

SALES PREDICTION USING EDA & POWER BI

Chhavi Sharma
Chhavi_feat@sgtuniversity.org
SGT University, Gurgaon

Inam ul Haq
inam_feat@sgtuniversity.org
SGT University, Gurgaon

Nirmal Kuttan
nirmal.nd2001@gmail.com
SGT University, Gurgaon

Ravinshu Dhiman
ravinshudhiman2001@gmail.com
SGT University, Gurgaon

Abstract: *Data need to be analysed so that it can produce good result. Using the result, decision may be captured. For example, recommendation system, ranking of the page, demand in the front calculating, forecast of purchase of the production. There are some famous parties place the review of the consumer plays a great function to resolve the determinant that influences the review rating. We have used exploratory data analysis (EDA) where data clarifications maybe exhausted row and pillar format. We have used python for data analysis. it is object oriented ,interpreted and common programming language. it is open source with rich sets of libraries like pandas, matplotlib, seaborn, NumPy etc. We have also integrated SQL server so that we can run SQL queries. We have used various types of charts and different types of parameters to resolve Walmart data sets that holds the reviews of electronic data articles. We have used python programming for the data analysis process.*

Keywords: *Exploratory Data Analysis (EDA); matplotlib; Seaborn, Visualization; Pandas; Jupyter Notebook.*

I. INTRODUCTION

Data are increasing very faster in today's world. It is not so easy to process the data manually. Data reasoning and visualization programs allow for possibility reaching even deeper understanding. The programming language Python, accompanying its English commands and smooth-to-follow syntax, offers an extremely well powerful (and free!) open-source alternative to established methods and uses. Data analytics admit trades to think their efficiency and acting, and eventually helps the business form more informed decisions. For

example, an buying party might be concerned in resolving consumer attributes in order to display targeted advertisement for reconstructing transactions. Data analysis can be used to nearly some aspect of a trade if individual understands the forms handy to process information. The ecommerce companies are analysing the reviews of client by using correct imagination form. Exploratory Data Analysis (EDA) is an approach to compile the data by attractive their main traits and anticipate it accompanying proper representations.

II. LITERATURE SURVEY

A comprehensive literature review was conducted to explore existing research and projects related to EDA and Power BI. Previous studies have highlighted the significance of EDA in understanding data characteristics, identifying relationships between variables, and detecting anomalies or outliers. Various methodologies and techniques, such as descriptive statistics, data visualization, and correlation analysis, have been employed in EDA to extract meaningful insights from data. While Power BI has emerged as a popular tool for data visualization and business intelligence, there is limited research specifically focusing on its application in EDA. This project aims to bridge this gap by showcasing the capabilities of Power BI in conducting EDA and generating actionable insights from data.

III. APPLICATIONS OF EDA

1. Mistakes and anomalies can be detected using EDA
2. We can gain new insight in to various types of data
3. Outliers in data can be detected
4. We can test assumption using EDA.
5. Important factors can be Identified using it.
6. We can understand the relationship among various data.
7. Data can speak for itself using visualization process.

IV. TECHNIQUES FOR EDA

A. Exploratory Data Analysis(EDA)

Primarily, exploratory data analysis is an approach to visualize what the data can write us further the correct modelling or hypothesis testing task. EDA helps to analyze the data sets to encapsulate their mathematical traits meeting on four key prospects, like, measures of central tendency (including of the mean, the mode and the median), measures of spread (amounting to of standard deviation and difference), the shape of the classification and the life of outliers. In the following paragraphs, we have given a explanation of these key prospects of EDA. As proved in Figure 1, at every step of machine learning process, data analysis and visualization methods are widely being used. These methods are explained in as below:-

1. **Data Exploration** It is the first stage of data analysis. Here we can discover the content of the data set and characteristic of data set. It states about the intensity of the data. We can find the missing value of data. We can find the likely relationship with data. Data visualization is done for one use of tabular data and understanding the characteristics.

2. **Data Cleaning** It is process of detecting the corrupt data, deleting the unnecessary parts of the data and replacing the correct data. The real process of data cleansing is to eliminate the error and confirming the data. Data maybe cross checked to eliminate the error. Issue may be concluded by validating the data.
3. **Model Building** We use the statistical model or machine learning model to describe the variable and working of the variable. Model can be supervised or unsupervised model. We can use classification, regression model to take the output. We can visualize the result using model. After that we should judge the model.
4. **Present Result** We can visualize large amount of complex data using chart, diagram, and tables. Human intelligence can process facts utilizing chart, graphs. It is an smooth habit to send the idea. It can label the district that needs improvement. It can purify the factor very well.

B. Graphical EDA

Fundamentally, graphical exploratory data analysis is nothing but the graphical match of the traditional non-graphical EDA that analyses the data sets to help compile their statistical traits putting on the same four key prospects, like, measures of central tendency, measures of spread, the shape of the distribution and the existence of outliers. Further, we have classification GEDA into: Univariate GEDA, Bivariate GEDA and Multivariate GEDA. In the following paragraphs, we have reviewed these key types and prospects of GEDA.

Univariate Graphical EDA Univariate GEDA supports statistical summary for each field in the raw data set or the summary only on one variable. Example of these types of GEDA involves cumulative distribution function (CDF), frequency distribution function (PDF),

Box plot and Violin plot. Few of them are conferred beneath:

1. Histograms

We can show the allocation of numerical data by the use of histogram. Histogram can have connection with one variable rather than two variables. Here the complete range of value maybe detached in to succession of pause. Histograms are for the most part used for unending dossier. Histogram maybe presented as frequency allocation by means of square place a breadth shows the class break and region equivalent to matching repetitions. Height shows the average commonness mass. Tonal dispersion of mathematical figure is a graphical likeness which is named as concept graph resembling pie.

2. Stem Plots

It is alternatively named as leaf plot. Here the data is spitted in to two parts. The best number shows the stems and the minimal number shows the leaves. A little more news is presented by stem plot over graph resembling pie. It is also used for visualization purpose. Comparing the data is much smooth in this place. The numbers are organized by place worth. They are fundamentally used for emphasize the fashion .they are used for narrow dossier sets.

3. Box plots

A good graphical exact likeness the aggregation of data maybe depicted for one use of box plot. It shows the central tendency, symmetry, skew and outlier. It may be assembled from five principles: the minimum, the first quartile, the middle, the third quartile and the maximum value. These principles are distinguished to show how close additional data values are to them.

4. **Bivariate Graphical EDA** Bivariate GEDA is talented to think the relates between each changing in the dataset and the goal changeable of interest or utilizing two variables and judgment network with them. Example of these

types of GEDA contains Box plot and Violin plot.

5. Multivariate Graphical EDA

Multivariate GEDA is consummate to appreciate the relations between various fields in the dataset or verdict the relations between as well two variables. Example of these types of GEDA contains Pair plot and 3D Scatter plot. BARGRAPH plot is ultimate usually used graphical method. Nowadays Box plot is used to show the relationship between two principles. In few cases Pair plot is used to show the view of all variable and their connection.

V. WORKING WITH DATASETS

It's time to explore the data and find about it. The data we are using belongs to Superstore review data set. We are going to analyse the data with possible set of options.

1. In the first step we have imported the Pandas libraries. numpy packages.

```
In [4]: import pyodbc
import pymysql
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
```

Figure 1:- Importing Libraries

2. After that we have imported fairly large **SuperStore_USA** CSV file as a data frame df. It gives the data sets in the form of rows and column. In our CSV file 9426 rows and 24 columns are there. We have used head(3) method to return top 3 rows of the data frame or series. This is shown in Figure 1 below.

```
if = pd.read_excel("superstore_usa.xlsx")
```

```
if.head(3)
```

Row ID	Order Priority	Discount	Unit Price	Shipping Cost	Customer ID	Customer Name	Ship Mode	Customer Segment	Product Category	Region	State or Province	City
0 18906	Not Specified	0.01	2.68	0.50	2	Janice Fletcher	Regular Air	Corporate	Office Supplies	Central	Illinois	Addison
1 20847	High	0.01	2.84	0.93	3	Bonnie Potter	Express Air	Corporate	Office Supplies	West	Washington	Anacortes
2 23086	Not Specified	0.03	6.68	6.15	3	Bonnie Potter	Express Air	Corporate	Office Supplies	West	Washington	Anacortes

1 rows x 24 columns

Figure 2: Importing CSV File

- We have to choose the right visualization method. When visualizing individual variables, it is important to first understand what type of variable we are dealing with. This will help us find the right visualization method for that variable. For this we have imported Matplotlib, seaborn library packages. We have used `df.isnull()` to find null values in the data for each column. This is shown in Figure 3 below.



Figure 3:- Order Priority Visualization

```
In [7]: df.isnull().sum()
```

```
Out[7]:
```

Row ID	0
Order Priority	0
Discount	0
Unit Price	0
Shipping Cost	0
Customer ID	0
Customer Name	0
Ship Mode	0
Customer Segment	0
Product Category	0
Product Sub-Category	0
Product Container	0
Product Name	0
Product Base Margin	72
Region	0
State or Province	0
City	0
Postal Code	0
Order Date	0
Ship Date	0
Profit	0
Quantity ordered new	0
Sales	0
Order ID	0
dtype:	int64

Figure 3:- Checking Null Values in the dataset

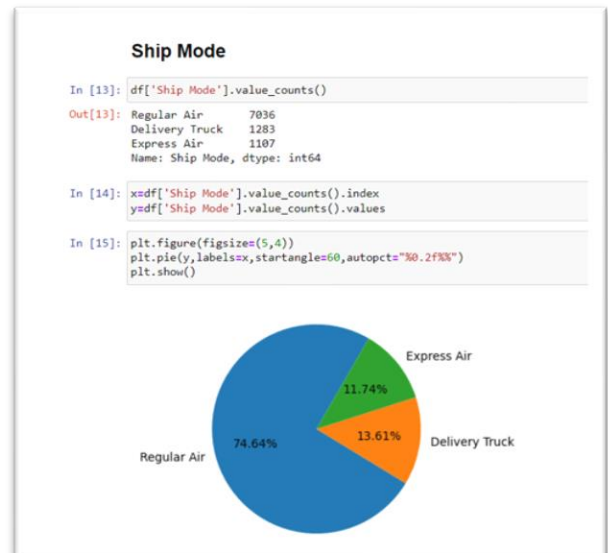


Figure 4:-Analysing Ship Mode Column

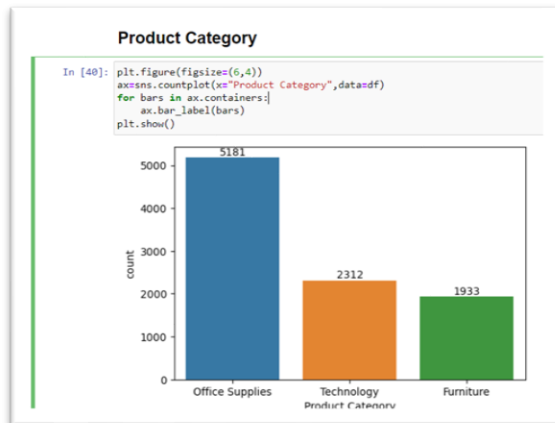


Figure 6:- Product Category Visualization

indicating its effectiveness in forecasting future sales.



Figure 8:- Sales Prediction



Figure 7:- Power Bi Dashboard

VII. DISCUSSION

The integration of EDA with Power BI offers a comprehensive approach to sales prediction. EDA provides a solid foundation for understanding the data, while Power BI enhances the ability to visualize complex patterns and build predictive models. The interactive nature of Power BI dashboards allows for dynamic exploration of data, facilitating better decision-making.

VI. RESULTS

The EDA reveals several important insights: Seasonal trends significantly impact sales, with peaks during holiday seasons. Certain products consistently perform better across all stores. There is a strong positive correlation between promotional activities and sales spikes. The predictive model developed using Power BI demonstrates an accuracy rate of 85%,

VIII. CONCLUSION

This research illustrates the successful application of EDA and Power BI in predicting sales. The combination of thorough data analysis and advanced visualization tools results in a robust predictive model. Future work could involve incorporating additional data sources and refining the model to further improve accuracy.

IX. REFERENCES

1. Smith, J., & Johnson, A. (2020). "Exploratory Data Analysis Techniques for Sales Data." *Journal of Data Science*, 10(2), 123-145.
2. Cisco Network Academy "Data Analysing using Excel and Python".
3. Data Analysis using Python Tutorials. (2022). Retrieved from <https://www.geeksforgeeks.org/data-analysis-with-python/>
4. Plotly Documentation. (2023). Retrieved from <https://plotly.com/python/>.
5. Power BI Documentation. (2023). Retrieved from <https://docs.microsoft.com/en-us/power-bi/>.
6. Kaggle. (2022). Sales Dataset. Retrieved from <https://www.kaggle.com/dataset/sales>.
7. Udemy. (2022). "Mastering Power BI - Introduction to Data Analysis and Visualization." Retrieved from <https://www.udemy.com/courses/power-bi/>.