Customer Segmentation

•Algorithm: K-Means Clustering

•Features: Total Sales, Orders, Average Profit

Outcome: Targeted marketing

3. Anomaly DetectionAlgorithms: Isolation Forest, DBSCAN

•Detect: Unusual profit margins, high discounts

4. Product Recommendation

•Algorithm: Association Rules (Apriori), Collaborative Filtering

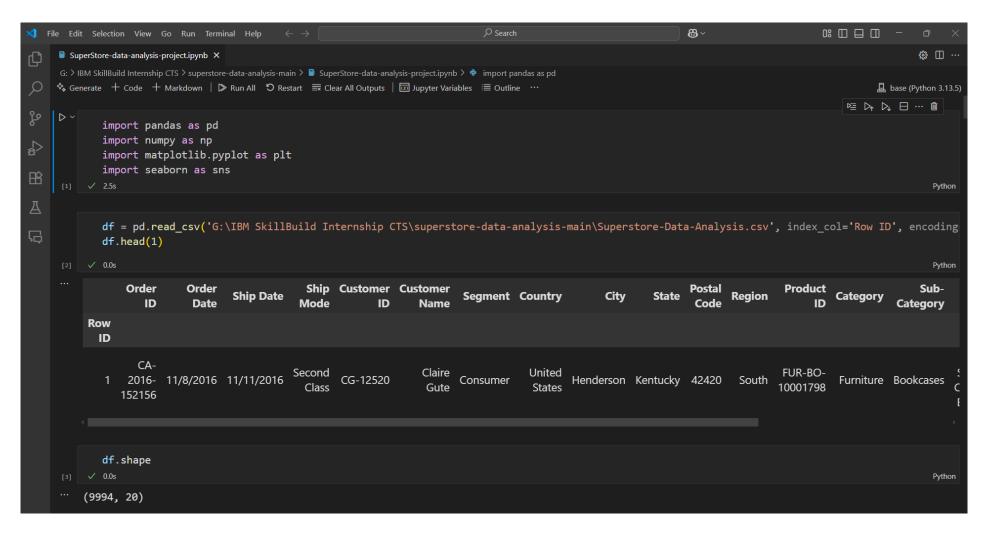
- 4. Al Integration (Optional)
- Use LSTM (Deep Learning) for time-series sales forecasting
- Integrate NLP to generate insights like: "Top 5 loss-making products in the West region in 2024".

5. Tools and Libraries

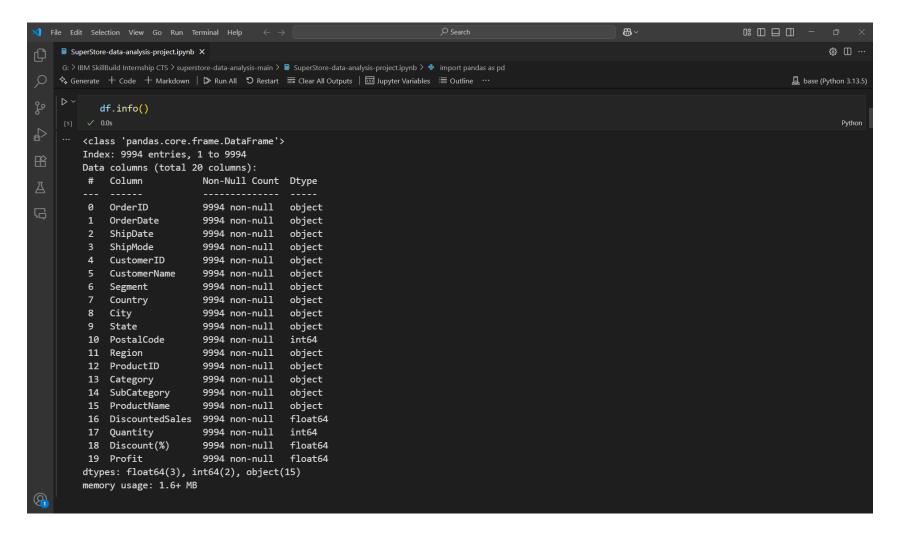
Python Libraries

- pandas, numpy, matplotlib, seaborn for EDA
- •scikit-learn ML algorithms
- xgboost, lightgbm advanced models
- •statsmodels or fbprophet time series forecasting
- •flask or streamlit for web deployment

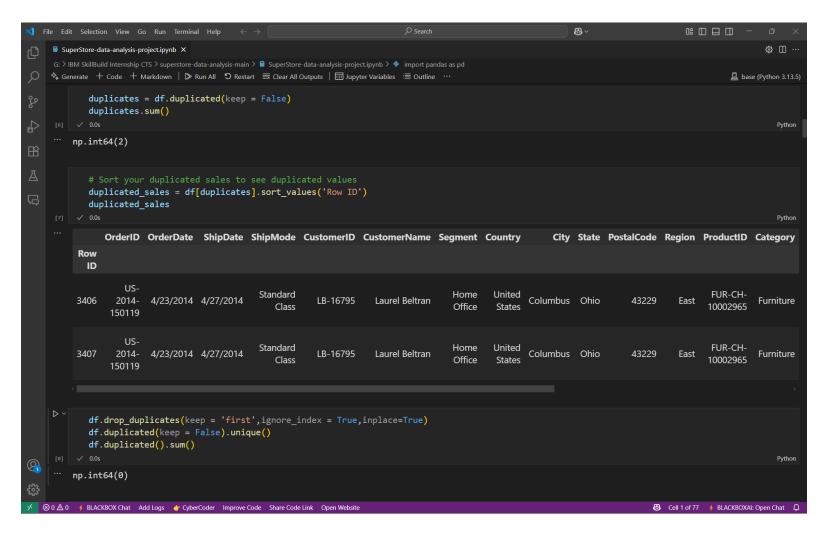




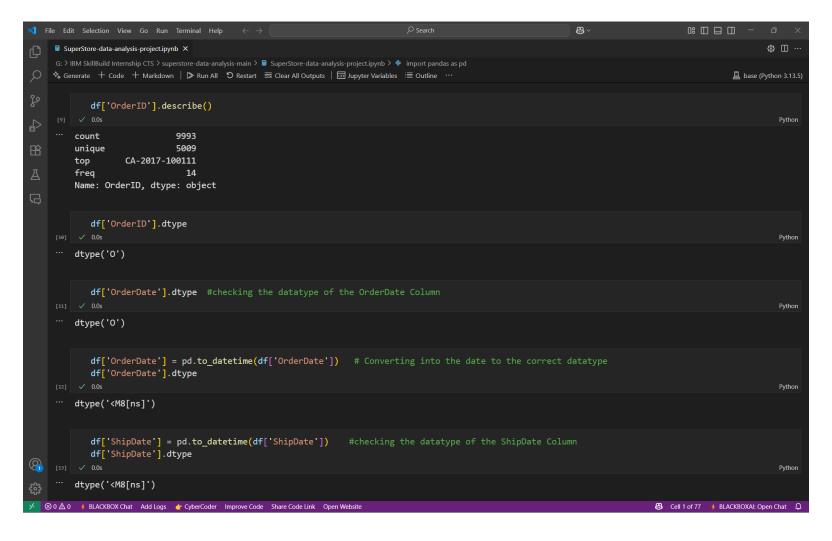




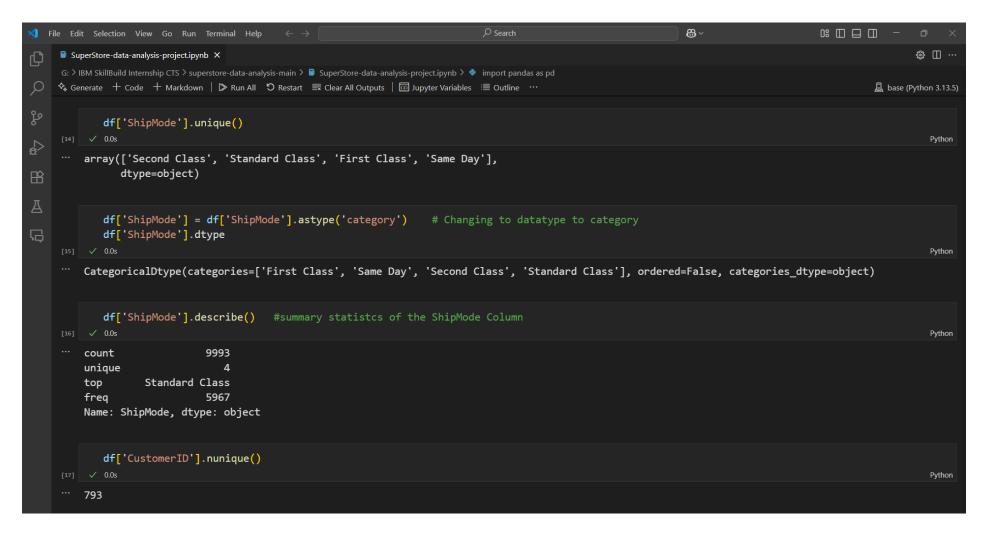




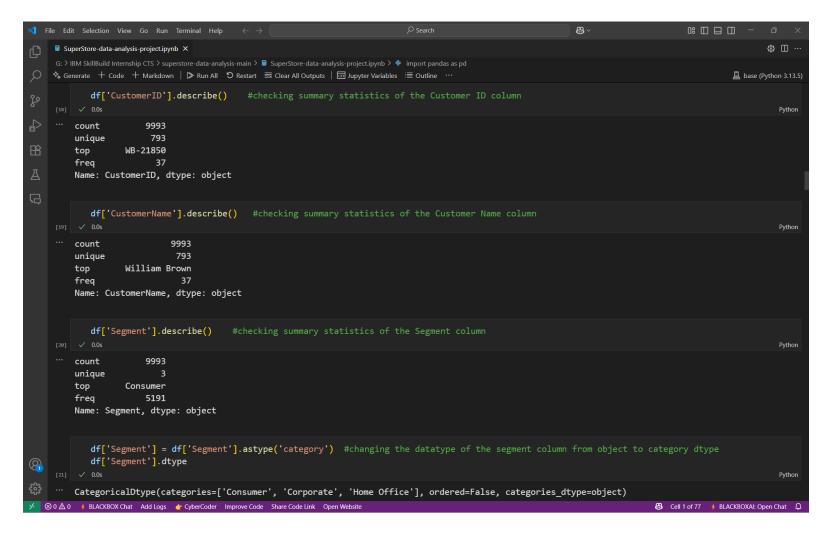








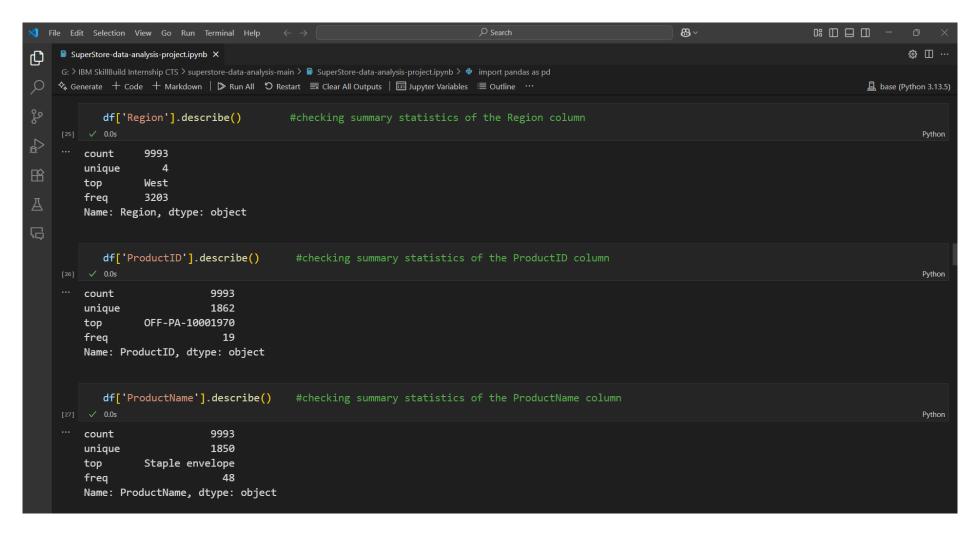




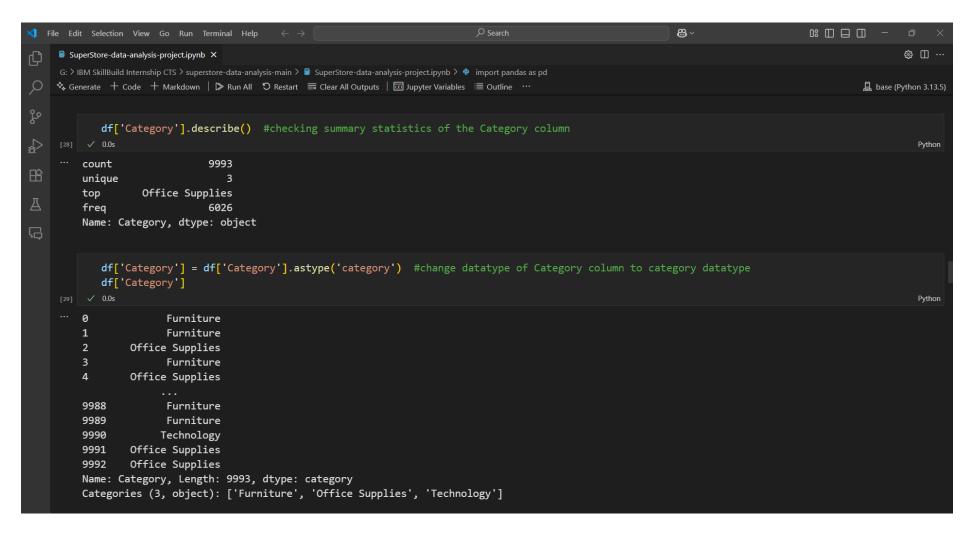


```
88 ~
                                                                                                                                               ∰ Ⅲ …
     SuperStore-data-analysis-project.ipynb X
     G: > IBM SkillBuild Internship CTS > superstore-data-analysis-main > 📳 SuperStore-data-analysis-project.ipynb > 🍨 import pandas as pd
🔘 🍫 Generate 🕂 Code 🕂 Markdown | ⊳ Run All 💆 Restart 🚍 Clear All Outputs | 👼 Jupyter Variables :≣ Outline ...
                                                                                                                                                         base (Python 3.13.5)
            df['City'].describe()
                                           #checking summary statistics of the City column
                              9993
         count
                               531
         unique
                    New York City
         freq
                               915
         Name: City, dtype: object
            df['State'].describe()
                                          #checking summary statistics of the State column
         count
                           9993
                             49
         unique
                    California
         top
         freq
                           2001
         Name: State, dtype: object
            df['PostalCode'] = df['PostalCode'].astype('str') # Changing PostalCode from int to str
            df['PostalCode'].describe()
                                              #checking summary statistics of the PostalCode column
     [24] 		 0.0s
                                                                                                                                                                   Python
        count
                     9993
         unique
                      631
         top
                    10035
                      263
         Name: PostalCode, dtype: object
```

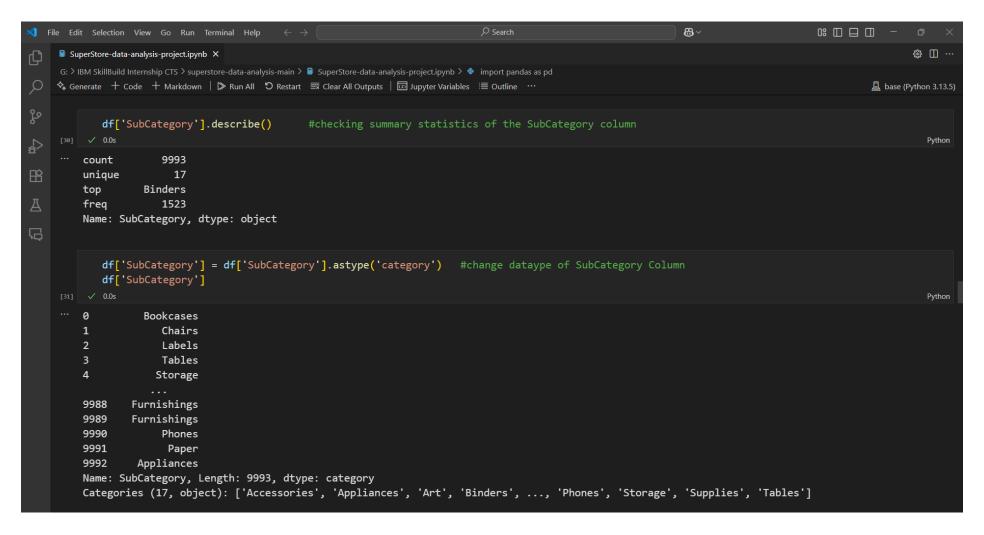




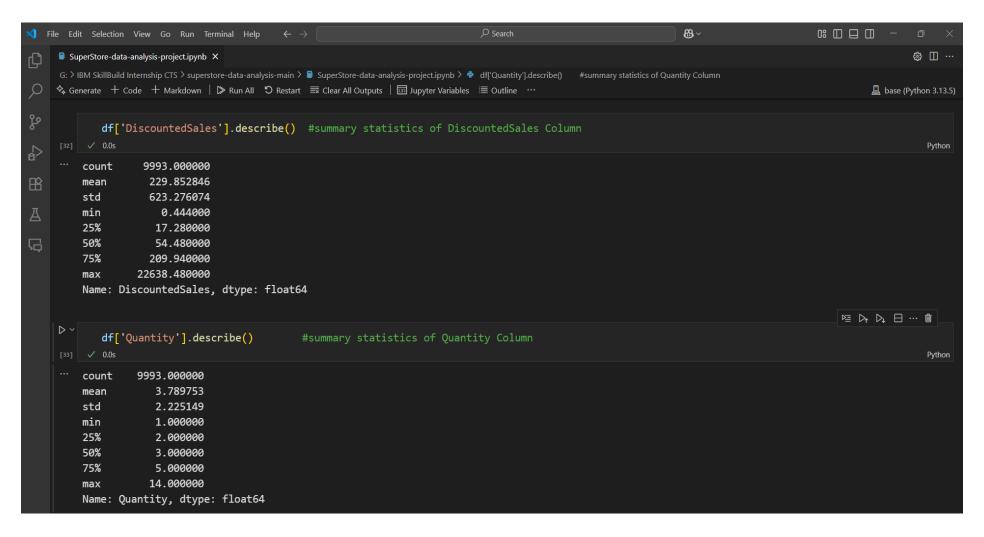




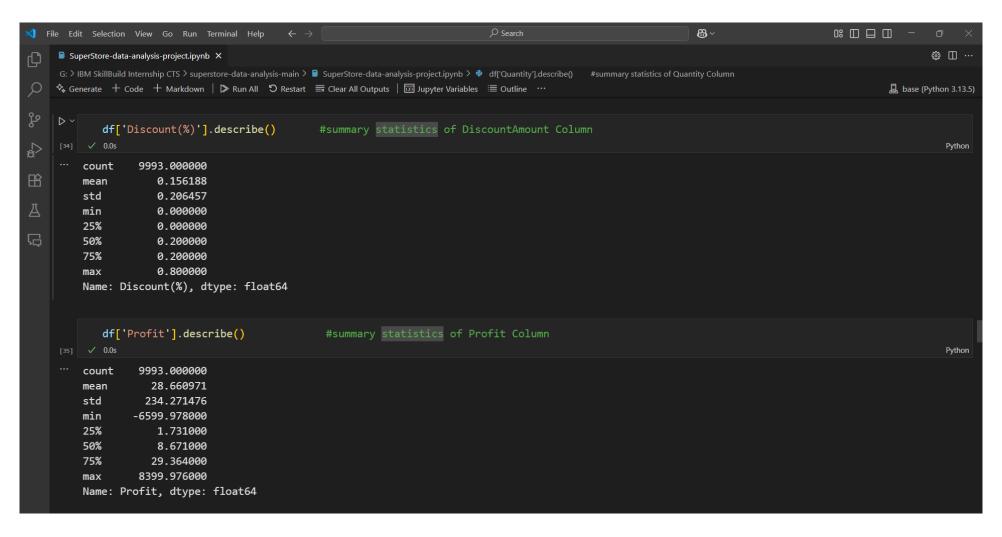




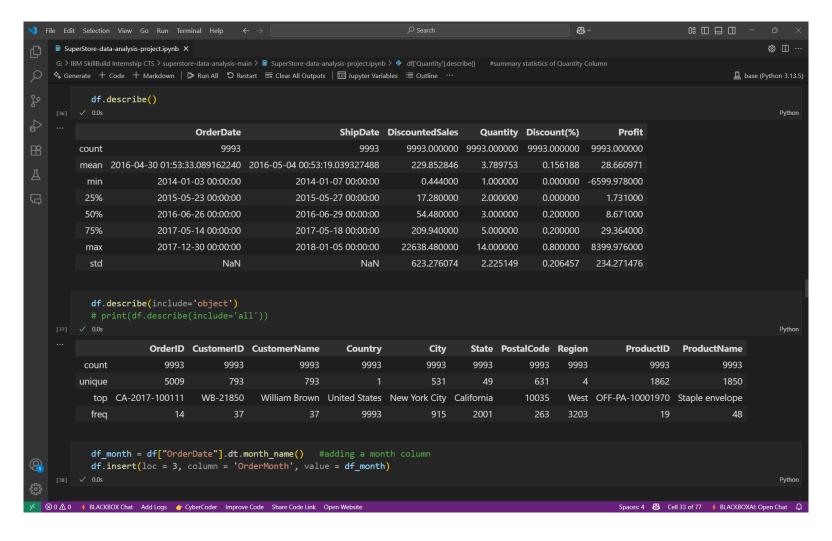




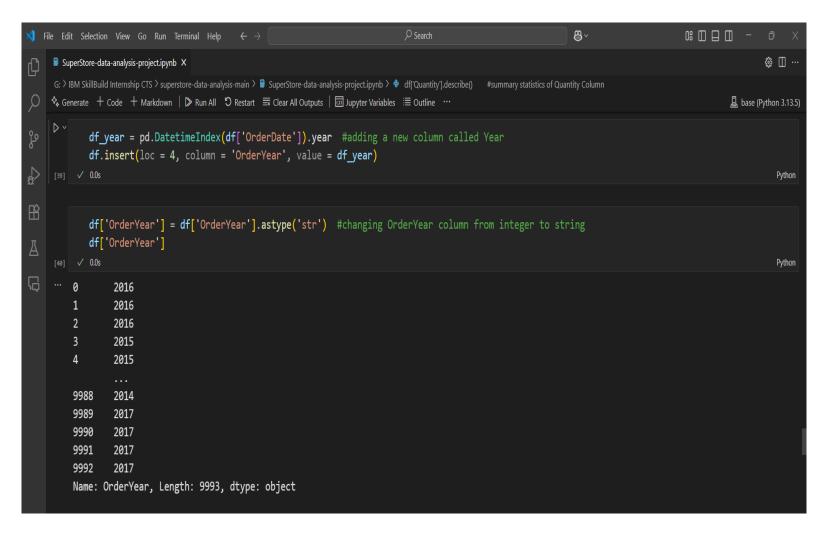












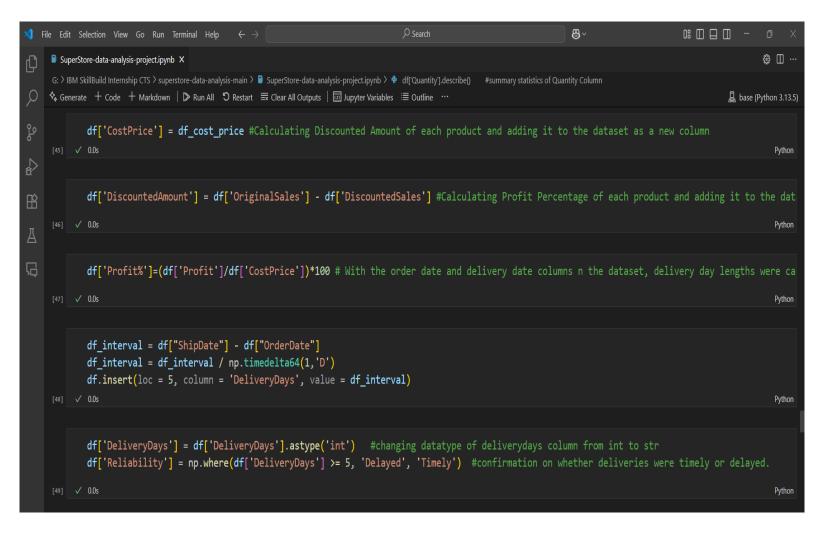


```
File Edit Selection View Go Run Terminal Help \leftarrow \Rightarrow
                                                                                                                             88 v
                                                                                                                                                                          ₩ 🛮 ..
  ■ SuperStore-data-analysis-project.ipynb X
   G: > IBM SkillBuild Internship CTS > superstore-data-analysis-main > 🛢 SuperStore-data-analysis-project.ipynb > 🏺 df['Quantity'].describe() #summary statistics of Quantity Column
 🔖 Generate + Code + Markdown | ▶ Run All 🤊 Restart 🗮 Clear All Outputs | 🖾 Jupyter Variables 🗎 Outline \cdots
                                                                                                                                                                  base (Python 3.13.5)
           original_sales = (1/(1-df['Discount(%)']))*df['DiscountedSales']
           original sales
                                                                                                                                                                             Python
                  261.96
                  731.94
                   14.62
                1741.05
                   27.96
       9988
                   31.56
       9989
                   91.96
                  323.22
       9990
       9991
                   29.60
                  243.16
        9992
       Length: 9993, dtype: float64
           df['OriginalSales'] = original_sales
```

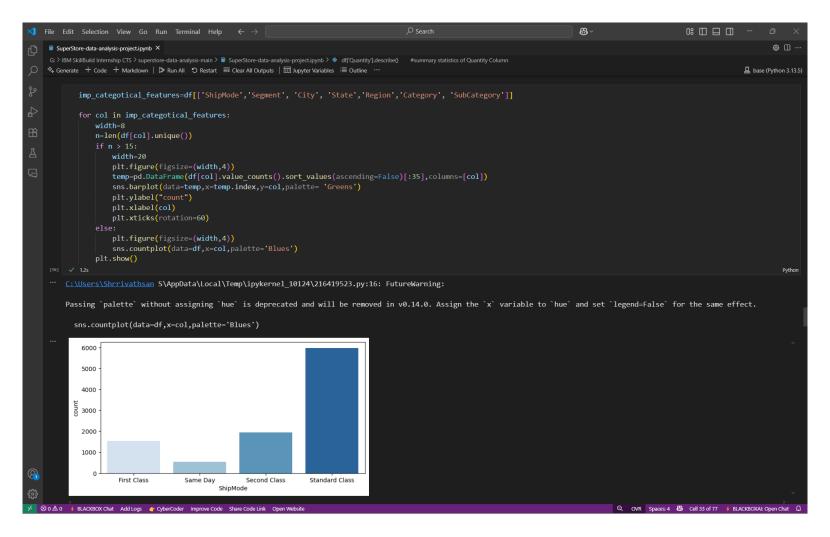


```
	imes File Edit Selection View Go Run Terminal Help 	extbf{\leftarrow} 	o
                                                                                                                                 & ~
                                                                                                                                                            ■ SuperStore-data-analysis-project.ipynb ×
                                                                                                                                                                               ₩ Ш …
     G: > IBM SkillBuild Internship CTS > superstore-data-analysis-main > 🛢 SuperStore-data-analysis-project.ipynb > 🍷 df['Quantity'].describe() #summary statistics of Quantity Column
    🍫 Generate 🕂 Code 🕂 Markdown | ⊳ Run All 💙 Restart 🚍 Clear All Outputs | 园 Jupyter Variables 🗏 Outline …
                                                                                                                                                                       base (Python 3.13.5)
             df['UnitPrice'] = original_sales/df['Quantity']
             df['UnitPrice']
                                                                                                                                                                                  Python
                   130.98
                   243.98
                     7.31
                   348.21
                    13.98
                    10.52
                    45.98
                   161.61
                     7.40
                   121.58
          Name: UnitPrice, Length: 9993, dtype: float64
             df cost price = df['DiscountedSales'] - df['Profit']
              df_cost_price
                     220.0464
                    512.3580
                      7.7486
                   1340.6085
                     19.8516
                     21.1452
                      76.3268
                     239.1828
                     16.2800
                    170.2120
          Length: 9993, dtype: float64
                                                                                                                                      OVR Spaces: 4 & Cell 33 of 77  
BLACKBOXAI: Open Chat Q
🔨 🛇 0 🛆 0 🦂 BLACKBOX Chat Add Logs 👉 CyberCoder Improve Code Share Code Link Open Website
```

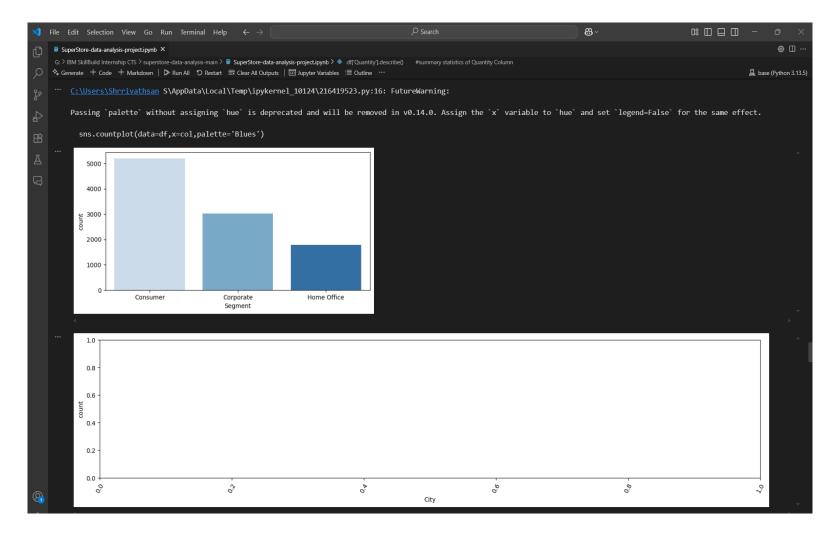




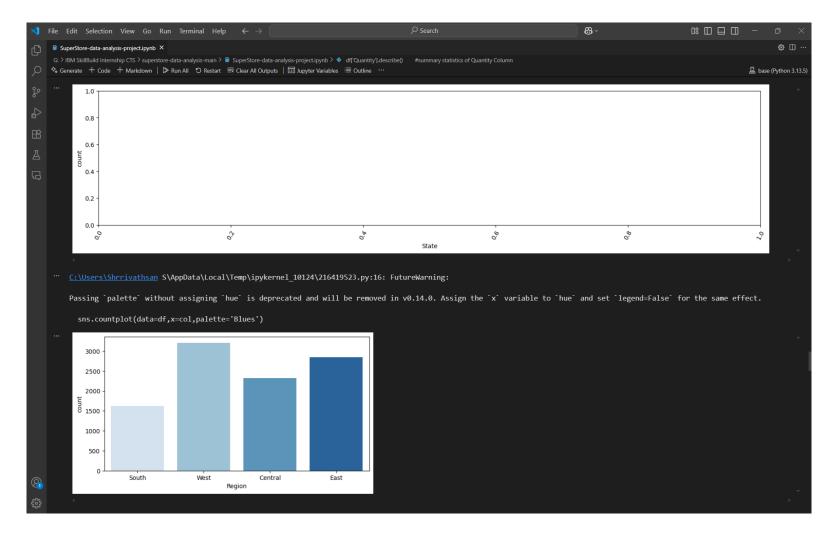




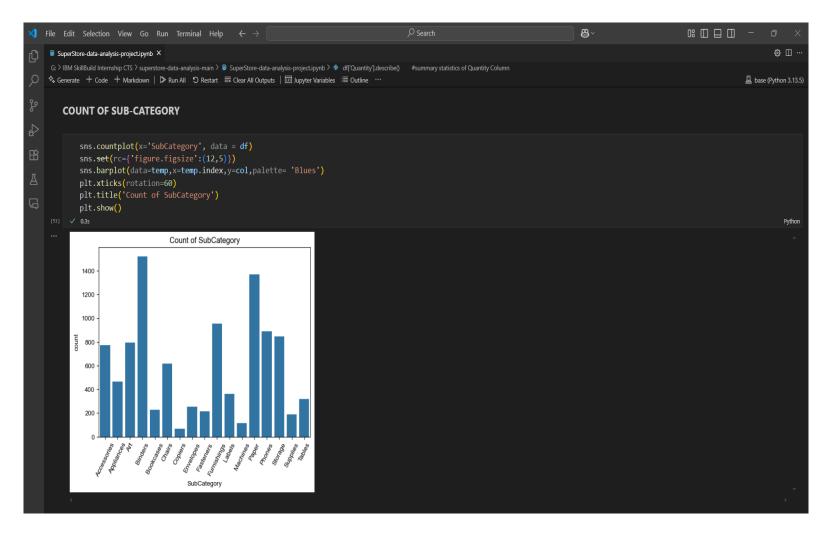




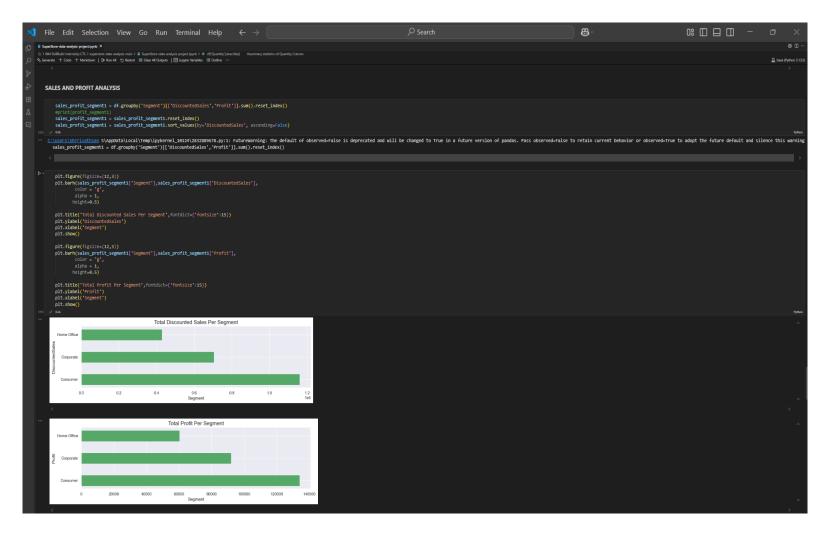




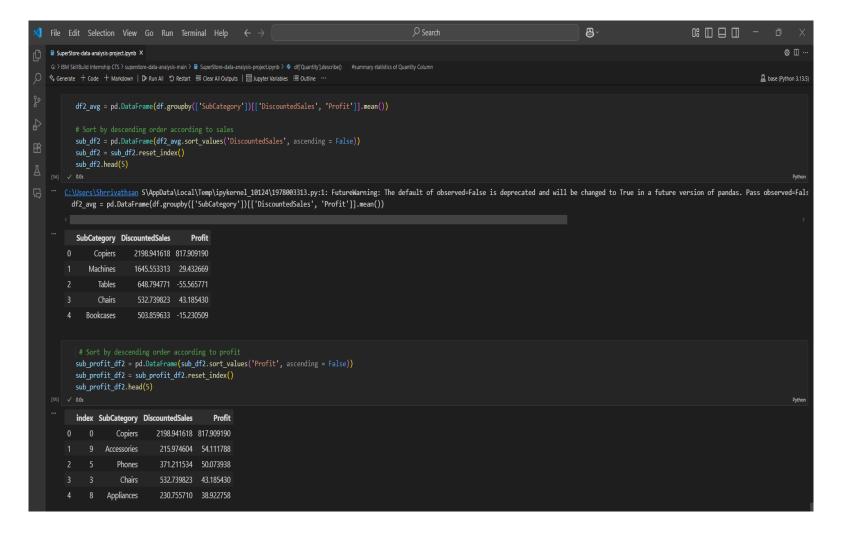




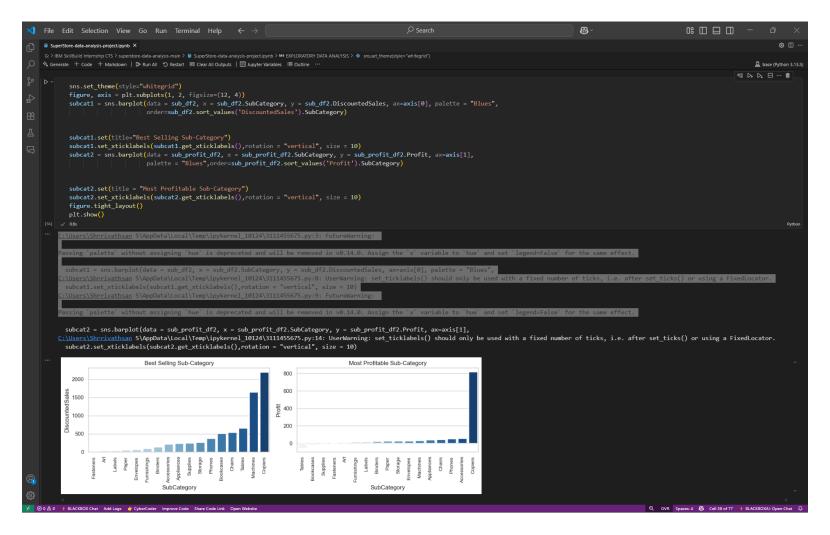




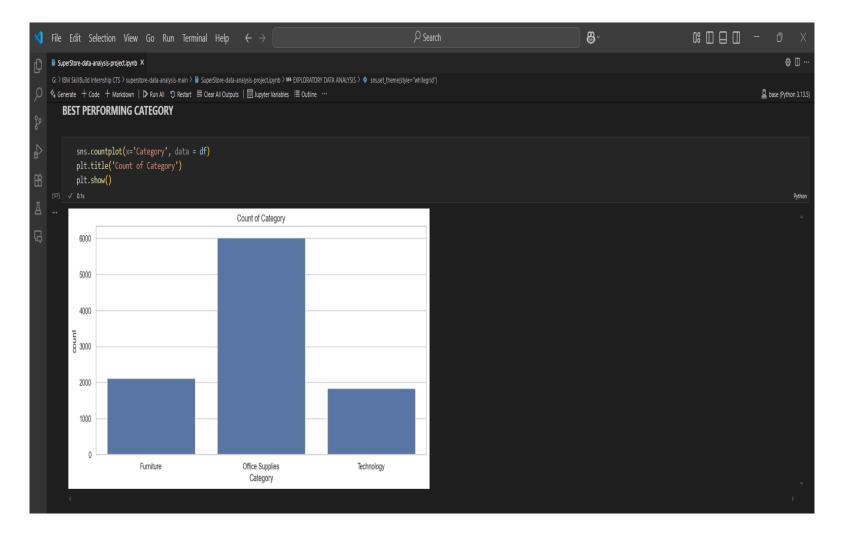




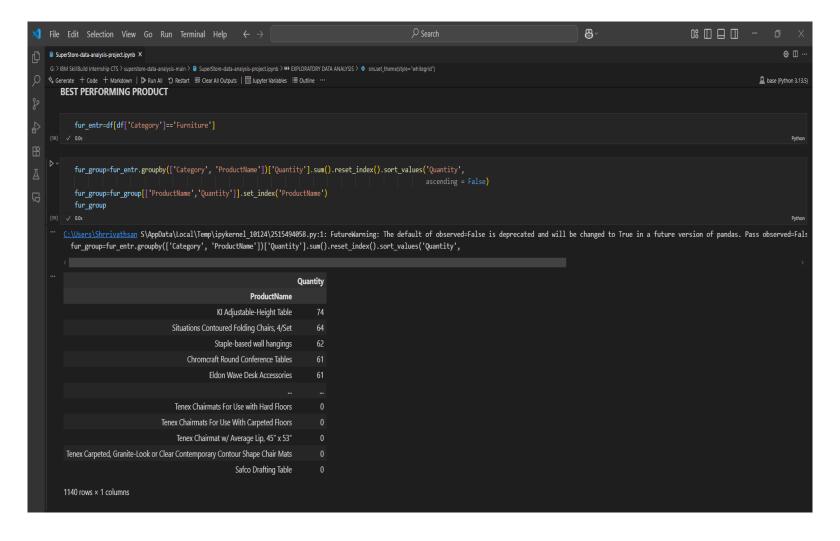




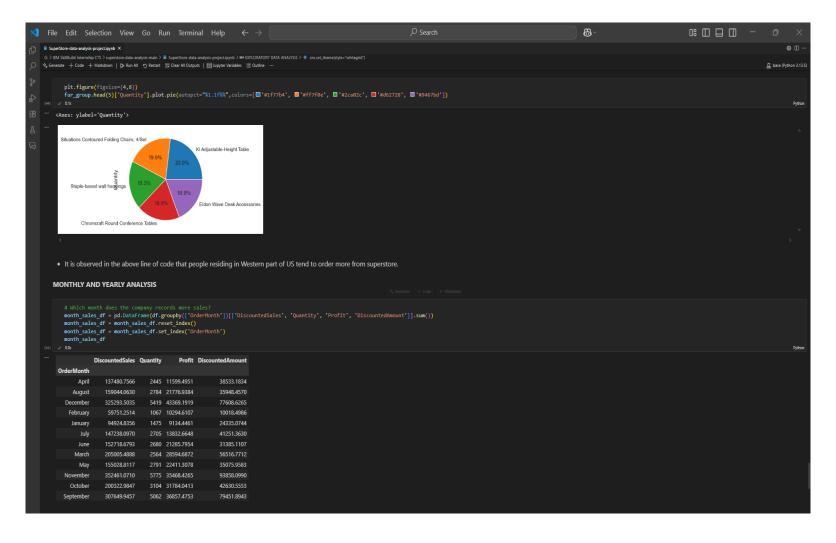




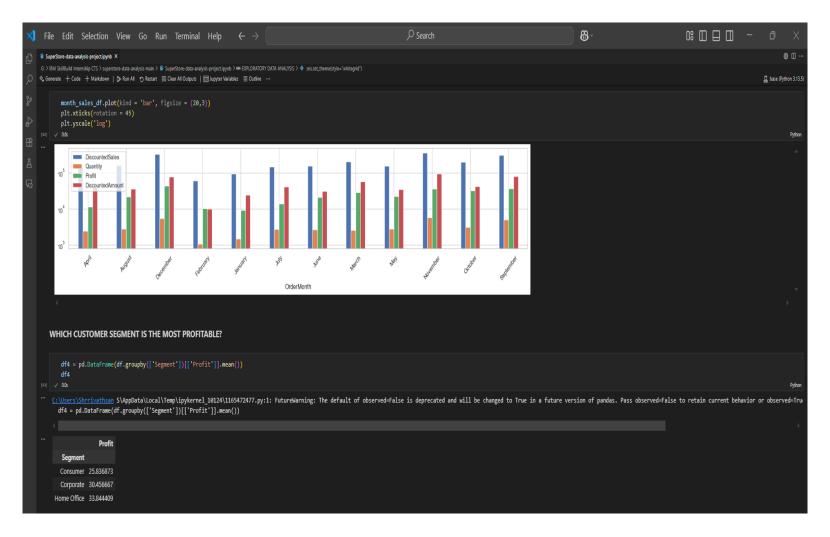




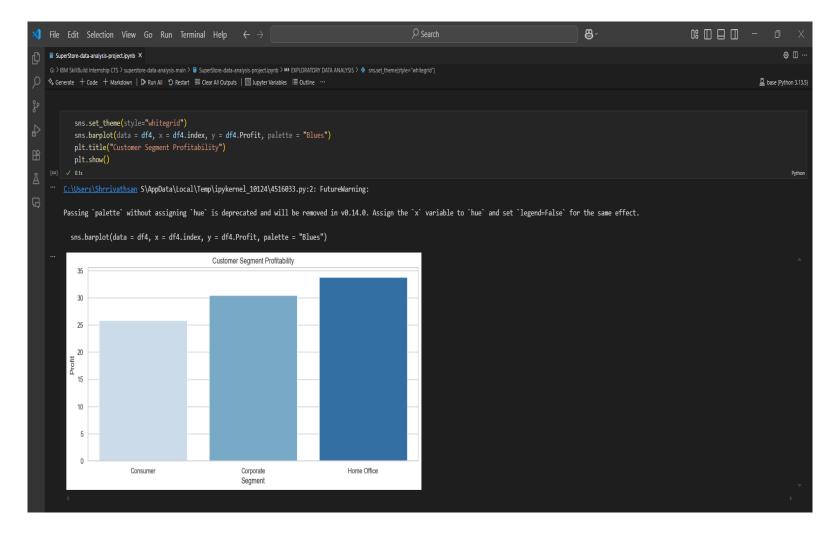




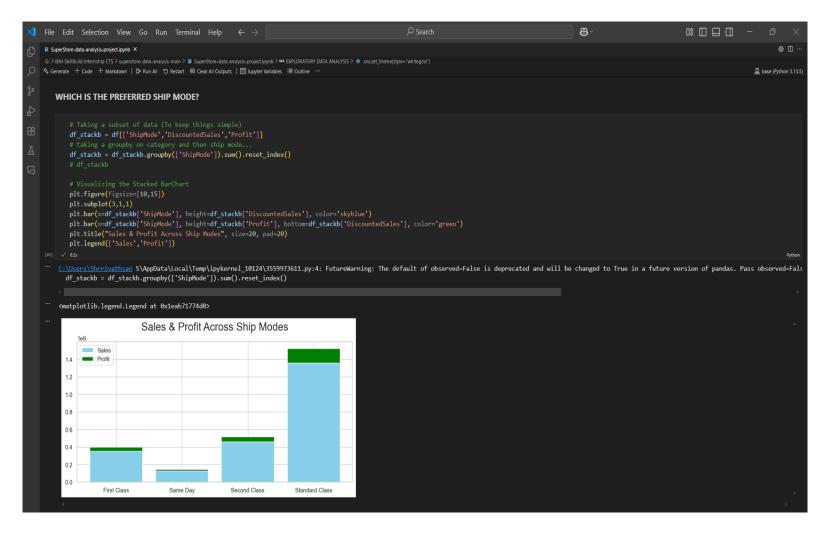




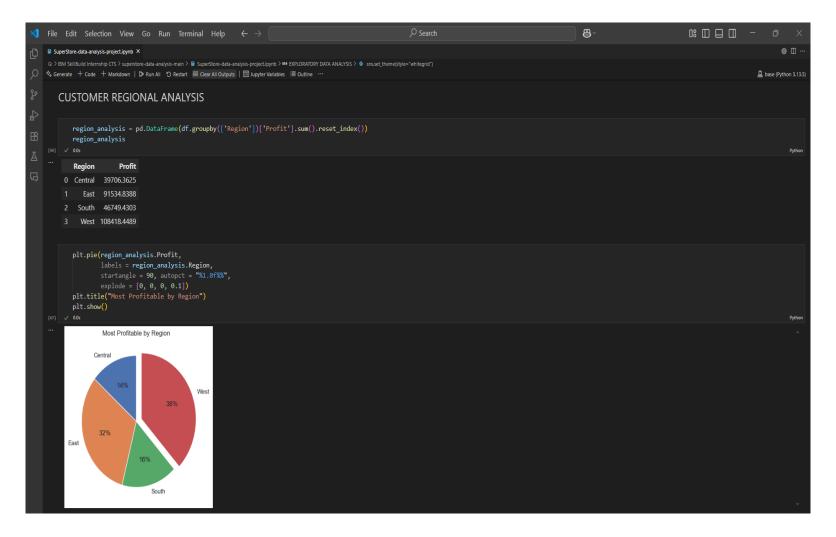














CONCLUSION

The SuperStore data analysis using AI and data science techniques provides deep insights into business performance, customer behavior, and operational bottlenecks. Through exploratory data analysis, machine learning models, and forecasting techniques, we have demonstrated how businesses can leverage data to make informed decisions.

Key takeaways include:

- Identification of high-performing product categories and loss-making areas.
- Discovery of customer segments through clustering, enabling targeted marketing.
- Accurate sales forecasting using time-series models to support inventory and resource planning.
- Use of predictive analytics to understand profit drivers and optimize discount strategies.
- The integration of AI methods such as machine learning, clustering, and forecasting has shown significant potential to improve efficiency, profitability, and customer satisfaction in a retail environment. Moreover, visual analytics and dashboards make these insights accessible and actionable for business stakeholders.
- With the growing accessibility of advanced AI tools, the future holds even more promise. Businesses can further evolve by incorporating real-time analytics, personalized recommendation engines, and geospatial intelligence into their decision-making processes.
- In conclusion, data science and Al not only empower SuperStore to understand its present but also to anticipate its future turning raw data into strategic advantage.



FUTURE SCOPE

As businesses increasingly adopt data-driven strategies, the scope of Al and data science in retail analytics is expanding rapidly. Below are some key future directions and enhancements that can be applied to the SuperStore analysis:

1. Advanced AI for Personalized Recommendations

Objective: Use customer purchase history, behavior, and segmentation data to recommend products.

Methods: Collaborative filtering, content-based filtering, or hybrid recommendation systems (e.g., matrix factorization, deep learning-based recommenders). Impact: Increase cross-selling, upselling, and customer satisfaction.

2. Real-Time Analytics & Streaming Data

Objective: Monitor sales, inventory, and customer behavior in real time.

Tools: Apache Kafka, Spark Streaming, and real-time dashboards using Power BI or Grafana.

Impact: Enables faster business decisions and immediate response to market trends.

3. Integration with ERP and CRM Systems

Objective: Combine SuperStore data with broader enterprise systems for holistic analytics.

Use Case: Unified dashboards for sales, supply chain, customer feedback, and finance.

Impact: Greater operational efficiency and 360° business insight.



FUTURE SCOPE

4. Predictive Inventory and Demand Planning

Objective: Predict future demand per region or product category using Al.

Techniques: Deep learning (LSTM/GRU), Prophet, or AutoML.

Impact: Reduce overstocking or stockouts, optimize warehouse management.

5. Geospatial Analysis

Objective: Analyze location-based sales patterns using maps and GIS data.

Tools: QGIS, Plotly, GeoPandas, Tableau Map Layers.

Impact: Optimize delivery routes, open new stores in high-potential regions.

6. Sentiment Analysis from Customer Reviews

Objective: Extract customer sentiment from feedback or reviews.

Techniques: NLP using BERT, RoBERTa, or LLMs (ChatGPT API for summarization).

Impact: Improve product quality, customer service, and retention.



FUTURE SCOPE

7. Al Chatbots for Sales & Support

Objective: Deploy Al chatbots trained on SuperStore data to guide customers.

Use Case: Automated helpdesk, guided shopping assistants.

Tools: OpenAl, Dialogflow, Rasa.

Impact: Better customer engagement and reduced support costs.

8. Ethical Al & Bias Detection

Objective: Ensure fairness and explainability in Al models (e.g., customer segmentation, credit scoring).

Tools: SHAP, LIME, AI Fairness 360.

Impact: Trustworthy Al systems and regulatory compliance (e.g., GDPR).

9. Augmented Reality (AR) & Visual AI for Retail

Objective: Visual product recommendations and in-store AR experiences.

Techniques: Computer vision, AR SDKs, Visual Search.

Impact: Next-generation retail experience.



REFERENCES

- ■ Books:
- 1. "Data Science for Business" by Foster Provost & Tom Fawcett
- 2. "Python for Data Analysis" by Wes McKinney
- 3. "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow" by Aurélien Géron
- Research Papers:
- 1. Aggarwal, C. C. (2015). Data Mining: The Textbook Comprehensive methods in clustering, classification, and forecasting.
- 2. Hyndman, R. J., & Athanasopoulos, G. (2018). Forecasting: Principles and Practice For ARIMA and time series modeling.
- 3. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system Use for boosting models on tabular data.



IBM CERTIFICATIONS

In recognition of the commitment to achieve professional excellence



Shrrivathsan S

Has successfully satisfied the requirements for:

Getting Started with Artificial Intelligence



Issued on: Jul 16, 2025 Issued by: IBM SkillsBuild

Verify: https://www.credly.com/badges/dd6aafea-ee76-47bb-aaeb-7b4a4f7bb010





IBM CERTIFICATIONS

In recognition of the commitment to achieve professional excellence Shrrivathsan S Has successfully satisfied the requirements for: Journey to Cloud: Envisioning Your Solution Issued on: Jul 16, 2025 Issued by: IBM SkillsBuild Verify: https://www.credly.com/badges/20c87343-62c2-4132-bedb-6f8181b7200e



IBM CERTIFICATIONS

IBM SkillsBuild

Completion Certificate



This certificate is presented to

Shrrivathsan S

for the completion of

Lab: Retrieval Augmented Generation with LangChain

(ALM-COURSE_3824998)

According to the Adobe Learning Manager system of record

Completion date: 19 Jul 2025 (GMT)

Learning hours: 20 mins



THANK YOU

