

PURCHASE PATTERN ANALYTICS

27 NOVEMBER 2024



TEAM MEMBERS

TEAM ID PTID-CDA-OCT-24-
2I7

PROJECT ID CDACL005

NIMMALA BHARATH
GOUD

ARYAN KUNDHARA

ASAOLU
OYINLOLAJOY

ASAOLU
OYINLOLAJOY





PROBLEM STATEMENT

- IN THE CONTEXT OF A RETAIL ENVIRONMENT, THE CHALLENGE IS TO LEVERAGE MARKET BASKET ANALYSIS (MBA) ON A TRANSACTION DATASET TO UNVEIL RELATIONSHIPS BETWEEN PRODUCTS. BY CONDUCTING EXPLORATORY DATA ANALYSIS, ADDRESSING DATA QUALITY ISSUES, EMPLOYING DATA VISUALIZATION TECHNIQUES, AND IMPLEMENTING THE APRIORI ALGORITHM, THE OBJECTIVE IS TO EXTRACT VALUABLE INSIGHTS. THIS INCLUDES UNDERSTANDING CUSTOMER PURCHASING BEHAVIOR, IDENTIFYING FREQUENTLY CO-OCCURRING PRODUCTS, AND DERIVING ACTIONABLE RECOMMENDATIONS TO ENHANCE SALES STRATEGIES AND CUSTOMER SATISFACTION.





TASKS ASSIGNED

- - ACQUIRE AND LOAD THE TRANSACTION DATASET.
- - BEGIN EXPLORATORY DATA ANALYSIS (EDA):
- - IDENTIFY KEY STATISTICS (MEAN, MEDIAN, MODE, ETC.).
- - VISUALIZE DATA DISTRIBUTIONS AND PATTERNS.
- - DETECT ANY OUTLIERS OR MISSING VALUES.



METHODOLOGY

Data Collection:

- STEPS:
- ESTABLISH CONNECTION TO THE DATABASE:
- HOST: 18.136.157.135
- USER: DM_TEAM16
- PASSWORD: 2O_HIHIFETRE
- IMPORT THE DATASET INTO A SUITABLE ANALYTICAL TOOL (PYTHON, R, OR SQL-BASED ENVIRONMENT).
- VERIFY THE STRUCTURE AND TYPES OF DATA FIELDS.

DATA CLEANING AND TRANSFORMATION

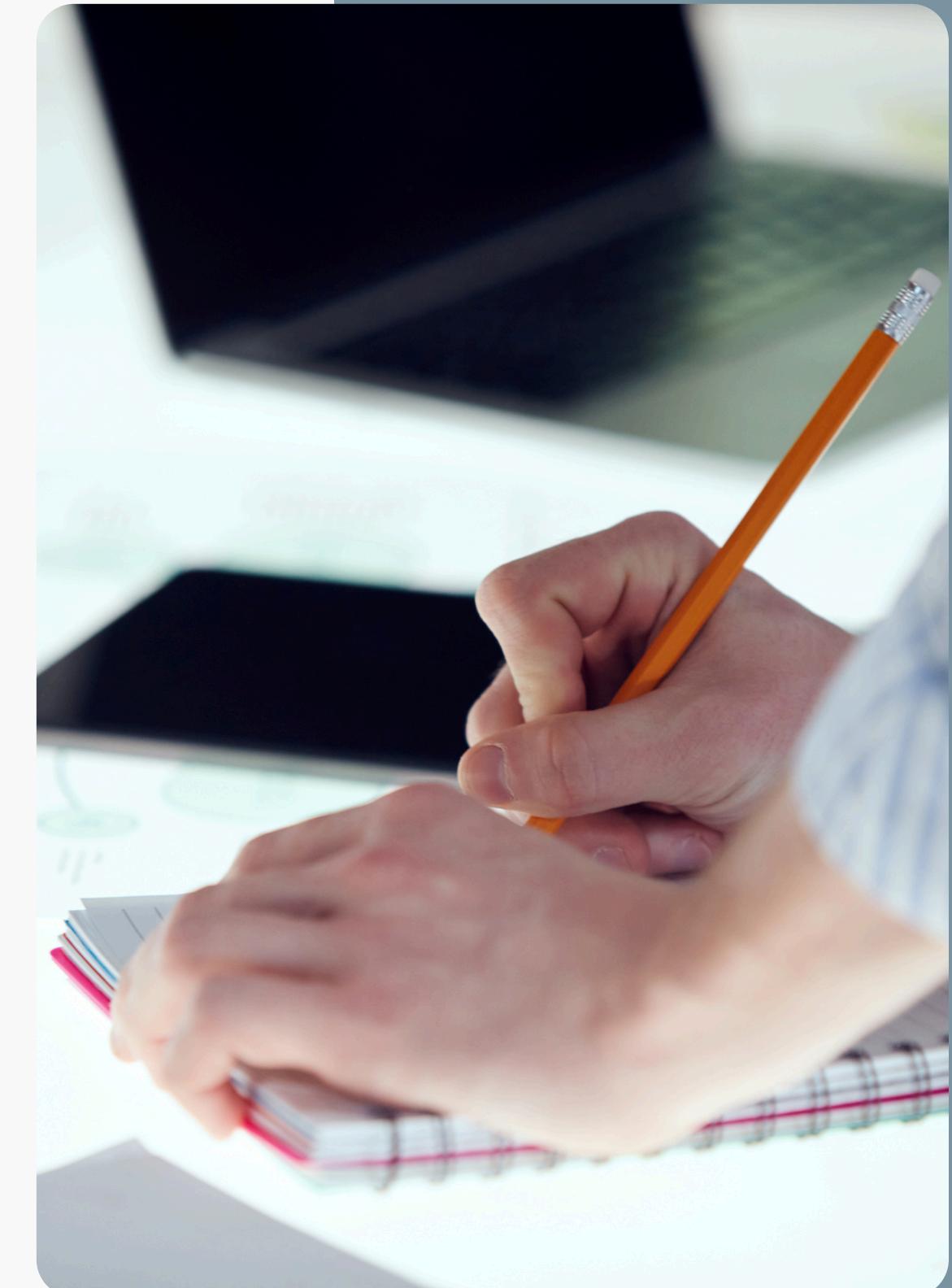
- HANDLE MISSING VALUES:
- IMPUTE WITH STATISTICAL METHODS (MEAN, MEDIAN) OR REMOVE ROWS/COLUMNS BASED ON RELEVANCE.
- CORRECT OUTLIERS OR ERRONEOUS ENTRIES.
- CONVERT DATA INTO TRANSACTIONAL FORMAT FOR MARKET BASKET ANALYSIS:
- EXAMPLE: FROM TABULAR FORMAT TO ITEMSET LISTS.

ADVANCED ANALYSIS WITH APRIORI ALGORITHM

- CREATE VISUALIZATIONS SHOWING TRENDS AND PATTERNS IN THE DATA.
- IMPLEMENT THE APRIORI ALGORITHM TO IDENTIFY:
- FREQUENT ITEMSETS.
- ASSOCIATION RULES (E.G., CONFIDENCE, LIFT)

TOOLS AND TECHNIQUES

- LIBRARIES: PANDAS, NUMPY, MATPLOTLIB, SEABORN, AND MLXTEND (FOR APRIORI).
- ALGORITHM: APRIORI FOR MBA.
- VISUALIZATION PLATFORMS: TABLEAU OR POWER BI FOR ADVANCED GRAPHS.

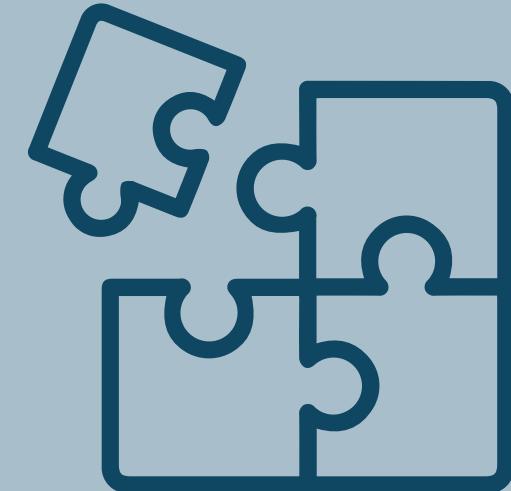


PROJECT OBJECTIVES



ANALYSIS PHASE

- WE DELVE INTO EXPLORING OUR DATASET TO GAIN INSIGHTS AND IDENTIFY PATTERNS. WE BEGIN BY EXAMINING THE STRUCTURE OF OUR DATA USING THE `.INFO()` METHOD, WHICH PROVIDES A SUMMARY OF THE DATA TYPES AND NON-NULL COUNTS FOR EACH COLUMN.



STRATEGY DEVELOPMENT

- IT'S ESSENTIAL TO ENSURE THAT OUR DATASET IS IN THE RIGHT FORMAT AND FREE OF INCONSISTENCIES. IN THIS CHAPTER, WE'LL WALK THROUGH THE STEPS TAKEN TO PREPARE OUR DATA FOR ANALYSIS. THIS INCLUDES CONVERTING DATE COLUMNS, CLEANING UP NUMERIC FORMATS, CHECKING FOR MISSING VALUES, AND GENERATING DESCRIPTIVE STATISTICS



IMPLEMENTATION PLAN

- IMPORTING LIBRARIES
- READING OUR DATA
- DATA PREPARATION
- DATA FILTERING AND CLEANING
- TOTAL SALES ANALYSIS
- UNIQUE ITEM ANALYSIS
- MARKET BASKET ANALYSIS
- ASSOCIATION RULE MINING
- CONFIDENCE
- LIFT ANALYSIS
- ANALYSIS OF ASSOCIATION RULE METRICS



1| Importing libraries

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import plotly.express as px
import plotly.graph_objects as go
from plotly.subplots import make_subplots # Importing make_subplots
from mlxtend.frequent_patterns import apriori,association_rules
import warnings
```



READ OUR DATA

2| Reading our data

```
df = pd.read_csv("purchase patteren.csv")
df.head()
```

	BillNo	Itemname	Quantity	Present_Date	Price	CustomerID	Country
0	536365	WHITE HANGING HEART T-LIGHT HOLDER	6	01-12-2010 08:26	2.55	17850.0	United Kingdom
1	536365	WHITE METAL LANTERN	6	01-12-2010 08:26	3.39	17850.0	United Kingdom
2	536365	CREAM CUPID HEARTS COAT HANGER	8	01-12-2010 08:26	2.75	17850.0	United Kingdom
3	536365	KNITTED UNION FLAG HOT WATER BOTTLE	6	01-12-2010 08:26	3.39	17850.0	United Kingdom
4	536365	RED WOOLLY HOTTIE WHITE HEART.	6	01-12-2010 08:26	3.39	17850.0	United Kingdom

```
#Get Summary
print(df.shape)
print(df.info())
```

```
(516778, 7)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 516778 entries, 0 to 516777
Data columns (total 7 columns):
 #   Column           Non-Null Count   Dtype  
 ---  -- 
 0   BillNo          516778 non-null    object 
 1   Itemname        515323 non-null    object 
 2   Quantity        516778 non-null    int64  
 3   Present_Date    516778 non-null    object 
 4   Price           516778 non-null    float64
 5   CustomerID     382811 non-null    float64
 6   Country         516778 non-null    object 
dtypes: float64(2), int64(1), object(4)
memory usage: 27.6+ MB
None
```



DATA PREPARATION

3| Data preparation

Before diving into the analysis, it's essential to ensure that our dataset is in the right format and free of inconsistencies. In this chapter, we'll walk through the steps taken to prepare our data for analysis. This includes converting date columns, cleaning up numeric formats, checking for missing values, and generating descriptive statistics.

```
: df["Present_Date"] = pd.to_datetime(df["Present_Date"], format='%d.%m.%Y %H:%M')

df["Year/Month"] = df["Present_Date"].dt.to_period("M")

: missing = df.isnull().sum()

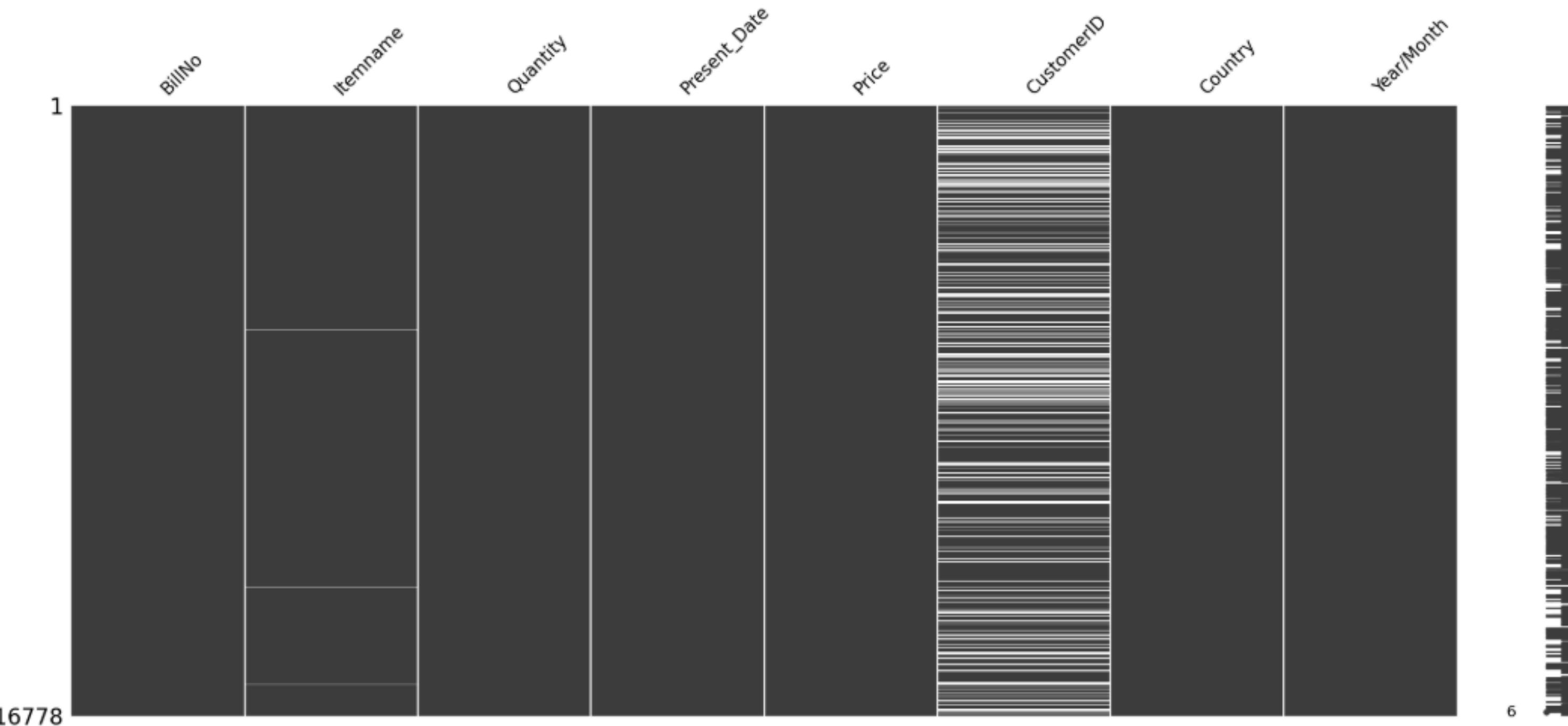
print(missing)
```

```
BillNo          0
Itemname        0
Quantity         0
Present_Date    0
Price            0
CustomerID      0
Country          0
Year/Month       0
TotalPrice       0
dtype: int64
```

DATA PREPARATION

```
import missingno as msno  
  
# Visualize missing values as a matrix  
msno.matrix(df)
```

<Axes: >



6



KEY STATISTICS

```
: #Statistics  
df.describe()
```

	Quantity	Present_Date	Price	CustomerID
count	516778.00000	516778	516778.000000	382811.000000
mean	10.16247	2011-07-04 08:18:39.177790464	3.841504	15310.347702
min	-9600.00000	2010-12-01 08:26:00	-11062.060000	12346.000000
25%	1.00000	2011-03-28 09:59:00	1.250000	13928.500000
50%	3.00000	2011-07-19 14:29:00	2.080000	15249.000000
75%	10.00000	2011-10-19 12:11:00	4.130000	16820.000000
max	80995.00000	2011-12-09 12:50:00	13541.330000	18287.000000
std	161.91653	NaN	42.113493	1722.483516



TASKS ASSIGNED

- DATA CLEANING AND TRANSFORMATION
- DATASET INTO A SUITABLE FORMAT FOR MARKET BASKET ANALYSIS.
- DATA VISUALIZATION AND APRIORI ALGORITHM IMPLEMENTATION
- CREATE VISUALIZATIONS TO REPRESENT PRODUCT OCCURRENCES,
ASSOCIATIONS, AND TRENDS
- IMPLEMENT THE APRIORI ALGORITHM FOR MARKET BASKET ANALYSIS
- IDENTIFY FREQUENT ITEM SETS AND ASSOCIATION RULES
- VISUALIZE THE RESULTS OF THE APRIORI ALGORITHM





4| Data Filtering and Cleaning

A crucial aspect of data analysis is ensuring the data's integrity by filtering out irrelevant or erroneous entries and handling missing values appropriately. In this chapter, we focus on filtering out non-positive values, removing rows with missing item names, filling in missing customer IDs, and calculating total prices per transaction. These steps are vital to ensure the accuracy and reliability of our subsequent analysis.

```
df = df[(df["Quantity"] > 0) & (df["Price"] > 0)]
```

```
#drop rows without item
```

```
df = df[df["Itemname"].notnull()]
```

```
# Filling missing customer IDs
```

```
df = df.fillna('#NV')
```

```
df["TotalPrice"] = df["Quantity"] * df["Price"]
```

```
df.head()
```

	BillNo	Itemname	Quantity	Present_Date	Price	CustomerID	Country	Year/Month	TotalPrice
0	536365	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00	2.55	17850.0	United Kingdom	2010-12	15.30
1	536365	WHITE METAL LANTERN	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
2	536365	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00	2.75	17850.0	United Kingdom	2010-12	22.00
3	536365	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34
4	536365	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00	3.39	17850.0	United Kingdom	2010-12	20.34



EXPLORATORY DATA ANALYSIS

5| EDA

In this section, we delve into exploring our dataset to gain insights and identify patterns. We begin by examining the structure of our data using the `.info()` method, which provides a summary of the data types and non-null counts for each column.

```
[3]: df.info()
```

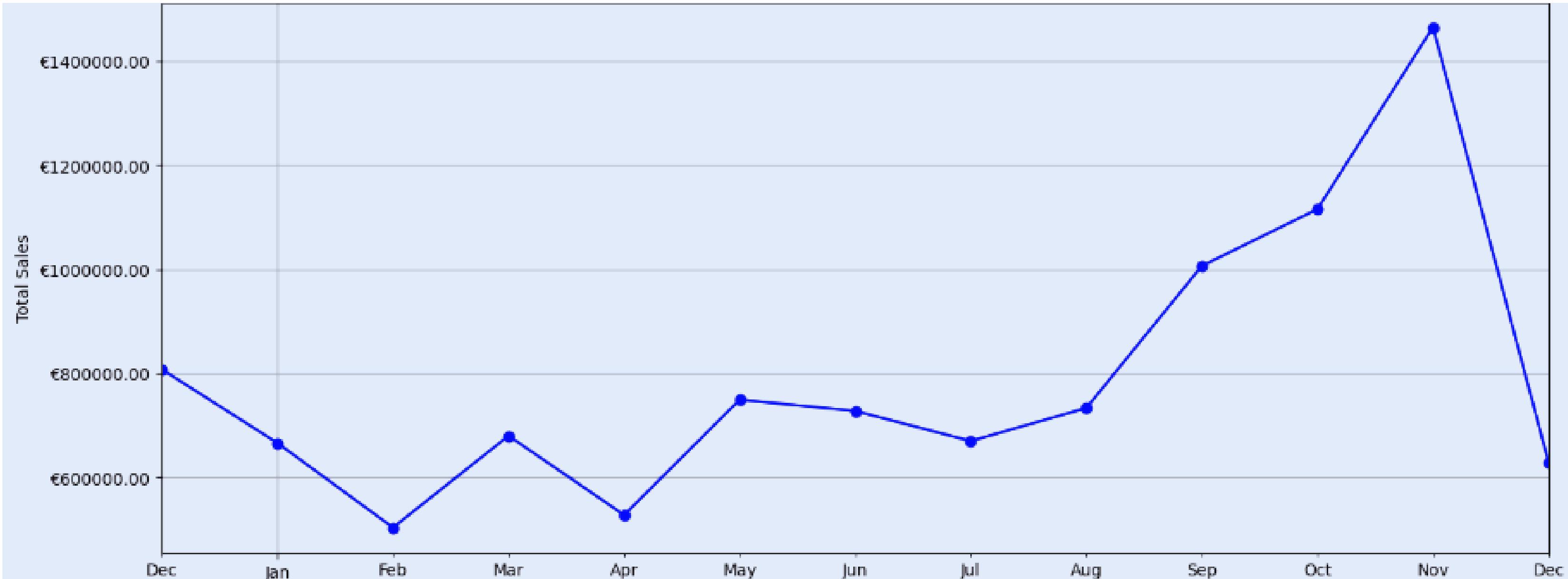
```
<class 'pandas.core.frame.DataFrame'>
Index: 514270 entries, 0 to 516777
Data columns (total 9 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   BillNo       514270 non-null   object 
 1   Itemname     514270 non-null   object 
 2   Quantity     514270 non-null   int64  
 3   Present_Date 514270 non-null   datetime64[ns]
 4   Price        514270 non-null   float64
 5   CustomerID   514270 non-null   object 
 6   Country      514270 non-null   object 
 7   Year/Month   514270 non-null   period[M]
 8   TotalPrice   514270 non-null   float64
dtypes: datetime64[ns](1), float64(2), int64(1), object(4), period[M](1)
memory usage: 39.2+ MB
```



TOTAL SALES ANALYSIS

WE GROUP THE DATA BY MONTH AND YEAR, CALCULATING THE TOTAL SUM OF SALES TO UNDERSTAND THE SALES TREND OVER TIME. THE RESULTING VISUALIZATION DEPICTS THE TOTAL SALES PER MONTH.

```
IMPORT MATPLOTLIB.TICKER AS TICKER  
# GROUPING THE DATA BY MONTH AND YEAR, AND CALCULATING THE TOTAL SUM OF SALES  
MONTHLY_SALES = DF.GROUPBY('YEAR/MONTH')[TOTALPRICE].SUM()  
  
PLT.FIGURE(FIGSIZE=(15,6))  
MONTHLY_SALES.PLOT(KIND='LINE', MARKER='O', COLOR='B')  
PLT.TITLE('TOTAL SALES PER MONTH')  
PLT.XLABEL('MONTH')  
PLT.YLABEL('TOTAL SALES')  
  
FORMATTER = TICKER.FORMATSTRFORMATTER('€%.2F')  
PLT.GCA().YAXIS.SET_MAJOR_FORMATTER(FORMATTER)  
  
PLT.GRID(TRUE)  
PLT.SHOW()
```



SALES PERFORMANCE ANALYSIS

KEY OBSERVATIONS

- DOWNWARD START: SALES START DECLINING AFTER DECEMBER BUT HIT A LOW POINT IN FEBRUARY 2011.
- FLUCTUATIONS: FROM MARCH TO AUGUST, SALES SHOW MINOR FLUCTUATIONS BUT REMAIN RELATIVELY STABLE.
- STEADY GROWTH: FROM SEPTEMBER ONWARDS, SALES DEMONSTRATE CONSISTENT GROWTH, PEAKING SIGNIFICANTLY IN NOVEMBER 2011.
- SHARP DROP: SALES DECLINE STEEPLY FROM NOVEMBER TO DECEMBER

PERFORMANCE ANALYSIS

- SEASONALITY: THE SIGNIFICANT PEAK IN NOVEMBER SUGGESTS SEASONAL TRENDS, POSSIBLY INFLUENCED BY HOLIDAY SHOPPING OR END-OF-YEAR ACTIVITIES.
- LOW MONTHS: FEBRUARY EXPERIENCED THE LOWEST PERFORMANCE, POSSIBLY DUE TO POST-HOLIDAY ECONOMIC SLOWDOWNS.
- GROWTH OPPORTUNITIES: THE STEADY GROWTH FROM SEPTEMBER TO NOVEMBER HIGHLIGHTS A POTENTIAL MARKETING OR SALES CAMPAIGN EFFECTIVENESS.



RECOMMENDATIONS

1. FOCUS ON NOVEMBER: IDENTIFY STRATEGIES THAT DROVE THE NOVEMBER PEAK AND REPLICATE THEM IN OTHER MONTHS.
2. BOOST EARLY-YEAR SALES: INTRODUCE PROMOTIONS OR CAMPAIGNS IN Q1 (JANUARY–MARCH) TO COUNTER POST-HOLIDAY SLUMPS.
3. MONITOR TRENDS: ASSESS MARKET FACTORS INFLUENCING SALES STABILITY FROM MARCH TO AUGUST FOR IMPROVEMENT.
4. PREPARE FOR DECEMBER: ADDRESS THE SHARP DECLINE BY ALIGNING MARKETING STRATEGIES WITH YEAR-END CUSTOMER BEHAVIORS.

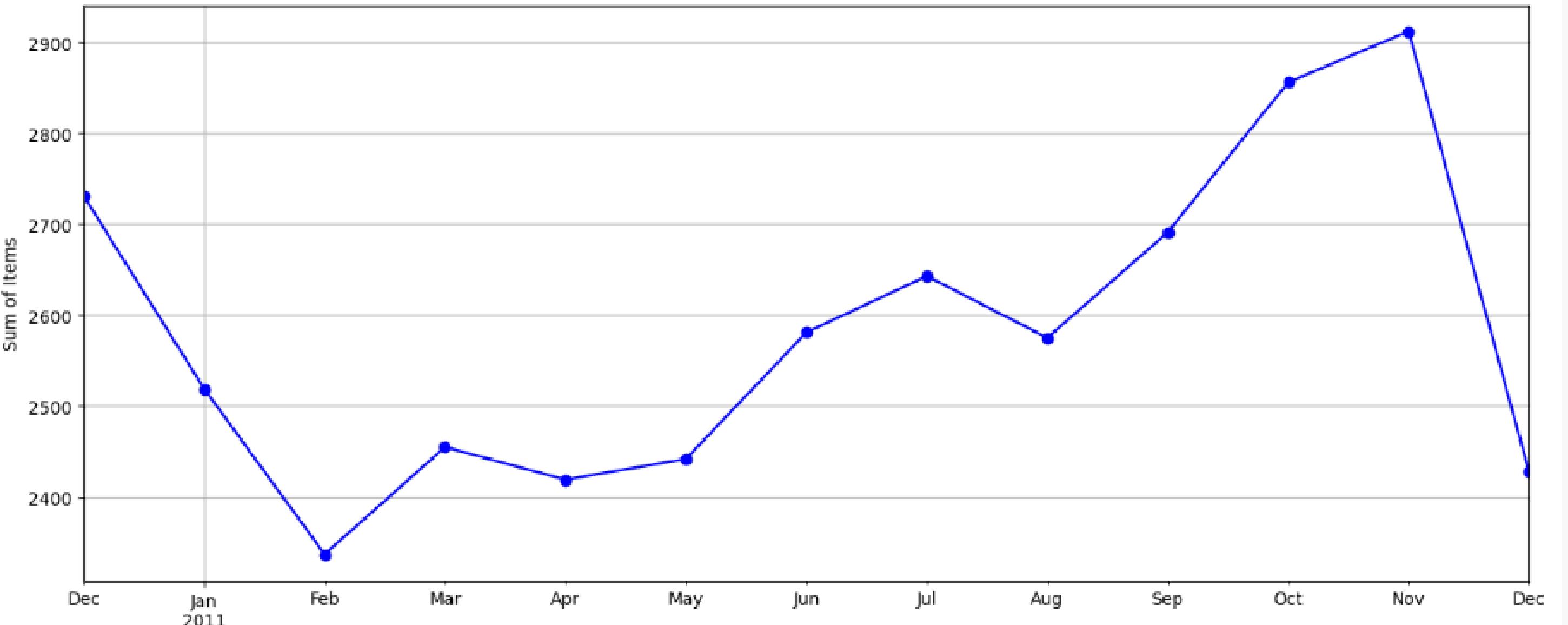
Unique Item Analysis

Next, we analyze the uniqueness of items sold per month by grouping the data and calculating the count of unique items. The line plot visualizes the sum of unique items per month.

```
monthly_item = df.groupby("Year/Month")["Itemname"].nunique()

plt.figure(figsize=(15,6))
monthly_item.plot(kind="line",marker="o",color="b")
plt.title("Sum of Unique Items per Month")
plt.xlabel("Month")
plt.ylabel("Sum of Items")
plt.grid(True)
plt.show()
```

Sum of Unique Items per Month



UNIQUE ITEM ANALYSIS

KEY OBSERVATIONS:

1. DECLINING START (DECEMBER–FEBRUARY):
 - THE SUM OF UNIQUE ITEMS SOLD DROPS SIGNIFICANTLY FROM DECEMBER TO FEBRUARY.
 - THIS INDICATES REDUCED VARIETY IN CUSTOMER PURCHASES, LIKELY DUE TO THE POST-HOLIDAY LULL.
2. STABILIZATION (MARCH–AUGUST):
 - FROM MARCH TO AUGUST, THE COUNT OF UNIQUE ITEMS SOLD STABILIZES, SHOWING MINOR FLUCTUATIONS.
 - THIS SUGGESTS A STEADY BUT LIMITED DEMAND FOR UNIQUE ITEMS DURING THESE MONTHS.
3. STRONG GROWTH (SEPTEMBER–NOVEMBER):
 - A SHARP INCREASE IN UNIQUE ITEM SALES IS OBSERVED FROM SEPTEMBER TO NOVEMBER.
 - NOVEMBER ACHIEVES THE PEAK FOR THE YEAR, INDICATING HEIGHTENED DEMAND FOR DIVERSE ITEMS, LIKELY DRIVEN BY SEASONAL FACTORS (E.G., HOLIDAY PREPARATIONS OR PROMOTIONS).
4. SHARP DECLINE IN DECEMBER:
 - SIMILAR TO OVERALL SALES TRENDS, THE COUNT OF UNIQUE ITEMS DROPS ABRUPTLY IN DECEMBER, POTENTIALLY DUE TO INVENTORY SELL-OUTS OR CUSTOMER FOCUS SHIFTING TO FEWER HIGH-DEMAND ITEMS.

PERFORMANCE INSIGHTS:

- SEASONAL INFLUENCES: THE SHARP GROWTH IN UNIQUE ITEMS FROM SEPTEMBER TO NOVEMBER ALIGNS WITH THE HOLIDAY SHOPPING PERIOD.
- LIMITED DEMAND IN EARLY YEAR: THE DECLINE IN FEBRUARY SUGGESTS THAT Q1 EXPERIENCES LIMITED MARKET ACTIVITY FOR UNIQUE ITEMS.
- STABILIZED MIDDLE MONTHS: MARCH TO AUGUST PROVIDES A STABLE BASELINE, SHOWING A CONSISTENT BUT LESS DYNAMIC VARIETY IN PURCHASES.

RECOMMENDATIONS:

1. CAPITALIZE ON Q4 TRENDS:
 - BOOST THE AVAILABILITY OF UNIQUE ITEMS DURING SEPTEMBER TO NOVEMBER TO MAXIMIZE REVENUE FROM INCREASED DEMAND.
 - STRATEGIZE TARGETED CAMPAIGNS FOR PROMOTING DIVERSE PRODUCTS DURING THIS PERIOD.
2. ADDRESS EARLY-YEAR SLUMP:
 - EXPLORE PROMOTIONAL STRATEGIES OR NEW PRODUCT LAUNCHES IN Q1 TO STIMULATE INTEREST IN UNIQUE ITEMS.
3. INVENTORY OPTIMIZATION:
 - ENSURE ADEQUATE INVENTORY OF DIVERSE PRODUCTS DURING PEAK MONTHS AND AVOID OVERSTOCKING IN SLOWER MONTHS.
4. ANALYZE CUSTOMER PREFERENCES:
 - CONDUCT DEEPER ANALYSIS TO IDENTIFY HIGH-PERFORMING UNIQUE ITEMS IN SEPTEMBER–NOVEMBER AND ADAPT OFFERINGS ACCORDINGLY.



Market Basket Analysis

In this section, we embark on a crucial aspect of retail analytics: Market Basket Analysis. This technique allows us to uncover hidden patterns and associations among items purchased together, providing valuable insights for optimizing product placement, marketing strategies, and inventory management.

6|MBA - Preprocessing

Before diving into analysis, we need to preprocess our transactional data. We select pertinent columns, such as 'BillNo' and 'Itemname', essential for identifying item associations. To facilitate analysis, we encode the items into a binary format using one-hot encoding, creating a basket representation where each row represents a transaction (BillNo) and each column represents an item. This step lays the groundwork for further analysis.

```
# Preprocessing the data
data_processed = df[['BillNo', 'Itemname']]
data_encoded = pd.get_dummies(data_processed, columns=['Itemname'])
data_encoded.columns = data_encoded.columns.str.replace("Itemname_", "")
basket = data_encoded.groupby('BillNo').sum()

basket.head()
```



*Boombbox Ipod Classic	*USB Office Mirror Ball	10 COLOUR SPACEBOY	12 COLOURED PARTY BALLOONS	12 DAISY PEGS IN WOOD BOX	12 EGG HOUSE PAINTED WOOD	12 HANGING EGGS HAND PAINTED	12 IVORY ROSE PEG PLACE SETTINGS	12 MESSAGE CARDS WITH ENVELOPES	12 PENCIL SMALL TUBE WOODLAND	...	ZINC STAR T- LIGHT HOLDER	ZINC SWEETHEART SOAP DISH	ZINC SWEETHEART WIRE LETTER RACK	Z H
------------------------------	----------------------------------	--------------------------	-------------------------------------	--	------------------------------------	--	--	---	--	-----	------------------------------------	---------------------------------	---	--------

BillNo

536365	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536366	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536367	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536368	0	0	0	0	0	0	0	0	0	0	0	0	0	0
536369	0	0	0	0	0	0	0	0	0	0	0	0	0	0

5 rows × 4006 columns



To simplify the analysis and focus solely on item presence, we convert the occurrence values to binary format. This binary transformation facilitates the identification of item associations, regardless of the quantity purchased in each transaction.

```
basket[basket > 0] = 1
```



APPLYING APRIORI ALGORITHM

GENERATING FREQUENT ITEMSETS USING THE APRIORI ALGORITHM, WE AIM TO EXTRACT FREQUENT ITEMSETS – COMBINATIONS OF ITEMS THAT FREQUENTLY CO-OCCUR IN TRANSACTIONS. THESE ITEMSETS PROVIDE VALUABLE INSIGHTS INTO CUSTOMER PURCHASING BEHAVIORS, REVEALING COMMONLY ASSOCIATED PRODUCTS.

```
frequent_itemsets = apriori(basket,min_support=0.02,use_colnames=True)  
frequent_itemsets.head()
```

	support	itemsets
0	0.023364	(3 STRIPEY MICE FELTCRAFT)
1	0.023671	(4 TRADITIONAL SPINNING TOPS)
2	0.048160	(6 RIBBONS RUSTIC CHARM)
3	0.021319	(60 CAKE CASES DOLLY GIRL DESIGN)
4	0.030675	(60 CAKE CASES VINTAGE CHRISTMAS)



ASSOCIATION RULE MINING

IN THIS PHASE, WE EMPLOY ASSOCIATION RULE MINING TECHNIQUES TO EXTRACT MEANINGFUL PATTERNS AND RELATIONSHIPS BETWEEN ITEMS IN OUR DATASET. BY ANALYZING THESE RULES, WE CAN UNCOVER ACTIONABLE INSIGHTS THAT DRIVE BUSINESS DECISIONS, SUCH AS PRODUCT PLACEMENT STRATEGIES AND TARGETED MARKETING CAMPAIGNS.

CONFIDENCE

CONFIDENCE REPRESENTS THE LIKELIHOOD THAT THE CONSEQUENT ITEM (ITEM BOUGHT AFTER) IS PURCHASED GIVEN THE ANTECEDENT ITEM(S) (ITEM BOUGHT BEFORE).

A CONFIDENCE OF 0.6 IMPLIES THAT THE CONSEQUENT ITEM IS PURCHASED IN 60% OF TRANSACTIONS WHERE THE ANTECEDENT ITEM(S) ARE ALSO PRESENT, INDICATING A STRONG POSITIVE RELATIONSHIP.

```
RULES_CONFIDENZ = ASSOCIATION_RULES(FREQUENT_ITEMSETS, METRIC="CONFIDENCE", MIN_THRESHOLD=0.8, NUM_ITEMSETS= 10000000)  
PRINT(RULES_CONFIDENZ.SHAPE)  
RULES_C = RULES_CONFIDENZ.ROUND(3)  
RULES_C.HEAD()
```

(16, 14)

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	representativity	leverage	conviction	zhangs_metric	jaccard	certainty	kulczynski
0	(ALARM CLOCK BAKELIKE GREEN)	(ALARM CLOCK BAKELIKE RED)	0.049	0.053	0.032	0.655	12.467	1.0	0.030	2.743	0.968	0.465	0.635	0.635
1	(ALARM CLOCK BAKELIKE RED)	(ALARM CLOCK BAKELIKE GREEN)	0.053	0.049	0.032	0.616	12.467	1.0	0.030	2.478	0.971	0.465	0.596	0.635
2	(PINK REGENCY TEACUP AND SAUCER)	(GREEN REGENCY TEACUP AND SAUCER)	0.038	0.050	0.031	0.822	16.520	1.0	0.029	5.332	0.976	0.547	0.812	0.721
3	(GREEN REGENCY TEACUP AND SAUCER)	(PINK REGENCY TEACUP AND SAUCER)	0.050	0.038	0.031	0.621	16.520	1.0	0.029	2.538	0.989	0.547	0.606	0.721





LIFT ANALYSIS

- LIFT MEASURES HOW MUCH MORE LIKELY THE CONSEQUENT ITEM(S) ARE PURCHASED WHEN THE ANTECEDENT ITEM(S) ARE PRESENT COMPARED TO WHEN THEY ARE NOT.
- A LIFT VALUE OF 1.0 INDICATES THAT THE ITEMS IN THE CONSEQUENT ARE BOUGHT TOGETHER AS OFTEN AS WOULD BE EXPECTED BY CHANCE.
- RULES_LIFT = ASSOCIATION_RULES(FREQUENT_ITEMSETS, METRIC="LIFT", MIN_THRESHOLD= 2.5, NUM_ITEMSETS=100000000)
- PRINT(RULES_LIFT.SHAPE)
- RULES_L = RULES_LIFT.ROUND(3)
- RULES_L.HEAD()

(174, 14)

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	representativity	leverage	conviction	zhangs_metric	jaccard	certainty	kulczynski	
0	(60 TEATIME FAIRY CAKE CASES)	(PACK OF 72 RETROSPOT CAKE CASES)	0.041	0.065	0.022	0.545	8.337		1.0	0.020	2.055	0.917	0.265	0.513	0.443
1	(PACK OF 72 RETROSPOT CAKE CASES)	(60 TEATIME FAIRY CAKE CASES)	0.065	0.041	0.022	0.340	8.337		1.0	0.020	1.454	0.942	0.265	0.312	0.443
2	(ALARM CLOCK BAKELIKE GREEN)	(ALARM CLOCK BAKELIKE PINK)	0.049	0.039	0.021	0.425	10.953		1.0	0.019	1.672	0.956	0.313	0.402	0.483
3	(ALARM CLOCK BAKELIKE PINK)	(ALARM CLOCK BAKELIKE GREEN)	0.039	0.049	0.021	0.542	10.953		1.0	0.019	2.073	0.945	0.313	0.518	0.483





SUPPORT ANALYSIS

- SUPPORT QUANTIFIES THE FREQUENCY WITH WHICH A RULE OCCURS IN THE DATASET.
- A SUPPORT OF 0.05 MEANS THAT THE RULE OCCURS IN AT LEAST 5% OF TRANSACTIONS, SIGNIFYING ITS SIGNIFICANCE IN THE DATASET.
- RULES_SUPPORT = ASSOCIATION_RULES(FREQUENT_ITEMSETS, METRIC="SUPPORT", MIN_THRESHOLD=0.03,NUM_ITEMSETS=100000000)
- PRINT(RULES_SUPPORT.SHAPE)
- RULES_S = RULES_SUPPORT.ROUND(3)
- RULES_S.HEAD()

(16, 14)

	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	representativity	leverage	conviction	zhangs_metric	jaccard	certainty	kulczynski
0	(ALARM CLOCK BAKELIKE GREEN)	(ALARM CLOCK BAKELIKE RED)	0.049	0.053	0.032	0.655 12.467	1.0	0.030	2.743	0.968	0.465	0.635	0.635	
1	(ALARM CLOCK BAKELIKE RED)	(ALARM CLOCK BAKELIKE GREEN)	0.053	0.049	0.032	0.616 12.467	1.0	0.030	2.478	0.971	0.465	0.596	0.635	
2	(PINK REGENCY TEACUP AND SAUCER)	(GREEN REGENCY TEACUP AND SAUCER)	0.038	0.050	0.031	0.822 16.520	1.0	0.029	5.332	0.976	0.547	0.812	0.721	
3	(GREEN REGENCY TEACUP AND SAUCER)	(PINK REGENCY TEACUP AND SAUCER)	0.050	0.038	0.031	0.621 16.520	1.0	0.029	2.538	0.989	0.547	0.606	0.721	



ANALYSIS OF ASSOCIATION RULE METRICS

- IN THIS SECTION, WE VISUALLY INSPECT THE DISTRIBUTION OF KEY METRICS ASSOCIATED WITH OUR ASSOCIATION RULES, NAMELY CONFIDENCE, LIFT, AND SUPPORT. THESE HISTOGRAMS PROVIDE INSIGHTS INTO THE FREQUENCY AND DISTRIBUTION OF THESE METRICS ACROSS THE EXTRACTED ASSOCIATION RULES.

```
# CREATE SUBPLOTS IN PLOTLY
```

```
FIG = MAKE_SUBPLOTS(ROWS=3, COLS=1, SUBPLOT_TITLES=('DISTRIBUTION OF CONFIDENCE',
    'DISTRIBUTION OF LIFT',
    'DISTRIBUTION OF SUPPORT'))
```

```
# CONFIDENCE HISTOGRAM
```

```
FIG.ADD_TRACE(
    GO.HISTOGRAM(X=RULES_C['CONFIDENCE'], NBINSX=20, HISTNORM='DENSITY', NAME='CONFIDENCE'),
    ROW=1, COL=1
)
```

```
FIG.ADD_TRACE(
    GO.SCATTER(X=RULES_C['CONFIDENCE'], Y=PD.SERIES(RULES_C['CONFIDENCE']).VALUE_COUNTS(NORMALIZE=True).SORT_INDEX(),
    MODE='LINES', LINE=DICT(COLOR='RED'), NAME='KDE'),
    ROW=1, COL=1
)
```

```
# LIFT HISTOGRAM
```

```
FIG.ADD_TRACE(
    GO.HISTOGRAM(X=RULES_L['LIFT'], NBINSX=20, HISTNORM='DENSITY', NAME='LIFT'),
    ROW=2, COL=1
)
```

```
FIG.ADD_TRACE(
    GO.SCATTER(X=RULES_L['LIFT'], Y=PD.SERIES(RULES_L['LIFT']).VALUE_COUNTS(NORMALIZE=True).SORT_INDEX(),
    MODE='LINES', LINE=DICT(COLOR='RED'), NAME='KDE'),
    ROW=2, COL=1
)
```





```
# SUPPORT HISTOGRAM
FIG.ADD_TRACE(
    GO.HISTOGRAM(X=RULES_S['SUPPORT'], NBINSX=20, HISTNORM='DENSITY', NAME='SUPPORT'),
    ROW=3, COL=1
)
FIG.ADD_TRACE(
    GO.SCATTER(X=RULES_S['SUPPORT'], Y=PD.SERIES(RULES_S['SUPPORT']).VALUE_COUNTS(NORMALIZE=True).SORT_INDEX(),
        MODE='LINES', LINE=DICT(COLOR='RED'), NAME='KDE'),
    ROW=3, COL=1
)

# UPDATE LAYOUT
FIG.UPDATE_LAYOUT(HEIGHT=900, WIDTH=800, TITLE_TEXT="DISTRIBUTION OF CONFIDENCE, LIFT, AND SUPPORT",
    SHOWLEGEND=False)

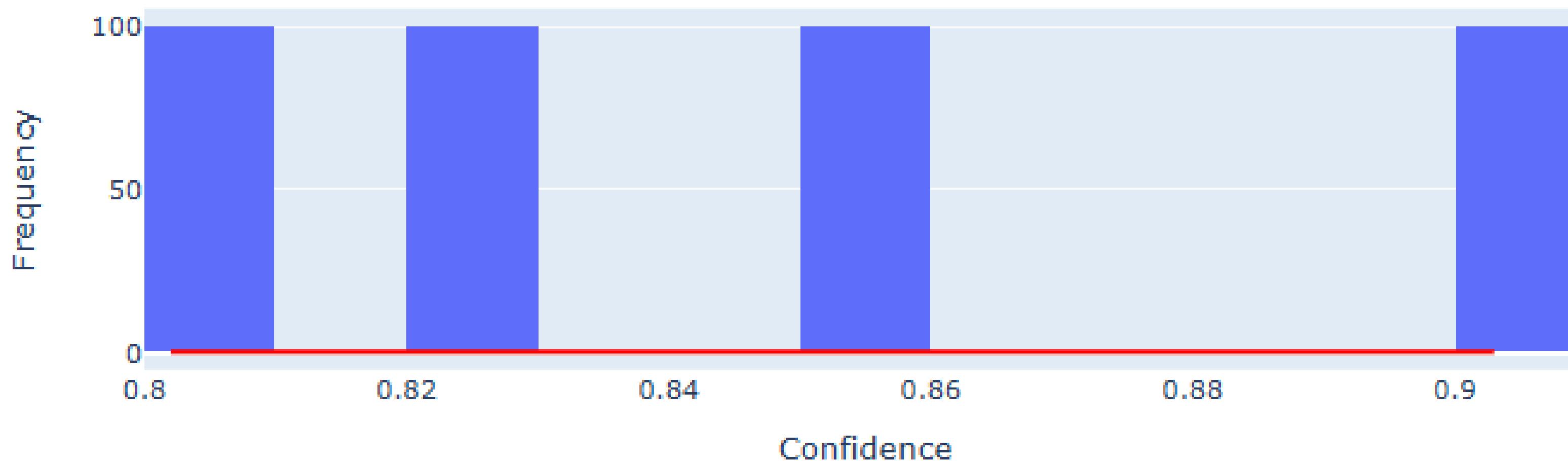
# UPDATE AXES TITLES
FIG.UPDATE_XAXES(TITLE_TEXT="CONFIDENCE", ROW=1, COL=1)
FIG.UPDATE_YAXES(TITLE_TEXT="FREQUENCY", ROW=1, COL=1)
FIG.UPDATE_XAXES(TITLE_TEXT="LIFT", ROW=2, COL=1)
FIG.UPDATE_YAXES(TITLE_TEXT="FREQUENCY", ROW=2, COL=1)
FIG.UPDATE_XAXES(TITLE_TEXT="SUPPORT", ROW=3, COL=1)
FIG.UPDATE_YAXES(TITLE_TEXT="FREQUENCY", ROW=3, COL=1)

FIG.SHOW()
```





Distribution of Confidence

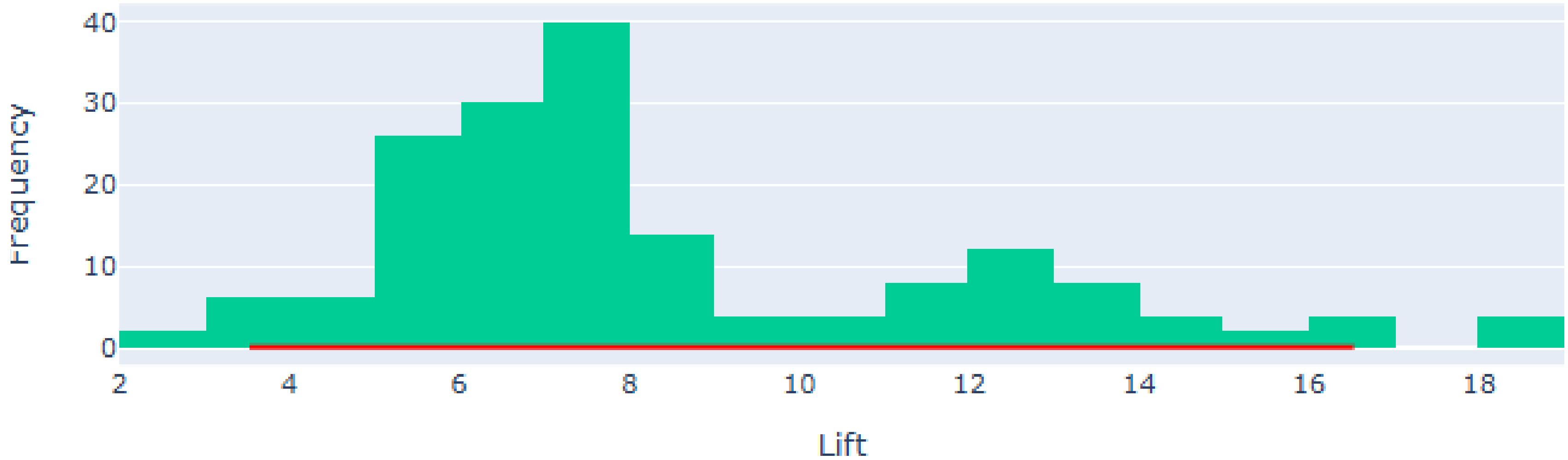


Confidence

- STRONG ASSOCIATION BETWEEN REGENCY TEAWARE: THERE'S A HIGH CONFIDENCE LEVEL (0.826 TO 0.905) BETWEEN DIFFERENT COMBINATIONS OF REGENCY TEACUPS AND SAUCERS. THIS INDICATES THAT CUSTOMERS WHO PURCHASE ONE ITEM FROM THIS CATEGORY ARE HIGHLY LIKELY TO PURCHASE ANOTHER.
- REGENCY TEAWARE AND ROSES: THERE'S A STRONG ASSOCIATION BETWEEN REGENCY TEAWARE AND ROSES. THIS SUGGESTS THAT CUSTOMERS WHO PURCHASE REGENCY TEAWARE ARE LIKELY TO ALSO PURCHASE ROSES.
- JUMBO STORAGE BAG AND RED RETROSPOT: THERE'S A HIGH CONFIDENCE LEVEL (0.802) BETWEEN THE JUMBO STORAGE BAG AND THE RED RETROSPOT. THIS INDICATES A STRONG ASSOCIATION BETWEEN THESE TWO ITEMS.
- ACTIONABLE INSIGHTS:
- CROSS-SELLING OPPORTUNITIES: USE THESE INSIGHTS TO CROSS-SELL RELATED PRODUCTS TO CUSTOMERS. FOR EXAMPLE, IF A CUSTOMER PURCHASES A GREEN REGENCY TEACUP, RECOMMEND A PINK REGENCY SAUCER.
- PRODUCT BUNDLING: CONSIDER BUNDLING COMPLEMENTARY PRODUCTS, SUCH AS THE GREEN AND RED BAKELITE ALARM CLOCKS OR THE DIFFERENT REGENCY TEAWARE ITEMS.
- TARGETED MARKETING: USE THESE INSIGHTS TO CREATE TARGETED MARKETING CAMPAIGNS FOR SPECIFIC CUSTOMER SEGMENTS. FOR EXAMPLE, YOU COULD TARGET CUSTOMERS WHO HAVE PURCHASED REGENCY TEAWARE WITH PROMOTIONS FOR ROSES.
- INVENTORY MANAGEMENT: ENSURE ADEQUATE STOCK LEVELS FOR FREQUENTLY PURCHASED ITEMS TO AVOID STOCKOUTS.

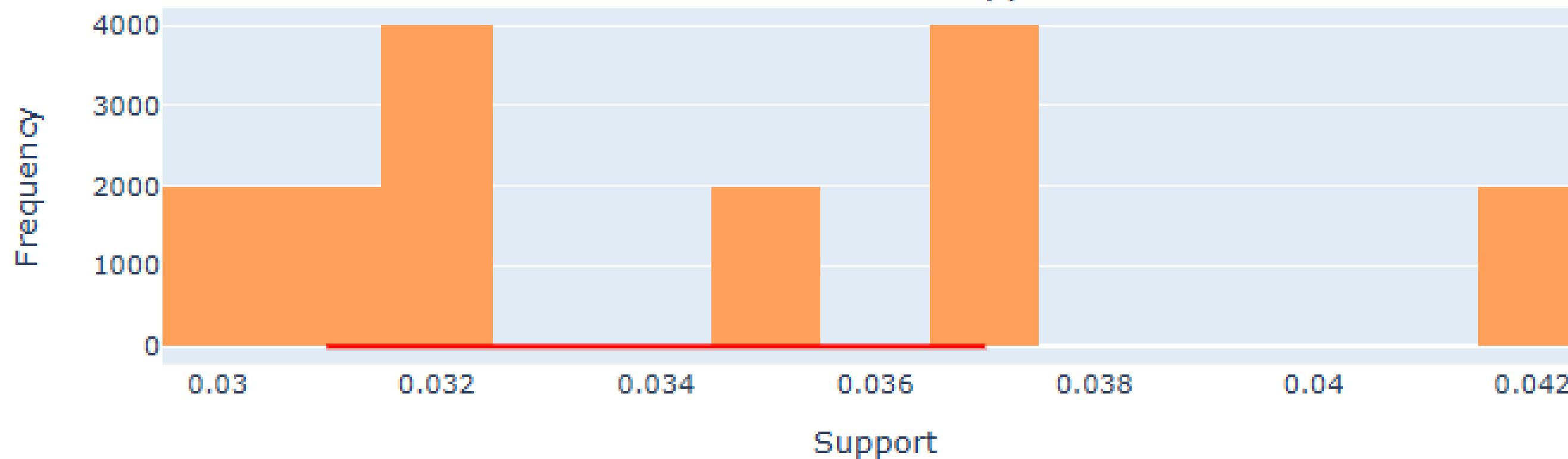


Distribution of Lift



- STRONG ASSOCIATION BETWEEN ALARM CLOCKS: THERE'S A STRONG ASSOCIATION BETWEEN GREEN AND RED BAKELITE ALARM CLOCKS. THIS SUGGESTS A POTENTIAL FOR BUNDLING THESE ITEMS OR CROSS-SELLING TO CUSTOMERS WHO PURCHASE ONE.
- REGENCY TEAWARE AFFINITY: GREEN AND PINK REGENCY TEACUPS AND SAUCERS ARE OFTEN PURCHASED TOGETHER, INDICATING A STRONG CUSTOMER PREFERENCE FOR THESE COLOR COMBINATIONS.
- PRODUCT RECOMMENDATIONS: LEVERAGE THESE INSIGHTS TO RECOMMEND COMPLEMENTARY PRODUCTS TO CUSTOMERS BASED ON THEIR PURCHASE HISTORY. FOR INSTANCE, IF A CUSTOMER BUYS A GREEN BAKELITE ALARM CLOCK, RECOMMEND THE RED ONE AS WELL.
- TARGETED MARKETING CAMPAIGNS: CREATE TARGETED MARKETING CAMPAIGNS FOR SPECIFIC CUSTOMER SEGMENTS BASED ON THEIR PURCHASING BEHAVIOR. FOR EXAMPLE, YOU COULD TARGET CUSTOMERS WHO HAVE PURCHASED GREEN REGENCY TEACUPS WITH PROMOTIONS FOR PINK SAUCERS.
- INVENTORY MANAGEMENT: OPTIMIZE INVENTORY LEVELS TO ENSURE SUFFICIENT STOCK OF FREQUENTLY PURCHASED ITEMS.
-

Distribution of Support



- POPULAR PRODUCTS: THE HIGHER THE SUPPORT VALUE, THE MORE FREQUENTLY THE ITEMSET IS PURCHASED. FOR EXAMPLE, THE ITEMSET OF "GREEN REGENCY TEACUP AND SAUCER" AND "PINK REGENCY TEACUP AND SAUCER" HAS A RELATIVELY HIGH SUPPORT VALUE, INDICATING THAT THESE ITEMS ARE FREQUENTLY PURCHASED TOGETHER.
- INVENTORY MANAGEMENT: UNDERSTANDING THE SUPPORT OF DIFFERENT ITEMSETS CAN HELP OPTIMIZE INVENTORY LEVELS. ITEMS WITH HIGH SUPPORT VALUES SHOULD BE STOCKED IN SUFFICIENT QUANTITIES TO AVOID STOCKOUTS.
- PRODUCT BUNDLING: ITEMS WITH HIGH SUPPORT VALUES CAN BE BUNDLED TOGETHER TO CREATE ATTRACTIVE OFFERS FOR CUSTOMERS. FOR EXAMPLE, YOU COULD BUNDLE THE GREEN AND PINK REGENCY TEACUPS AND SAUCERS.

Data Analysis

- QUANTITY OVERVIEW:
- MEAN: THE AVERAGE QUANTITY PER TRANSACTION IS APPROXIMATELY 10 ITEMS, BUT THE DATA SHOWS EXTREME OUTLIERS:
 - MINIMUM: -9600, WHICH IS LIKELY AN ERROR (NEGATIVE QUANTITIES ARE UNUSUAL UNLESS REPRESENTING RETURNS OR ERRORS).
 - MAXIMUM: 80995, WHICH IS AN EXTREMELY HIGH VALUE AND MAY NEED VALIDATION.
- PRICE ANALYSIS:
- MEAN: THE AVERAGE PRICE IS 3.84, BUT THE DATA INCLUDES NEGATIVE AND ABNORMALLY HIGH VALUES:
 - MINIMUM: -11062.06, WHICH MAY INDICATE INCORRECT PRICING OR REFUND DATA.
 - MAXIMUM: 13541.33, WHICH MIGHT REPRESENT LUXURY PRODUCTS OR INCORRECT ENTRIES.
- THE STANDARD DEVIATION (42.11) INDICATES A WIDE VARIATION IN PRODUCT PRICES.
- CUSTOMERID DATA:
 - **AROUND 134,000 CUSTOMERID VALUES ARE MISSING (516778 - 382811)**, WHICH MIGHT COMPLICATE CUSTOMER-LEVEL INSIGHTS.
- DATE RANGE:
 - THE TRANSACTIONS SPAN FROM 2010-12-01 TO 2011-12-09, COVERING ABOUT ONE YEAR

INSIGHTS

1. SEASONAL TRENDS

- PEAK IN NOVEMBER: INDICATES STRONG SEASONAL DEMAND, LIKELY TIED TO HOLIDAYS OR END-OF-YEAR ACTIVITIES.
- LOW PERFORMANCE IN Q1: FEBRUARY RECORDS THE LOWEST SALES, OFTEN DUE TO POST-HOLIDAY SLOWDOWNS.
- GROWTH OPPORTUNITIES: STEADY GROWTH FROM SEPTEMBER TO NOVEMBER SIGNALS EFFECTIVE CAMPAIGNS OR SEASONAL PRODUCT RELEVANCE.

2. UNIQUE ITEM PATTERNS

- HIGH VARIETY IN Q4: THE DEMAND FOR DIVERSE ITEMS PEAKS FROM SEPTEMBER TO NOVEMBER, SHOWING CUSTOMER INTEREST IN UNIQUE OFFERINGS.
- STABLE MIDDLE PERIOD: MARCH TO AUGUST REFLECTS CONSISTENT BUT LESS DYNAMIC DEMAND.

3. ASSOCIATION RULE ANALYSIS

- FREQUENT ITEMSETS: ITEMS LIKE REGENCY TEACUPS AND ROSES, AND GREEN/RED BAKELITE ALARM CLOCKS, SHOW STRONG ASSOCIATIONS.
- CROSS-SELLING POTENTIAL: HIGH CONFIDENCE BETWEEN CERTAIN PRODUCT PAIRS, E.G., JUMBO STORAGE BAGS AND RED RETROSPOT.

4. INVENTORY CHALLENGES

- OUTLIERS IN DATA: EXTREME VALUES IN QUANTITY (-9600 AND 80995) AND PRICE (-11062.06 AND 13541.33) SUGGEST DATA ERRORS OR UNIQUE EVENTS, EMPHASIZING THE NEED FOR VALIDATION.
- STOCKING NEEDS: ITEMS WITH HIGH SUPPORT VALUES (FREQUENT PURCHASE) REQUIRE ROBUST INVENTORY MANAGEMENT TO PREVENT STOCKOUTS.

RECOMMENDATIONS

1. STRATEGIC CAMPAIGNS

- Q1 PROMOTIONS: IMPLEMENT EARLY-YEAR DISCOUNTS OR INTRODUCE NEW PRODUCTS TO COUNTERACT THE POST-HOLIDAY SLUMP.
- BOOST NOVEMBER SALES: REPLICATE SUCCESSFUL STRATEGIES FROM NOVEMBER DURING OTHER POTENTIAL PEAK MONTHS.
- TARGETED MARKETING: FOCUS CAMPAIGNS ON CUSTOMER SEGMENTS PURCHASING HIGH-SUPPORT ITEMS.

2. PRODUCT BUNDLING

- BUNDLE HIGHLY ASSOCIATED ITEMS (E.G., GREEN AND PINK REGENCY TEACUPS) TO ENCOURAGE LARGER TRANSACTIONS.
- INTRODUCE HOLIDAY-THEMED BUNDLES IN Q4 TO CAPITALIZE ON HEIGHTENED DEMAND.

3. INVENTORY MANAGEMENT

- MAINTAIN SUFFICIENT STOCK FOR HIGH-SUPPORT ITEMS LIKE GREEN REGENCY TEACUP & SAUCER TO AVOID STOCKOUTS.
- USE TRENDS FROM MARKET BASKET ANALYSIS TO ALIGN INVENTORY WITH EXPECTED DEMAND.

4. DATA VALIDATION

- ADDRESS DATA ANOMALIES (E.G., EXTREME NEGATIVE QUANTITIES AND PRICES) TO IMPROVE ANALYSIS RELIABILITY.
- ENRICH THE DATASET WITH CUSTOMERID DETAILS TO ENHANCE CUSTOMER SEGMENTATION AND PERSONALIZATION.

5. VISUALIZATION AND MONITORING

- TO MONITOR SALES TRENDS, PRODUCT PERFORMANCE, AND INVENTORY LEVELS IN REAL TIME.
- REGULARLY REVIEW METRICS SUCH AS LIFT, CONFIDENCE, AND SUPPORT TO REFINE MARKETING AND INVENTORY STRATEGIES.

6. SEASONAL PRODUCT INSIGHTS

- INCREASE THE AVAILABILITY OF DIVERSE PRODUCTS DURING SEPTEMBER TO NOVEMBER TO CAPTURE HEIGHTENED DEMAND.
- ANALYZE LOW-PERFORMING MONTHS TO IDENTIFY ANY OVERLOOKED OPPORTUNITIES OR NECESSARY ADJUSTMENTS.

CONCLUSION



- THE MARKET BASKET ANALYSIS (MBA) PROVIDED VALUABLE INSIGHTS INTO CUSTOMER PURCHASING PATTERNS AND PRODUCT RELATIONSHIPS WITHIN THE TRANSACTION DATASET. THE FINDINGS REVEAL SIGNIFICANT SEASONAL TRENDS, INCLUDING PEAKS IN NOVEMBER AND SLUMPS IN FEBRUARY, EMPHASIZING THE IMPORTANCE OF TIMING IN SALES STRATEGIES. THE ASSOCIATION RULE ANALYSIS IDENTIFIED STRONG RELATIONSHIPS BETWEEN SPECIFIC PRODUCTS, PRESENTING OPPORTUNITIES FOR CROSS-SELLING, PRODUCT BUNDLING, AND TARGETED MARKETING.
- DATA ANOMALIES, SUCH AS EXTREME OUTLIERS IN QUANTITY AND PRICE, HIGHLIGHT THE NECESSITY FOR ROBUST DATA VALIDATION AND CLEANING PROCESSES TO ENSURE ACCURATE ANALYSIS. EFFECTIVE INVENTORY MANAGEMENT, INFORMED BY HIGH-SUPPORT ITEMSETS, IS CRITICAL TO MEETING CUSTOMER DEMAND AND AVOIDING STOCKOUTS.
- BY LEVERAGING THESE INSIGHTS, THE BUSINESS CAN IMPLEMENT STRATEGIC CAMPAIGNS, OPTIMIZE INVENTORY, AND ENHANCE CUSTOMER SATISFACTION, ULTIMATELY DRIVING REVENUE GROWTH AND BUILDING STRONGER CUSTOMER RELATIONSHIPS. TO MAXIMIZE IMPACT, IT IS ESSENTIAL TO CONTINUALLY MONITOR PERFORMANCE, ADAPT STRATEGIES BASED ON TRENDS, AND REFINE OPERATIONS TO ALIGN WITH EVOLVING MARKET DYNAMICS.



The background image shows a panoramic aerial view of the Dubai Marina skyline during sunset. The sky is filled with large, billowing clouds illuminated by the warm orange and yellow light of the setting sun. In the foreground, the calm waters of the Persian Gulf reflect the sunlight. On the left, the iconic skyscrapers of the Marina Bay area stand tall, with the Burj Khalifa being the most prominent. To the right, the Ferris wheel of the Dubai Eye is visible across the water. The overall atmosphere is one of luxury and modernity.

Thank you