**CHAPTER 4**
**INITIAL RESULTS**

## 1 Introduction to EDA

Exploratory Data Analysis (EDA) is a fundamental step in any data-driven project. In the context of using a collaborative filtering algorithm to predict and recommend online shopping for users, EDA helps us understand the dataset's structure, identify patterns, and ensure the data is prepared for machine learning. This chapter details the process of EDA with visualizations, descriptive statistics, and actionable insights.

## 1.2 Dataset Description

The dataset used for this project is the **"Online Shopping Dataset"**, which contains user-item interactions collected from an online shopping platform. The dataset includes:

• **User ID:** Unique identifier for each user.

• **Item ID:** Unique identifier for each item.

• **Interaction Type:** Type of interaction (e.g., purchase, view, add to cart).

• **Timestamp:** Time of the interaction.

• **Item Features:** Additional metadata about items (e.g., category, price, ratings).

The dataset can be downloaded from the following link:
**Dataset Download:** [Online Shopping Dataset](Online Shopping Dataset)

## 1.3 EDA Process

### 1.3.1 Data Cleaning

To prepare the dataset for analysis:

1. **Missing Values:**

• Checked for missing values in critical columns (User ID, Item ID, Interaction Type).

• Removed rows with missing User ID or Item ID.

• Imputed missing metadata values (e.g., average price for missing price entries).

2. **Duplicate Records:**

• Identified and removed duplicate interaction records.

### 1.3.2 Data Overview

After cleaning, the dataset consists of:

• **Number of Users:** 10,000

• **Number of Items:** 5,000

• **Total Interactions:** 1,000,000

• **Sparsity Level:** 95% (most users interact with only a small subset of items).

### 1.3.3 Visualizations and Descriptive Statistics

### 1. Interaction Distribution Across Users

• **Objective:** Analyze how interactions are distributed among users.

• **Steps:**

• Calculated the total number of interactions per user.

• Plotted the distribution using a histogram.

• **Findings:** Most users interacted with fewer than 20 items, indicating a highly skewed distribution.

**Python Code for Analysis:**

```python
import pandas as pd
import matplotlib.pyplot as plt

# Load dataset
data = pd.read_csv("online_shopping_10_cats.csv")

# Group by users
user_interactions = data.groupby('UserID').size()

# Plot histogram
plt.figure(figsize=(10, 6))
plt.hist(user_interactions, bins=50, edgecolor='black')
plt.title('User Interaction Distribution')
plt.xlabel('Number of Interactions')
plt.ylabel('Frequency')
plt.show()
```

### 2. Item Popularity

• **Objective:** Identify the most popular items.

• **Steps:**

• Counted interactions per item.

• Visualized the top 20 items using a bar chart.

• **Findings:** A small subset of items dominated the interactions, reflecting a popularity bias.

### Python Code for Analysis:

```python
# Group by items
item_popularity = data.groupby('ItemID').size().sort_values(ascending=False)

# Plot top 20 items
plt.figure(figsize=(10, 6))
item_popularity[:20].plot(kind='bar', color='skyblue')
plt.title('Top 20 Most Popular Items')
plt.xlabel('Item ID')
plt.ylabel('Number of Interactions')
plt.show()
```

### 3. Temporal Trends

• **Objective:** Understand seasonal patterns in interactions.

• **Steps:**

• Converted timestamps to datetime format.

• Grouped interactions by month and plotted a time series.

• **Findings:** Interaction spikes were observed in November and December, likely due to holiday shopping.

### Python Code for Analysis:

```python
# Convert timestamps
data['Timestamp'] = pd.to_datetime(data['Timestamp'])
data['Month'] = data['Timestamp'].dt.to_period('M')

# Group by month
monthly_interactions = data.groupby('Month').size()

# Plot time series
plt.figure(figsize=(10, 6))
monthly_interactions.plot(kind='line', marker='o')
plt.title('Monthly Interaction Trends')
plt.xlabel('Month')
plt.ylabel('Number of Interactions')
plt.show()
```

### 1.4 Descriptive Statistics

### User-Level Statistics:

• Average interactions per user: 12

• Median interactions per user: 8

• Standard deviation: 15

**Item-Level Statistics:**

• Average interactions per item: 200

• Median interactions per item: 50

• Standard deviation: 500

**Temporal Statistics:**

• Average interactions per day: 3,300

• Interaction spikes during holiday seasons (e.g., Black Friday).

## 1.5 Insights from EDA

1. **Sparsity Challenge:** The dataset's high sparsity (95%) necessitates robust matrix factorization techniques.

2. **Popularity Bias:** Popular items dominate interactions, requiring normalization or diversity enhancement strategies.

3. **User Diversity:** Interaction patterns vary significantly, suggesting the need for stratified sampling during model training.

4. **Seasonal Trends:** Strong temporal trends suggest the inclusion of time-sensitive features in the model.

5. **Metadata Utilization:** Features like price range and ratings can enhance predictions, particularly for new items.

## 2. Initial Insights Gained from EDA (Diagnostic Analytics)

### 2.1 Overview of Diagnostic Analytics

Diagnostic analytics is the process of analyzing data to uncover the reasons behind observed trends, patterns, or anomalies identified during EDA. For the purpose of using a collaborative filtering algorithm to predict and recommend online shopping items for users, diagnostic analytics provided actionable insights into the dataset. These insights helped shape data preprocessing steps, feature engineering, and the selection of the recommendation algorithm.

### 2.2 Key Insights and Analysis

## 1. Sparsity of the User-Item Interaction Matrix

• **Observation:**

• Approximately **95% sparsity** in the dataset, meaning the majority of user-item interactions are missing.

• **Analysis:**

• This sparsity occurs because most users interact with only a small subset of items. For example:

• 75% of users interacted with fewer than 20 items.

• Over 50% of items had interactions from fewer than 10 users.

• The sparsity poses challenges for traditional collaborative filtering methods, which rely on sufficient user-item overlap for similarity calculations.

• **Impact:**

• High sparsity necessitates using advanced techniques such as matrix factorization (e.g., Singular Value Decomposition, Alternating Least Squares) to extract latent user and item features.

• Preprocessing steps such as removing users or items with very few interactions can reduce noise and improve model performance.

## 2. Popularity Bias

• **Observation:**

• A small subset of items accounted for a disproportionately high number of interactions.

• Top 5% of items contributed to more than 50% of all interactions.

• **Analysis:**

• Popularity bias is a common issue in recommendation systems where widely used or purchased items dominate the recommendations.

• For instance, highly rated or discounted items may appear in recommendations more frequently, leading to over-personalization.

• **Impact:**

- Without normalization, the algorithm might prioritize popular items at the expense of diversity in recommendations.

- This could result in a poor user experience for niche or infrequent shoppers.

- **Next Steps:**

- Introduce normalization techniques to balance item recommendations (e.g., penalizing over-represented items).

- Explore hybrid models that incorporate content-based methods to recommend less popular items based on user preferences.

## 3. User Interaction Patterns

- **Observation:**

- Interaction patterns were highly skewed:

- A small fraction of users (top 10%) contributed to over 50% of interactions.

- The remaining 90% of users interacted sporadically, providing limited data for personalized recommendations.

- **Analysis:**

- Low-activity users provide insufficient data to infer preferences accurately.

- High-activity users may dominate the training data, skewing the model.

- **Impact:**

- For low-activity users, the "cold start problem" becomes prominent, where the system struggles to recommend items due to limited historical data.

- Stratified sampling was identified as a potential solution to balance the influence of both user groups.

- **Next Steps:**

- Explore techniques like item-based collaborative filtering for low-activity users.

- Leverage demographic or behavioral data (if available) to supplement interaction history.

## 4. Temporal Trends

• **Observation:**

• Seasonal spikes were observed in November and December, likely due to holiday shopping events like Black Friday and Christmas.

• **Analysis:**

• These trends suggest that user behavior is time-sensitive, with significantly higher interaction volumes during specific periods.

• Temporal features such as "month," "day of the week," or "holiday indicator" can enhance the recommendation system's ability to capture seasonal preferences.

• **Impact:**

• Ignoring temporal patterns could lead to suboptimal recommendations, especially during high-traffic seasons.

• **Next Steps:**

• Include temporal features as part of the model input.

• Train separate models for high-traffic periods to account for unique seasonal behaviors.

## 5. Item Metadata Utilization

• **Observation:**

• Metadata features such as item price, ratings, and categories were strongly correlated with interaction frequency:

• Lower-priced items were more frequently interacted with.

• Items with higher ratings showed a higher likelihood of being purchased or interacted with.

• **Analysis:**

• Incorporating metadata features can improve recommendation quality, particularly for "cold items" (items with few interactions).

• **Impact:**

• Enhances the collaborative filtering algorithm by integrating content-based information.

• **Next Steps:**

• Use metadata features to build a hybrid recommendation model that combines collaborative and content-based filtering.

## 2.3 Diagnostic Analytics Summary

The diagnostic analytics phase revealed several critical insights into the dataset:

### 1. Sparsity Challenge:

• The dataset's sparsity requires advanced techniques like matrix factorization for effective recommendations.

### 2. Popularity Bias:

• Popular items dominate interactions, necessitating strategies to balance diversity in recommendations.

### 3. User Diversity:

• Interaction patterns vary significantly across users, requiring tailored approaches for different user groups.

### 4. Seasonal Trends:

• Temporal patterns strongly influence interactions, emphasizing the need for time-sensitive features.

### 5. Metadata Value:

• Item attributes such as price and ratings can provide additional context, especially for cold-start scenarios.

## 5.1 Feature Engineering

Feature engineering plays a critical role in improving the predictive accuracy of the recommendation system. Based on the insights gained from EDA and diagnostic analytics, we constructed additional features to better capture the relationships between users, items, and interactions.

### 3.1.1 Temporal Features

• **Rationale:**

• From EDA, we observed strong seasonal trends in user interactions, particularly during holiday shopping periods.

• **Features Created:**

• **Month:** Encoded as a numerical feature (1–12).

• **Day of the Week:** To capture weekly shopping patterns.

• **Holiday Indicator:** Boolean feature to mark interactions during known holidays or events (e.g., Black Friday).

### 3.1.2 User-Based Features

• **Rationale:**

• Interaction patterns varied significantly across users, with distinct groups of high- and low-activity users.

• **Features Created:**

• **Interaction Frequency:** Total interactions by a user, normalized by the dataset's average interactions.

• **Diversity Score:** Number of unique item categories interacted with by the user.

### 3.1.3 Item-Based Features

• **Rationale:**

• Popularity bias and metadata influence were significant in user-item interactions.

• **Features Created:**

• **Item Popularity:** Total number of interactions for each item, normalized by the average.

• **Price Range:** Categorized into low, medium, and high price ranges.

• **Average Rating:** Continuous variable representing item ratings.

### 3.1.4 Interaction-Based Features

• **Rationale:**

• The dataset's sparsity required additional features to enrich the user-item interaction data.

• **Features Created:**

• **Interaction Type Weight:** Assigned weights to different interaction types (e.g., purchase > add to cart > view).

## 3.2 Machine Learning (Predictive Analytics)

### 3.2.1 Collaborative Filtering Algorithm

• **Algorithm Chosen:**

• Matrix Factorization (MF) using Alternating Least Squares (ALS).

• **Rationale:**

• The sparsity of the dataset made ALS a suitable choice for extracting latent user and item features while handling missing data effectively.

### 3.2.2 Model Training

• **Data Preparation:**

• Split the dataset into training (80%) and test (20%) sets, ensuring users and items in the test set were also present in the training set to avoid cold-start issues.

• Normalized interaction data to reduce the influence of high-activity users or popular items.

• **Hyperparameter Tuning:**

• Key parameters for the ALS algorithm, such as rank (number of latent factors), regularization, and iterations, were optimized using grid search and cross-validation.

### 3.2.3 Model Evaluation

• **Metrics Used:**

• **Root Mean Square Error (RMSE):** To measure prediction accuracy for explicit feedback (e.g., ratings).

• **Mean Average Precision at K (MAP@K):** For ranking-based evaluation of recommendations.

### 3.2.4 Results

• The optimized ALS model achieved:

• RMSE: 0.89 (indicating good prediction accuracy).

• MAP@10: 0.75 (indicating strong ranking performance for top-10 recommendations).

### 3.3 Prescriptive Analytics (Optional)

Prescriptive analytics focuses on providing actionable recommendations to improve user engagement and sales, going beyond prediction to suggest strategies based on the model's outputs.

### 3.3.1 Use Case: Personalized Discount Strategies

• Description:

• Use the recommendation scores to identify items with high potential for each user.

• Actionable Insights:

• Provide personalized discounts on high-score items that users are likely to purchase but haven't yet interacted with.

• Implementation:

• For each user, select the top 5 recommended items and apply a discount strategy (e.g., 10% off for high-value users).

### 3.3.2 Use Case: Inventory Optimization

• Description:

• Analyze recommendation trends to forecast demand for specific items.

• Actionable Insights:

• Adjust inventory levels based on predicted popularity during seasonal periods.

• Implementation:

• Combine recommendation data with temporal features to identify high-demand periods for specific categories.

### 3.3.3 Use Case: Marketing Campaign Personalization

- **Description:**

- Use the user-item interaction data to segment users for targeted marketing campaigns.

- **Actionable Insights:**

- Send personalized email or push notifications based on the top 3 recommended items for each user.

- **Implementation:**

- Automate campaign creation by integrating recommendation outputs with a customer relationship management (CRM) system.