# PHASE-3 MARKET BASKET INSIGHTS

## Market Basket Analysis Project

### Overview

This notebook is part of a project focused on market basket analysis. We will begin by loading and preprocessing the dataset.

### Dataset Information

The dataset is stored in the file `Assignment-1_Data.xlsx` located at `/kaggle/input/market-basket-analysis/`. It contains information related to market transactions.

```
In[1]:
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
/kaggle/input/market-basket-analysis/Assignment-1_Data.xlsx
/kaggle/input/market-basket-analysis/Assignment-1_Data.csv
```

### Loading the Dataset

Let's start by loading the dataset into a DataFrame using pandas.

```
In[2]:
import pandas as pd

# Load the dataset
dataset_path = '/kaggle/input/market-basket-analysis/Assignment-1_Data.xlsx'
df = pd.read_excel(dataset_path)
```

## Initial Exploration

We'll perform an initial exploration of the dataset to understand its structure and characteristics.

```
In[3]:
print("Number of rows and columns:", df.shape)
print("\nData Types and Missing Values:")
print(df.info())
print("\nFirst few rows of the dataset:")
print(df.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
```

# Preprocessing

We'll preprocess the data to ensure it's ready for analysis.

In [4]:
```python
#Check Missing Values
print("Missing Values:")
print(df.isnull().sum())

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

```
Missing Values:
BillNo             0
Itemname        1455
Quantity           0
Date               0
Price              0
```

```
CustomerID    134041
Country            0
dtype: int64
```

In[5]:
```python
# Convert dataframe into transaction data
transaction_data = df.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 = '/kaggle/working/transaction_data.csv'
transaction_data.to_csv(transaction_data_path, index=False)
```

In[6]:
```python
# 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:
                                          Itemname
0  WHITE HANGING HEART T-LIGHT HOLDER, WHITE META...
1  HAND WARMER UNION JACK, HAND WARMER RED POLKA DOT
2  ASSORTED COLOUR BIRD ORNAMENT, POPPY'S PLAYHOU...
3  JAM MAKING SET WITH JARS, RED COAT RACK PARIS ...
4                         BATH BUILDING BLOCK WORD
```

Out[6]: (18192, 1)