

MARKET BASKET INSIGHTS

Problem Statement:

The problem is to analyze customer purchase data to discover associations between products, understand buying patterns, and extract insights to improve marketing and sales strategies.

Design Thinking Process:

Empathize: Understand the business and customer needs, gather data on sales and customer transactions.

Define: Define the problem and objectives, such as increasing cross-selling, optimizing product placement, or improving inventory management.

Ideate: Generate ideas on what kind of associations might exist and how they can be used to achieve the objectives.

Prototype: Prepare the dataset and choose appropriate tools and algorithm for analysis.

Test: Apply association analysis techniques and evaluate their effectiveness.

Implement: Translate insights into actionable strategies and monitor their impact on sales and customer behavior.

Phases of Development:

Data Collection: Gather transaction data, including information on products, customers, and purchase history.

Data Preprocessing:

- Handle missing data, outliers, and duplicates.
- Convert data into a suitable format for association analysis.
- Encode categorical variables.
- Remove irrelevant or infrequent items to reduce noise.

Association Analysis Techniques:

- Use algorithms like Apriori or FP-growth to discover frequent itemsets and association rules.

- Set appropriate support and confidence thresholds to filter out relevant rules.
- Generate association rules that reveal item co-occurrence patterns.

Interpretation:

- Examine the discovered association rules, which may include items that tend to be bought together.
- Identify strong and actionable rules based on high confidence and lift values

Business Implications:

- Implement strategies like bundling products frequently bought together.
- Optimize store layouts based on product associations.
- Create targeted marketing campaigns.
- Enhance inventory management to ensure product availability for frequently associated items.
- Monitor changes in customer behavior and adapt strategies accordingly.

Documentation:

- Maintain clear documentation of the entire process, including data sources, preprocessing steps, analysis results, and business actions taken.

Dataset

File name: Assignment-1_Data

List name: retaildata

File format: .xlsx

Number of Row: 522065

Number of Attributes: 7

BillNo: 6-digit number assigned to each transaction. Nominal.

Itemname: Product name. Nominal.

Quantity: The quantities of each product per transaction. Numeric.

Date: The day and time when each transaction was generated. Numeric.

Price: Product price. Numeric.

CustomerID: 5-digit number assigned to each customer. Nominal.

Country: Name of the country where each customer resides. Nominal

Data Preprocessing:

- **Data Collection:** Gather transaction data, typically in the form of a dataset with rows representing transactions and columns representing items
- **Data Cleaning:** Remove duplicates, missing values, and irrelevant information. Ensure consistent formatting and encoding of items.
- **Data Transformation:** Convert data into a suitable format, like a binary matrix where rows are transactions, and columns are items with 1s and 0s representing presence or absence.
- **Support Threshold:** Set a minimum support threshold, which determines how frequent an itemset must be to be considered in the analysis. This helps in reducing the size of the itemset space.

Association Analysis Techniques:

This is a classic algorithm for finding frequent itemsets. It uses a level-wise approach to iteratively discover itemsets with higher support.

- **FP-Growth Algorithm:** This is another frequent itemset mining algorithm that uses a tree structure to efficiently discover frequent itemsets without the need for multiple database scans.
- **Association Rule Generation:** After identifying frequent itemsets, association rules are generated. These rules consist of an antecedent (if) and a consequent (then) part. For example, "If a customer buys item A, then they are likely to buy item B."
- **Confidence and Lift:** These metrics help assess the strength of association rules. Confidence measures how often the rule is correct, while lift measures the significance of the rule's performance compared to random chance.
- **Pruning and Post-Processing:** Remove redundant or less interesting rules to keep the analysis focused on valuable insights.
- **Visualization:** Use data visualization techniques like heatmaps, scatter plots, or network graphs to represent associations and make them more interpretable.

Discovered association rules and their business implications:

Discovering association rules from market basket data like this can help businesses understand relationships between products in customer transactions. Association rules typically have the format “If {antecedent}, then {consequent}” and are based on the support, confidence, and lift metrics. Here’s how you might interpret the discovered association rules and their implications:

- **High Support Rules:** Rules with high support indicate that the combination of items in the antecedent and consequent frequently occurs together. This implies a strong association between these items.
- **High Confidence Rules:** High confidence indicates that when the antecedent items are in the basket, there's a high likelihood of the consequent items also being in the basket.
- **High Lift Rules:** Lift measures how much more likely it is for the consequent to be bought when the antecedent is purchased, compared to when it's bought without the antecedent. A lift greater than 1 suggests a positive association.

Business implications:

- **Product Placement:** Retailers can use these rules to optimize the placement of products in the store. If products frequently bought together are placed close to each other, it can increase sales.
- **Pricing Strategies:** Knowing which products are commonly purchased together can inform pricing strategies. For example, bundling two associated products might attract more sales.
- **Customer Segmentation:** Understanding the association rules can help in customer segmentation. Retailers can target specific customer groups with tailored promotions based on their common purchase patterns.
- **Cross-selling and Upselling:** By identifying associations, businesses can effectively cross-sell and upsell products. For instance, if a customer adds a laptop to their cart, a suggestion for a laptop bag or accessories can be made.
- **Market Expansion:** If an association rule suggests that a certain combination of products is popular in one country but not in another, it can inform market expansion strategies.

Jupyter Notebook:

```
import pandas as pd
import numpy as np

#for viz
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline

#to avoid warning
import warnings
warnings.filterwarnings('ignore')

#to display all feature if the number increase
pd.set_option('display.max_columns', None)

data=pd.read_excel('C:\\Users\\ELCOT\\Downloads\\archive11\\Assignment-1_Data
.xlsx')

import os

os.getcwd()

'C:\\Users\\ELCOT'

data
```

	BillNo	Itemname	Quantity	\
0	536365	WHITE HANGING HEART T-LIGHT HOLDER	6	
1	536365	WHITE METAL LANTERN	6	
2	536365	CREAM CUPID HEARTS COAT HANGER	8	
3	536365	KNITTED UNION FLAG HOT WATER BOTTLE	6	
4	536365	RED WOOLLY HOTTIE WHITE HEART.	6	
...	
522059	581587	PACK OF 20 SPACEBOY NAPKINS	12	
522060	581587	CHILDREN'S APRON DOLLY GIRL	6	
522061	581587	CHILDRENS CUTLERY DOLLY GIRL	4	
522062	581587	CHILDRENS CUTLERY CIRCUS PARADE	4	
522063	581587	BAKING SET 9 PIECE RETROSPOT	3	

	Date	Price	CustomerID	Country
0	2010-12-01 08:26:00	2.55	17850.0	United Kingdom
1	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
2	2010-12-01 08:26:00	2.75	17850.0	United Kingdom
3	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
4	2010-12-01 08:26:00	3.39	17850.0	United Kingdom
...

522059	2011-12-09	12:50:00	0.85	12680.0	France
522060	2011-12-09	12:50:00	2.10	12680.0	France
522061	2011-12-09	12:50:00	4.15	12680.0	France
522062	2011-12-09	12:50:00	4.15	12680.0	France
522063	2011-12-09	12:50:00	4.95	12680.0	France

[522064 rows x 7 columns] **Apriori Algorithm:**

```
# Display basic information about the dataset
print("Number of rows and columns:", data.shape)
print("\nData Types and Missing Values:")
print(data.info())
print("\nFirst few rows of the dataset:")
print(data.head())
```

Number of rows and columns: (522064, 7)

Data Types and Missing Values:

```
<class 'pandas.core.frame.DataFrame'>
```

RangeIndex: 522064 entries, 0 to 522063

Data columns (total 7 columns):

#	Column	Non-Null Count	Dtype
0	BillNo	522064 non-null	object
1	Itemname	520609 non-null	object
2	Quantity	522064 non-null	int64
3	Date	522064 non-null	datetime64[ns]
4	Price	522064 non-null	float64
5	CustomerID	388023 non-null	float64
6	Country	522064 non-null	object

dtypes: datetime64[ns](1), float64(2), int64(1), object(3)

memory usage: 27.9+ MB

None

First few rows of the dataset:

	BillNo	Itemname	Quantity	Date
0	536365	WHITE HANGING HEART T-LIGHT HOLDER	6	2010-12-01 08:26:00
1	536365	WHITE METAL LANTERN	6	2010-12-01 08:26:00
2	536365	CREAM CUPID HEARTS COAT HANGER	8	2010-12-01 08:26:00
3	536365	KNITTED UNION FLAG HOT WATER BOTTLE	6	2010-12-01 08:26:00
4	536365	RED WOOLLY HOTTIE WHITE HEART.	6	2010-12-01 08:26:00

	Price	CustomerID	Country
0	2.55	17850.0	United Kingdom
1	3.39	17850.0	United Kingdom
2	2.75	17850.0	United Kingdom
3	3.39	17850.0	United Kingdom
4	3.39	17850.0	United Kingdom

```

#Check Missing Values
print("Missing Values:")
print(data.isnull().sum())

#Drop Rows with Missing Values
data.dropna(inplace=True)

Missing Values:
BillNo          0
Itemname       1455
Quantity        0
Date            0
Price           0
CustomerID    134041
Country         0
dtype: int64

# Convert dataframe into transaction data
transaction_data = data.groupby(['BillNo', 'Date'])['Itemname'].apply(lambda
x: ', '.join(x)).reset_index()

#Drop Unnecessary Columns
columns_to_drop = ['BillNo', 'Date']
transaction_data.drop(columns=columns_to_drop, inplace=True)

# Save the transaction data to a CSV file
transaction_data_path = 'C:\\Users\\ELCOT\\Downloads\\archive11\\transaction_
data.csv'
transaction_data.to_csv(transaction_data_path, index=False)

# Display the first few rows of the transaction data
print("\nTransaction Data for Association Rule Mining:")
print(transaction_data.head())
transaction_data.shape

```

Transaction Data for Association Rule Mining:

	0		1 \
0	WHITE HANGING HEART T-LIGHT HOLDER	WHITE METAL LANTERN	
1	HAND WARMER UNION JACK	HAND WARMER RED POLKA DOT	
2	ASSORTED COLOUR BIRD ORNAMENT	POPPY'S PLAYHOUSE BEDROOM	
3	JAM MAKING SET WITH JARS	RED COAT RACK PARIS FASHION	
4	BATH BUILDING BLOCK WORD	None	
	2		3 \
0	CREAM CUPID HEARTS COAT HANGER	KNITTED UNION FLAG HOT WATER BOTTLE	
1	None	None	
2	POPPY'S PLAYHOUSE KITCHEN	FELTCRAFT PRINCESS CHARLOTTE DOLL	
3	YELLOW COAT RACK PARIS FASHION	BLUE COAT RACK PARIS FASHION	
4	None	None	

		4		5	\
0	RED WOOLLY HOTTIE WHITE HEART.		SET 7 BABUSHKA NESTING BOXES		
1		None		None	
2	IVORY KNITTED MUG COSY		BOX OF 6 ASSORTED COLOUR TEASPOONS		
3		None		None	
4		None		None	

	6	7
0	GLASS STAR FROSTED T-LIGHT HOLDER	None
1	None	None
2	BOX OF VINTAGE JIGSAW BLOCKS	BOX OF VINTAGE ALPHABET BLOCKS
3	None	None
4	None	None

	8	9	\
0	None	None	
1	None	None	
2	HOME BUILDING BLOCK	WORD LOVE BUILDING BLOCK	WORD
3	None	None	
4	None	None	

				10		11	12	13	14	15
\										
0				None		None	None	None	None	None
1				None		None	None	None	None	None
2	RECIPE	BOX	WITH	METAL	HEART	DOORMAT	NEW	ENGLAND	None	None
3				None		None	None	None	None	None
4				None		None	None	None	None	None

[illegible][illegible][illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

4	None	None	None	None	None	None	None	None	None	None	None	None	
	472	473	474	475	476	477	478	479	480	481	482	483	\
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	
	484	485	486	487	488	489	490	491	492	493	494	495	\
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	
	496	497	498	499	500	501	502	503	504	505	506	507	\
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	
	508	509	510	511	512	513	514	515	516	517	518	519	\
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	
	520	521	522	523	524	525	526	527	528	529	530	531	\
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	
	532	533	534	535	536	537	538	539	540	541	542	543	
0	None	None	None	None	None	None	None	None	None	None	None	None	
1	None	None	None	None	None	None	None	None	None	None	None	None	
2	None	None	None	None	None	None	None	None	None	None	None	None	
3	None	None	None	None	None	None	None	None	None	None	None	None	
4	None	None	None	None	None	None	None	None	None	None	None	None	

(18192, 544)

Split the 'Itemname' column into individual items

`items_df = transaction_data['Itemname'].str.split(', ', expand=True)`

Concatenate the original DataFrame with the new items DataFrame

```
transaction_data = pd.concat([transaction_data, items_df], axis=1)
```

```
# Drop the original 'Itemname' column
```

```
transaction_data = transaction_data.drop('Itemname', axis=1)
```

```
# Display the resulting DataFrame
```

```
print(transaction_data.head())
```

	0	1 \
0	WHITE HANGING HEART T-LIGHT HOLDER	WHITE METAL LANTERN
1	HAND WARMER UNION JACK	HAND WARMER RED POLKA DOT
2	ASSORTED COLOUR BIRD ORNAMENT	POPPY'S PLAYHOUSE BEDROOM
3	JAM MAKING SET WITH JARS	RED COAT RACK PARIS FASHION
4	BATH BUILDING BLOCK WORD	None

	2	3 \
0	CREAM CUPID HEARTS COAT HANGER	KNITTED UNION FLAG HOT WATER BOTTLE
1	None	None
2	POPPY'S PLAYHOUSE KITCHEN	FELTCRAFT PRINCESS CHARLOTTE DOLL
3	YELLOW COAT RACK PARIS FASHION	BLUE COAT RACK PARIS FASHION
4	None	None

	4	5 \
0	RED WOOLLY HOTTIE WHITE HEART.	SET 7 BABUSHKA NESTING BOXES
1	None	None
2	IVORY KNITTED MUG COSY	BOX OF 6 ASSORTED COLOUR TEASPOONS
3	None	None
4	None	None

	6	7 \
0	GLASS STAR FROSTED T-LIGHT HOLDER	None
1	None	None
2	BOX OF VINTAGE JIGSAW BLOCKS	BOX OF VINTAGE ALPHABET BLOCKS
3	None	None
4	None	None

	8	9 \
0	None	None
1	None	None
2	HOME BUILDING BLOCK WORD	LOVE BUILDING BLOCK WORD
3	None	None
4	None	None

	10	11	12	13	14	15
\						
0	None	None	None	None	None	None
1	None	None	None	None	None	None
2	RECIPE BOX WITH METAL HEART	DOORMAT NEW ENGLAND	None	None	None	None
3	None	None	None	None	None	None

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

4	None	None	None	None	None	None	None	None	None	None	None	None
	520	521	522	523	524	525	526	527	528	529	530	531 \
0	None	None	None	None	None	None	None	None	None	None	None	None
1	None	None	None	None	None	None	None	None	None	None	None	None
2	None	None	None	None	None	None	None	None	None	None	None	None
3	None	None	None	None	None	None	None	None	None	None	None	None
4	None	None	None	None	None	None	None	None	None	None	None	None
	532	533	534	535	536	537	538	539	540	541	542	543
0	None	None	None	None	None	None	None	None	None	None	None	None
1	None	None	None	None	None	None	None	None	None	None	None	None
2	None	None	None	None	None	None	None	None	None	None	None	None
3	None	None	None	None	None	None	None	None	None	None	None	None
4	None	None	None	None	None	None	None	None	None	None	None	None

Convert items to boolean columns

```
df_encoded = pd.get_dummies(transaction_data, prefix='', prefix_sep='').group
by(level=0, axis=1).max()
```

Save the transaction data to a CSV file

```
df_encoded.to_csv('transaction_data_encoded.csv', index=False)
```

Load transaction data into a DataFrame

```
df_encoded = pd.read_csv('transaction_data_encoded.csv')
```

```
from mlxtend.frequent_patterns import apriori, association_rules
```

Association Rule Mining

```
frequent_itemsets = apriori(df_encoded, min_support=0.007, use_colnames=True)
rules = association_rules(frequent_itemsets, metric="confidence", min_thres
ld=0.5)
```

Display information of the rules

```
print("Association Rules:")
print(rules.head())
```

Association Rules:

	Antecedents	consequents \
0	(CHOCOLATE BOX RIBBONS)	(6 RIBBONS RUSTIC CHARM)
1	(60 CAKE CASES DOLLY GIRL DESIGN)	(PACK OF 72 RETROSPOT CAKE CASES)
2	2 (60 TEATIME FAIRY CAKE CASES)	(PACK OF 72 RETROSPOT CAKE CASES)
3	(ALARM CLOCK BAKELIKE CHOCOLATE)	(ALARM CLOCK BAKELIKE GREEN)

4 (ALARM CLOCK BAKELIKE CHOCOLATE) (ALARM CLOCK BAKELIKE PINK)

Antecedent support consequent support support confidence lift \

0	0.012368	0.039193	0.007036	0.568889	14.515044
1	0.018525	0.054529	0.010059	0.543027	9.958409
2	0.034631	0.054529	0.017315	0.500000	9.169355
3	0.017150	0.042931	0.011379	0.663462	15.454151
4	0.017150	0.032652	0.009125	0.532051	16.294742

Leverage conviction zhangs_metric

0	0.006551	2.228676	0.942766
1	0.009049	2.068984	0.916561
2	0.015427	1.890941	0.922902
3	0.010642	2.843862	0.951613
4	0.008565	2.067210	0.955009

```
import matplotlib.pyplot as plt
import seaborn as sns
```

```
# Plot scatterplot for Support vs. Confidence
```

```
plt.figure(figsize=(12, 8))
```

```
sns.scatterplot(x="support", y="confidence", size="lift", data=rules, hue="lift", palette="viridis", sizes=(20, 200))
```

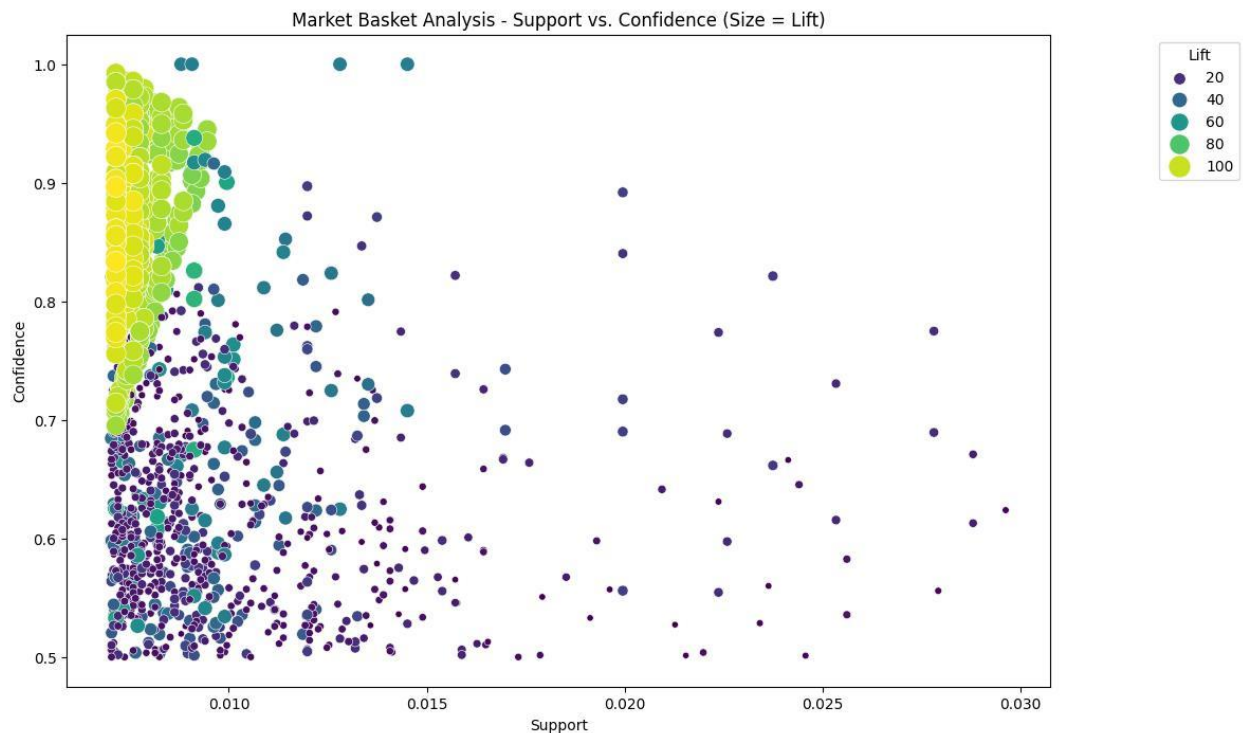
```
plt.title('Market Basket Analysis - Support vs. Confidence (Size = Lift)')
```

```
plt.xlabel('Support')
```

```
plt.ylabel('Confidence')
```

```
plt.legend(title='Lift', loc='upper right', bbox_to_anchor=(1.2, 1))
```

```
plt.show()
```



```
import plotly.express as px

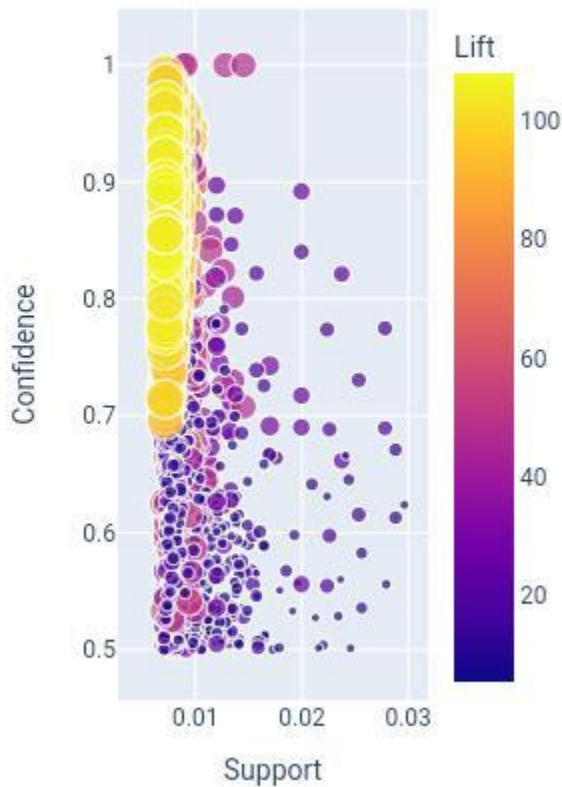
# Convert frozensets to lists for serialization
rules['antecedents'] = rules['antecedents'].apply(list)
rules['consequents'] = rules['consequents'].apply(list)

# Create an interactive scatter plot using plotly express
fig = px.scatter(rules, x="support", y="confidence", size="lift",
                 color="lift", hover_name="consequents",
                 title='Market Basket Analysis - Support vs. Confidence',
                 labels={'support': 'Support', 'confidence': 'Confidence'})

# Customize the layout
fig.update_layout(xaxis_title='Support', yaxis_title='Confidence', coloraxis_colorbar_title='Lift', showlegend=True)

# Show the interactive plot
fig.show()
```


Market Basket Analysis - Support vs. Cor



```
import networkx as nx
import matplotlib.pyplot as plt
import plotly.graph_objects as go

# Create a directed graph
G = nx.DiGraph()

# Add nodes and edges from association rules
for idx, row in rules.iterrows():
    G.add_node(tuple(row['antecedents']), color='skyblue')
    G.add_node(tuple(row['consequents']), color='orange')
    G.add_edge(tuple(row['antecedents']), tuple(row['consequents']), weight=row['support'])

# Set node positions using a spring layout
pos = nx.spring_layout(G)
```

```
# Create an interactive plot using plotly
```

```
edge_x = []
```

```
edge_y = []
```

```
for edge in G.edges(data=True):
```

```
    x0, y0 = pos[edge[0]]
```

```
    x1, y1 = pos[edge[1]]
```

```
    edge_x.append(x0)
```

```
    edge_x.append(x1)
```

```
    edge_x.append(None)
```

```
    edge_y.append(y0)
```

```
    edge_y.append(y1)
```

```
    edge_y.append(None)
```

```
edge_trace = go.Scatter(
```

```
    x=edge_x, y=edge_y,
```

```
    line=dict(width=0.5, color='#888'),
```

```
    hoverinfo='none',
```

```
    mode='lines')
```

```
node_x = []
```

```
node_y = []
```

```
for node in G.nodes():
```

```
    x, y = pos[node]
```

```
    node_x.append(x)
```

```
    node_y.append(y)
```

```
node_trace = go.Scatter(
```

```
    x=node_x, y=node_y,
```

```
    mode='markers',
```

```
    hoverinfo='text',
```

```
    marker=dict(
```

```
        showscale=True,
```

```
        colorscale='YlGnBu',
```

```
        size=10,
```

```
        colorbar=dict(
```

```
            thickness=15,
```

```
            title='Node Connections',
```

```
            xanchor='left',
```

```
            titleside='right'
```

```
        )
```

```
    )
```

```
)
```

```
# Customize the layout
```

```
layout = go.Layout(
```

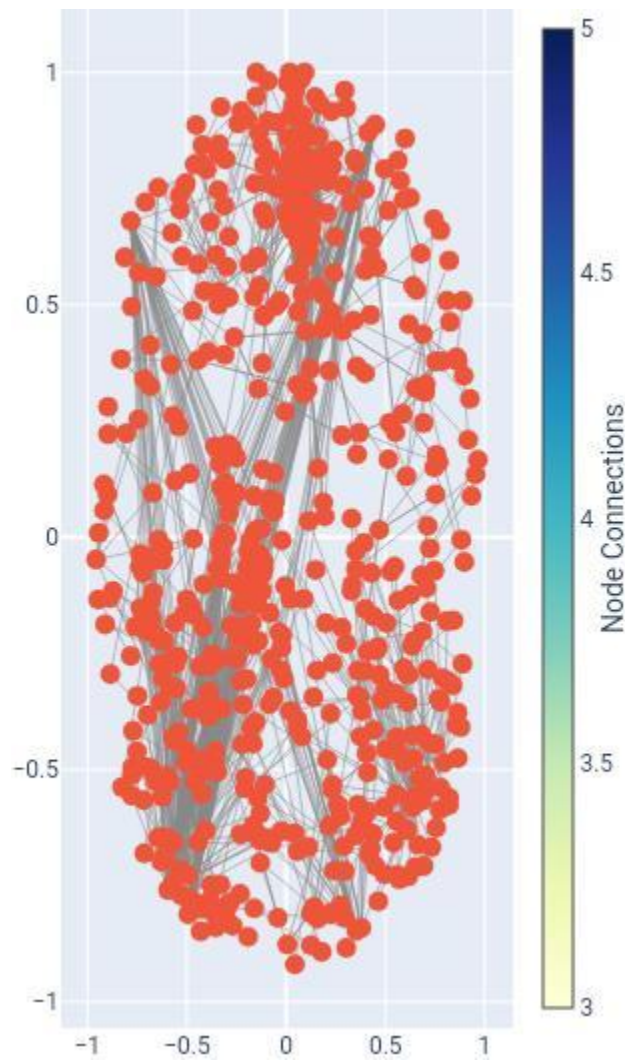
```
    showlegend=False,
```

```
    hovermode='closest',
```

```
    margin=dict(b=0, l=0, r=0, t=0),
```

```
)
```

```
# Create the figure  
fig = go.Figure(data=[edge_trace, node_trace], layout=layout)  
  
# Show the interactive graph  
fig.show()
```



Conclusion

Based on the results of these calculations can be used as a recommendation for retail owners to arrange the arrangement of product catalogs and take strategic steps to improve product marketing.. By utilizing the association rules which are discovered as a result of the analyses, the retailer can apply effective marketing and sales promotion strategies, he will be able to increase customer engagement and improve customer experience and identify customer behavior.