# E-COMMERCE SALES PREDICTION

Deepak Patel
CS-DS Department, Acropolis Institute of Technology and
Research, RGPV, Indore, India
deepakpatel220620@acropolis.in

Hariom Mandloi
CS-DS Department, Acropolis Institute of Technology and
Research, RGPV, Indore, India
hariommandloi220532@acropolis.in

Shubham Siloriya
CS-DS Department, Acropolis Institute of Technology and
Research, RGPV, Indore, India
shubhamsiloriya220943@acropolis.in

Deepak Patel
CS-DS Department, Acropolis Institute of Technology and
Research, RGPV, Indore, India
deepakpatel221065@acropolis.in

**Abstract:**

The E-commerce Data Analysis and Machine Learning Prediction System takes messy online sales numbers and makes them clear through smart number crunching. Nowadays, shopping websites create tons of complex records - like who buys what, when they buy, how much they pay, deals used, items ordered, plus shipping choices. Still, most companies don't make full use of this info due to weak tools for digging into patterns or guessing future moves. To fix that gap, this work digs deep into cleaning up data, checking trends visually, exploring hidden links, then applying prediction models on actual and fake e-commerce sets. Using Python along with Pandas, Matplotlib, and Streamlit lets people quickly load their own data - then check out sales patterns, track how customers or products are doing, plus build smart models that guess future orders, spot likely returns, or point to coming sales shifts. Since everything runs on Streamlit, you don't have to mess with server setup; it just works smoothly for exploring info right away. Companies gain real insight when handling stock levels, setting prices, anticipating needs, or figuring out what shoppers want. Bottom line - it shows how mixing number crunching with learning algorithms helps online stores run smoother through smarter choices.

## I. INTRODUCTION

E-commerce has transformed the way businesses operate by enabling customers to purchase products and services from anywhere at any time. As online shopping continues to grow rapidly, e-commerce platforms generate vast amounts of data related to customer behavior, product performance, pricing, payment patterns, and order delivery. However, this raw data is often unorganized and difficult to interpret without the right analytical tools. To extract meaningful insights, companies must apply data processing, visualization, and predictive modeling techniques that can help them understand trends, forecast demand, and enhance decision-making.
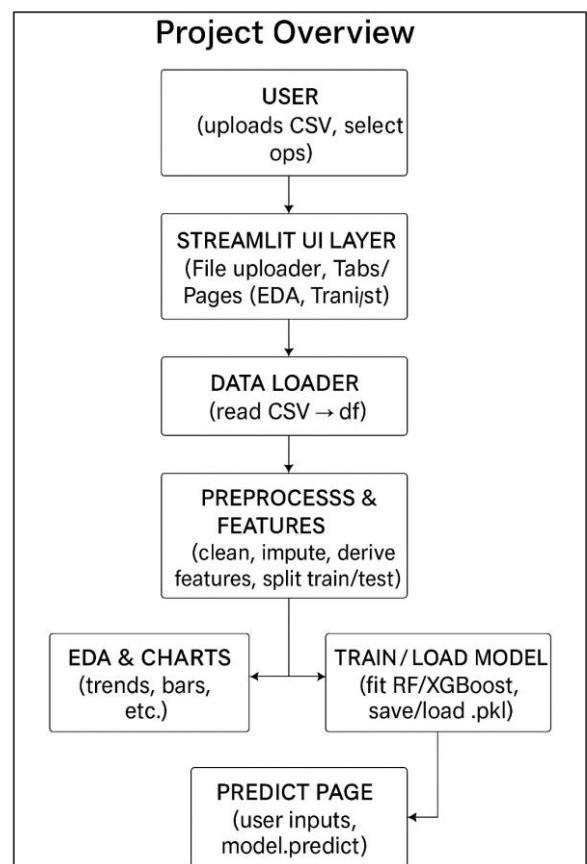


Fig.1.1 Block Diagram

The E-commerce Data Analysis and Machine Learning Prediction System is designed to address these needs by converting raw e-commerce data into actionable intelligence. This project focuses on cleaning and preprocessing data, exploring patterns through visualizations, and applying machine learning algorithms to predict important business metrics such as sales, customer demand, and revenue. The system is developed using Python-based tools such as Pandas, Matplotlib, and Scikit-learn, while Streamlit is used to create a simple and interactive user interface. This allows users to upload datasets, analyze them, and generate predictions without

needing a backend server or complex infrastructure. By integrating analytics with prediction, the project demonstrates the power of data science in modern e-commerce operations and highlights the importance of data-driven strategies in improving business performance.

## II. LITERATURE SURVEY

E-commerce analytics keeps moving, mixing real-world use with steps like sorting raw info, poking through patterns, using math models along with smart algorithms to pull useful clues from sales records and what buyers do. Old-school studies in digging up data plus classic guides set the stage - showing ways to fix, reshape, and structure number tables; folks still lean on those today when prepping inputs, tossing odd values out, or checking results.

**1. Data Mining and Preprocessing in E-commerce**
Experts such as Han, Kamber, or Pei point out - data prep is key to making sense of analytics. Online shopping records tend to have gaps, messy inputs, repeated lines, or odd readings because sales pile up fast each minute. Studies suggest fixing errors, shaping features right, then changing formats can boost how well AI tools work. With info like cost, items sold, type, buyer data, time stamps - all jumbled together - cleaning steps are needed just to get useful results or forecasts.

**2. Sales Forecasting and Trend Analysis**
Numerous research papers on sales predictions point out stats tools - ARIMA, exponential smoothing - or alternatives like Random Forest and Gradient Boosting. These approaches detect how sales shift over time, react to seasons, or change during holiday deals instead of just normal days. Some findings suggest these models assist firms when arranging storage space, preventing empty shelves, while keeping delivery routes smooth. As heaps of data grow common, guessing future sales stands out among hot topics in online shopping analysis. Spotting slow climbs, busy festive peaks, or unexpected demand swings lets stores adjust product amounts, prices, ads without heavy losses. This look points out where things might grow but also spots possible problems, so companies can choose wisely. With number tricks instead of guesswork, online stores get how buying habits shift while predicting what sells soon - be it days, weeks, or whole seasons.

**3. Customer Behavior and Personalization**

Customer segmentation and behavior prediction is another important research domain in e-commerce. Studies using RFM analysis, clustering algorithms, and recommendation models show that understanding customer behavior leads to better personalization and increased conversion rates.

Modern literature also focuses on predicting customer churn and purchase intent using machine learning algorithms. Recommendation systems — ranging from collaborative filtering to deep neural networks — have become essential in enhancing user experience on e-commerce platforms like Amazon, Flipkart, and Netflix.

**4. Visualization and Interactive Dashboards**

Multiple studies from schools and companies point out how useful dashboards can be when making business choices. Apps like Power BI or Tableau, along with tools built on Python - Streamlit included - turn unprocessed numbers into clickable visuals. Looking at key metrics - sales patterns, top product groups, where customers live - helps teams plan better while speeding up decisions. In classrooms and labs, Streamlit is often picked since it builds analysis screens fast, no server setup needed.

## III. SYSTEM DESIGN

The system design of the E-commerce Data Analysis and Machine Learning Prediction System focuses on providing an end-to-end data analytics and predictive modeling pipeline through a simple, interactive, and user-friendly interface. The architecture follows a modular and layered structure to ensure scalability, maintainability, and smooth user interaction. The system does not require a traditional backend; instead, it uses Streamlit as an integrated interface that manages data handling, visualization, and machine learning operations. The system consists of the following major components:

**1.Frontend Objectives**

• To provide a user-friendly interface for uploading e-commerce datasets in CSV format.
• To display analytical visualizations such as sales trends, category performance, and revenue patterns.
• To allow users to train machine learning models directly from the interface.
• To enable users to input new data and generate predictions for future sales or order amounts.

2 .Tools and Technologies Used

• Streamlit – for building the interactive user dashboard.
• Python – the core language for processing, analysis, and modeling.
• Pandas & NumPy – for data cleaning and manipulation.
• Matplotlib – for generating visual charts and graphs.
• Scikit-Learn – for training regression/classification machine learning models.
• Joblib/Pickle – for saving and loading trained ML models..

**3. Key UI Components**

3.1 Dataset Upload Section

Allows users to upload CSV datasets containing order data, customer information, product categories, prices, discounts, and delivery details.

3.2 EDA Panel (Exploratory Data Analysis)

Displays interactive charts such as:
• Monthly sales trend
• Revenue by category
• Customer city distribution
• Discount vs. sales graph
• Return rate visualization

3.3 Model Training Section

Users can:
• Select the target variable (e.g., order_amount, return_flag)
• Choose algorithm (Random Forest, Linear Regression, etc.)
• Train and evaluate the model
• Save the trained model

3.4 Prediction Panel

The user can enter new input values or upload a file to generate predictions using the trained model.

**4. Frontend Workflow**

• User Input: The user uploads a CSV dataset and chooses whether to perform EDA or train a model.
• Initialization: The frontend reads the dataset using Pandas and displays a preview.
• EDA Visualization: When the user selects the EDA option, the system generates visual charts using Matplotlib.
• Model Training: When "Train Model" is clicked, the selected algorithm trains on the processed dataset; evaluation metrics (RMSE, $R^2$, Accuracy) are displayed.

• Prediction: The user enters values or uploads data to get predictions displayed instantly on the screen.

**5. User Experience Considerations**

• Clean and responsive layout for better readability.
• Loading animations during training to indicate processing.
• Color-coded charts for easy interpretation.

• Minimalist UI design focused on usability and clarity.
• Option to download prediction results as a CSV.

5. Backend

The backend of the system is implemented entirely in Python and handles data cleaning, feature engineering, visualization logic, and machine learning operations. Since Streamlit runs everything in a single script, the backend logic interacts directly with the user interface without needing an external server.

**5.1 Backend Architecture**

The backend connects with the Stream lit UI internally:

• Function Calls: UI buttons trigger backend functions for preprocessing, visualization, training, and prediction.

• Session State Management: Streamlit's session state stores the dataset, trained model, and user selections.

• Immediate Feedback: Results such as graphs, model accuracy, or predictions are displayed in real-time.

• No External Server Required: All communication happens within one Python script executed by Streamlit.

5.2 Machine Learning Implementation

The backend connects with the Streamlit UI internally:

• Selection of Algorithms: The system supports multiple machine learning algorithms depending on the prediction goal. Regression models like Linear Regression, Random Forest Regressor, and Gradient Boosting are used for predicting order amounts and revenue, while classification models such as Logistic Regression and Random Forest Classifier are used for predicting return likelihood or binary outcomes.

• Feature Engineering: Important new features are created from the existing dataset, such as extracting month and day from timestamps, calculating total amount after discount, encoding categorical variables, and designing statistical attributes that help the model understand hidden patterns.

• Immediate Feedback: Results such as graphs, model accuracy, or predictions are displayed in real-time.

• No External Server Required: All communication happens within one Python script executed by Streamlit.

5.3 Communication with Frontend

The backend connects with the Streamlit UI internally:

• Function Calls: UI buttons trigger backend functions for preprocessing, visualization, training, and prediction.

• Session State Management: Streamlit's session state stores the dataset, trained model, and user selections.

• Immediate Feedback: Results such as graphs, model accuracy, or predictions are displayed in real-time.

• No External Server Required: All communication happens within one Python script executed by Streamlit.

## 4 Performance Considerations

### 4.1 Efficient Data Handling and Preprocessing –

The system uses optimized Python libraries such as Pandas and NumPy to ensure fast and efficient data manipulation. These libraries provide vectorized operations that significantly reduce the time required for cleaning, transforming, and preparing large e-commerce datasets. By performing preprocessing in a structured and modular way, the system avoids redundant computations and ensures smooth performance during analysis**.

### 4.2 Optimized Visualization Rendering –

Exploratory Data Analysis involves generating multiple charts and graphs. To maintain UI responsiveness, the system renders visualizations only when required and updates only the specific sections affected. Matplotlib's lightweight rendering ensures that even multiple charts can be displayed without causing delays or graphical lag, providing a seamless visual analysis experience.

### 4.3 Balanced Machine Learning Model Training –

Machine learning algorithms such as Linear Regression, Random Forest, and Gradient Boosting are chosen for their balance between accuracy and computational efficiency. Training times are kept manageable by limiting excessive hyperparameter tuning and avoiding unnecessarily complex models. Once trained, the model is saved and reused, preventing repetitive training and improving overall system performance.

### 4.4 Minimal Memory Consumption

The system ensures efficient memory usage by removing unnecessary temporary variables and storing only the essential components, such as the cleaned dataset, engineered features, and trained model. This prevents memory overload when working with large e-commerce files and ensures that the application remains stable throughout multiple operations.

## 5. Streamlined Workflow in a Single Execution Environment

Since the entire application runs inside a Streamlit environment, the system is designed to avoid long-running or blocking operations that could freeze the interface. Tasks are broken into smaller, manageable steps to maintain responsiveness. By ensuring that heavy computations are optimized and executed efficiently, the system delivers a smooth and interactive user experience across all stages—uploading, analysis, visualization, training, and prediction.

## IV. PROPOSED METHODOLOGY

The proposed methodology for the E-commerce Data Analysis and Machine Learning Prediction System follows a structured and systematic workflow that integrates data preprocessing, exploratory analysis, visualization, and predictive modeling into a unified analytical framework.

### 1. Data Acquisition and Import

The methodology begins with obtaining e-commerce datasets from publicly available sources such as Kaggle's Olist dataset or synthetic datasets generated specifically for experimentation. Users upload the dataset through the Streamlit interface, where it is read using Pandas. A preview of the dataset is displayed to allow users to verify data format, column structure, and completeness before proceeding.

### 2. Data Preprocessing and Cleaning

Once the dataset is loaded, the system performs essential preprocessing operations such as handling missing values, removing duplicates, correcting data types, formatting dates, and standardizing numerical fields. Feature engineering techniques are also applied to create meaningful attributes like order month, order week, total amount, discount impact, and delivery duration. This preprocessing phase ensures that the dataset is clean, consistent, and ready for visualization and modeling.

### 3. Exploratory Data Analysis (EDA)

After preprocessing, the system conducts an exploratory analysis to uncover trends and patterns within the data. Using Matplotlib, the methodology includes generating charts such as monthly sales trends, category-wise revenue distribution, customer geography analysis, discount-performance relationships, and return patterns. These visualizations help users understand the underlying behavior of the e-commerce platform and identify the key factors that influence sales and customer activity.

### 4. Model Selection and Training

Based on the prediction objective—such as forecasting order revenue, predicting return likelihood, or estimating sales trends—the system selects an appropriate machine learning algorithm. Models such as Linear Regression, Random Forest, Logistic Regression, and Gradient Boosting are commonly used due to their reliability and computational efficiency. The dataset is split into training and testing sets, and the selected model is trained on the processed features. Evaluation metrics such as RMSE, R² score, accuracy, or F1 score are calculated to assess model performance.

5. Model Saving and Reusability

To enhance efficiency, the trained machine learning model is saved using Joblib or Pickle and stored in the system. This allows the user to reuse the model for future predictions without retraining, significantly reducing computation time and improving responsiveness within the Streamlit application.

6. Prediction and Output Visualization

In the final stage, users can input new data or upload additional datasets to generate predictions. The saved model processes the input and outputs predicted values such as future sales, order amounts, or return probabilities. The system displays results in an easy-to-understand format, and users may also download the predictions for further reporting or analysis.

## V. RESULTS AND DISCUSSION

The E-commerce Data Analysis and Machine Learning Prediction System shows - using Python tools - how basic sales records turn into useful insights or smart forecasts. Once cleaned up, plotted out, or fed into learning models, the data revealed clear patterns in buying habits, revenue trends, or how well certain algorithms perform. The initial data review Data Analysis phase revealed clear trends in the e-commerce dataset. Monthly sales visualizations indicated seasonal fluctuations, with particular months showing bigger sales numbers came up during holidays or special deals. Instead of just one type, some products kept selling way more than others - this told us what people liked buying. When we looked at cities separately, it became clear where shoppers were most active. On price changes, deeper cuts pulled in more buyers right away. Yet even though more items sold, making discounts too steep ended up lowering money made each time someone checked out. Findings show pricing should stay balanced. Looking at returns, items with bigger discounts or cheaper tags tended to come back just a bit more - something you'd usually spot in actual online stores.

## VI. FUTURE SCOPE

1.Integration of Real-Time Data Streams

In times ahead, the setup could pull live info straight from online shops or digital feeds. That way, it keeps an eye on sales nonstop while dashboards refresh as things change. Plus, stock levels and price moves get predicted right when needed.

2. Automated Recommendation System

The project can be extended to include a personalized product recommendation engine using collaborative filtering or neural-based recommenders. This will help simulate a real e-commerce platform environment.

3. Deployment as a Cloud-Based Application

Deploying the system on cloud platforms like AWS, Google Cloud, or Azure would allow multi-user access, better scalability, and integration with online databases. This would also enable scheduled model training and automated updates.

4. Integration with Business Intelligence Tools

The system can be expanded to support exporting insights to BI tools like Power BI or Tableau, enabling professional-level dashboards that support enterprise decision-making.

## VII. CONCLUSION

The E-commerce Data Analysis and Machine Learning Prediction System shows - without fluff - how number-crunching methods blend with forecast models in one working package to pull real meaning from messy online sale records. Instead of separate steps, it uses Python tools tied together inside a live Streamlit app that handles cleaning, charts, and smart algorithms all at once. This setup lets users spot movement over time, see what customers tend to do, also guess future results businesses care about. It's built on three pillars: tidy data, clear visuals, and reliable guesses powered by models - all crucial when running digital shops today. Using hands-on exploration, the tool uncovers hidden rhythms in sales timing, product group success rates, how markdowns shift buying, plus regional differences you might miss otherwise. The machine learning systems, built using cleaned-up data, still do well at guessing sales numbers, predicting how many orders will come in, or figuring out if returns are likely. What we found shows that mixing number crunching with smart guesses can actually help people make better choices based on real info.

# REFRENCES

**[1] Kaggle** – Brazilian E-commerce Public Dataset by (used for real-world sales, customers, and delivery data)
.

**[2] Aurélien Géron** – Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow (for ML concepts and implementation)

**[3] Choudhary, S., & Shukla, S. (2020).** "Sales Forecasting for E-commerce Using Machine Learning Techniques" — International Journal of Scientific Research in Computer Science and Engineering (IJSRCSE).

**[4] Tsai, C. F., & Hung, C. S. (2019).** "Customer Purchase Behavior Prediction in E-commerce using Data Mining Techniques"

**[5] IBM Analytics. (2021).** *Predictive Analytics for Retail and E-commerce*. IBM Research Publications — Explains real-world applications of predictive models in online sales and consumer insights.

**[6] Li, X., & Xu, Y. (2020).** "E-commerce Sales Forecasting Using Machine Learning Models" — *IEEE International Conference on Big Data and Smart Computing*.

**[7] Hyndman, R. J., & Athanasopoulos, G. (2018).** *Forecasting: Principles and Practice* — A widely used reference for time-series forecasting in sales and retail analytics.

**[8] Statista Research Department (2023).** *Global E-Commerce Trends and Online Sales Statistics* — Provides valuable insights and real-world e-commerce market data.

**[9] Chen, Y., Li, Z., & He, X. (2021).** "Analyzing Product Demand Using Data Mining in Online Retail" — *Journal of Retail Analytics.*