

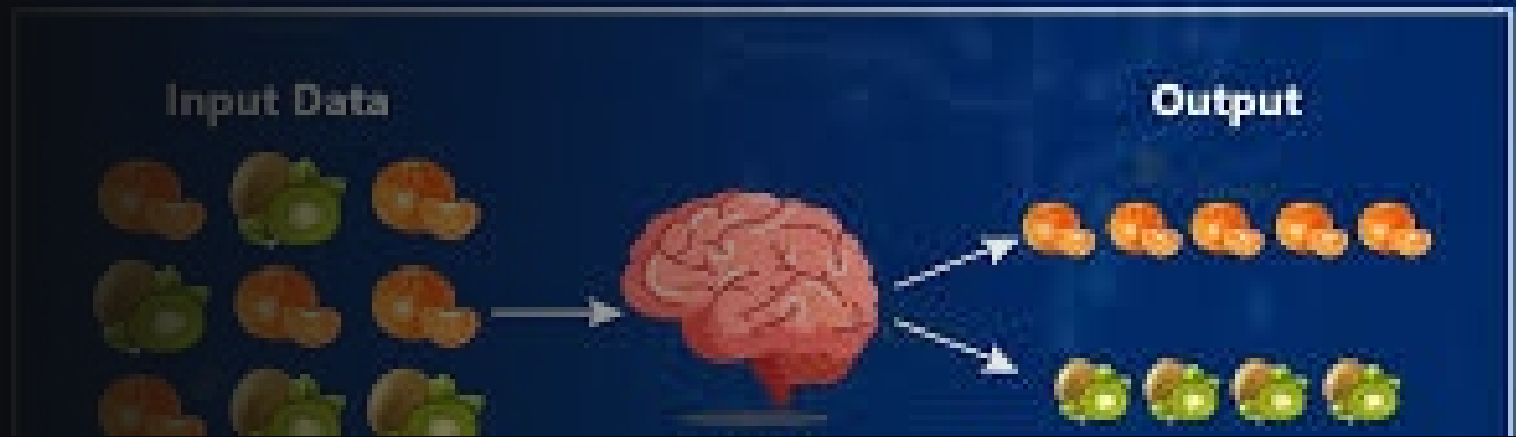
Unsupervised Learning Project

HOUSE



Unsupervised Machine Learning

By:
Rana Al-suffi





The dataset that I used was
wholesale_data.csv.



I used Jupyter notebook



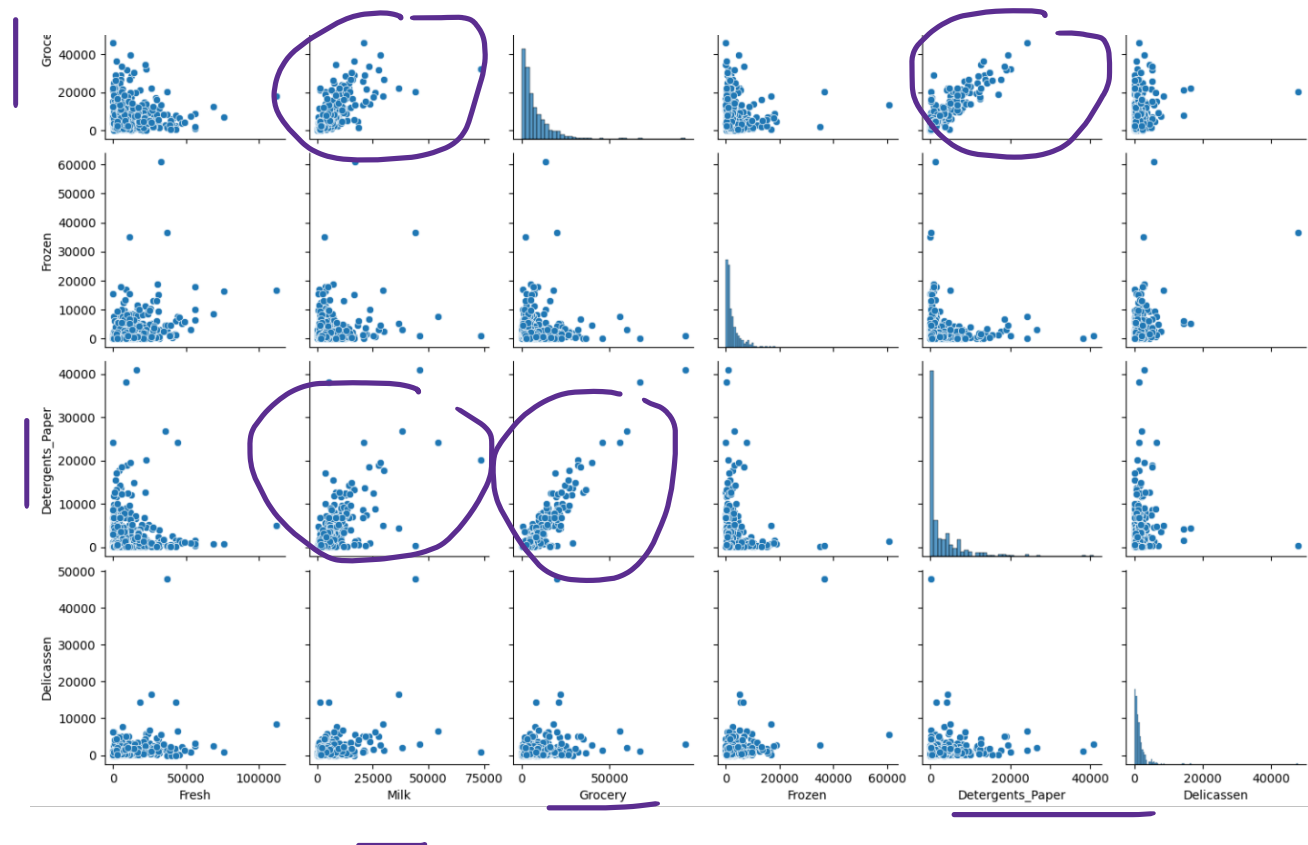
Libraries include pandas, NumPy,
matplotlib, scipy, sklearn, and
seaborn

For this
project:

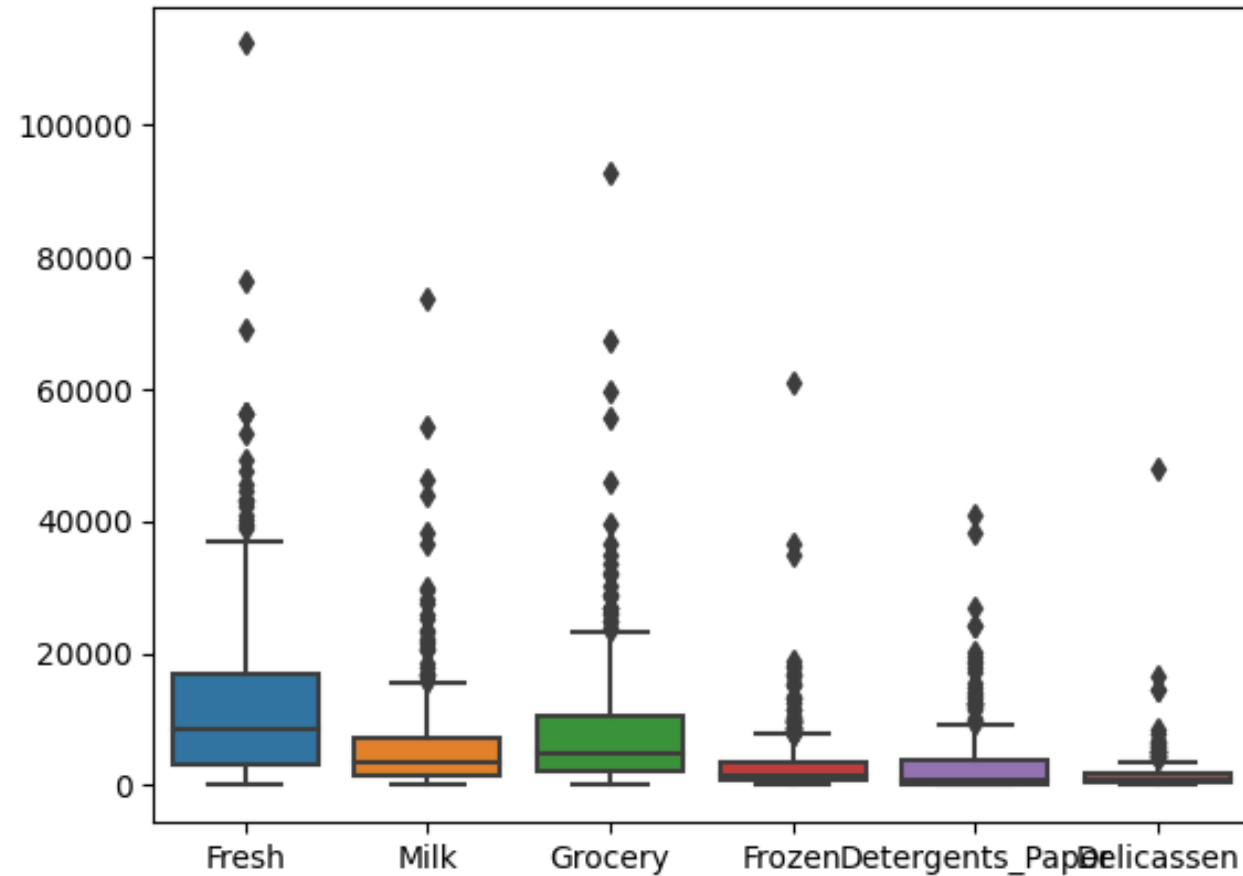


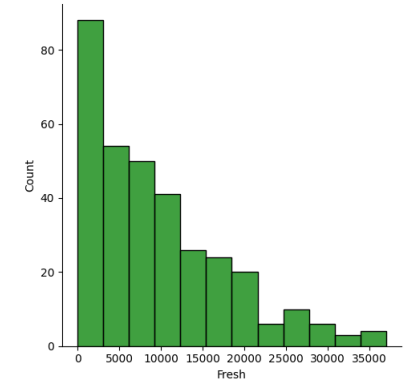
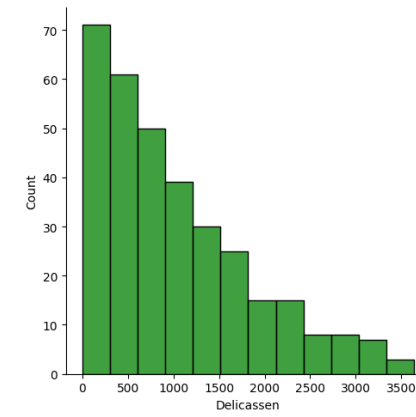
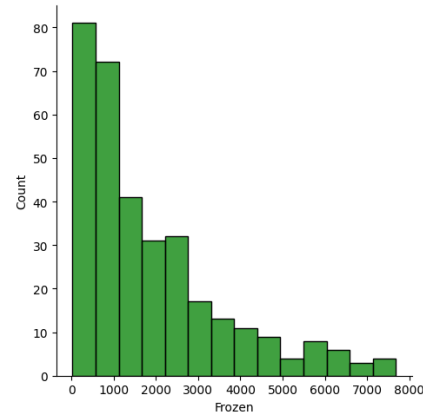
EDA - Exploratory Data Analysis & Pre-processing

- It looks as if all columns has a positive correlation with each other.
- However, the relationship between Detergents, Milk and grocery are higher correlation than other columns.
- Customers intend to purchase these products more frequently than others.

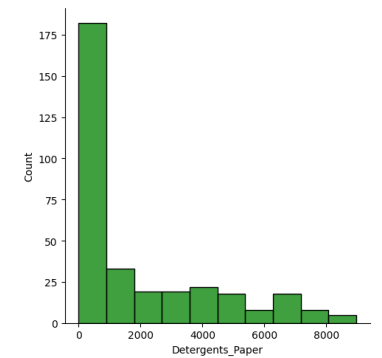
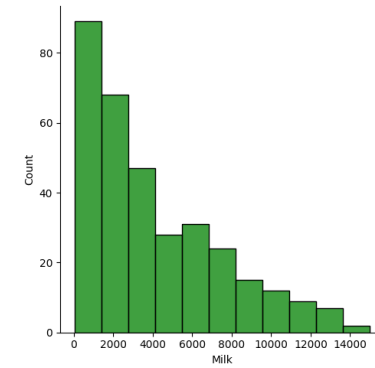
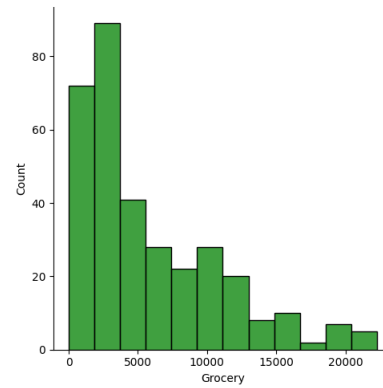


- From the data, it is obvious that there are outliers and in multiple columns.
- I used boxplot to take another look.
- And it looks like there are too many outliers.



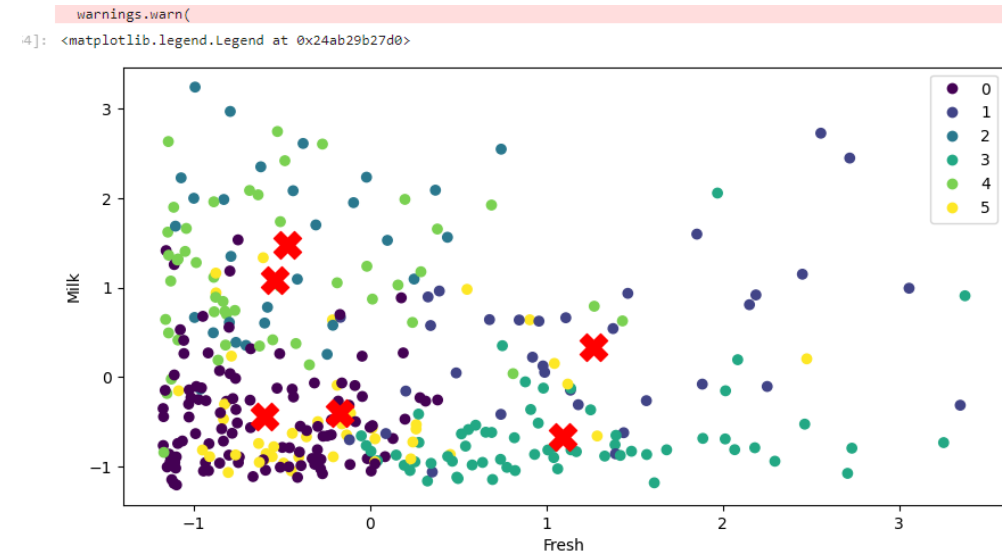


- I used IQR to remove the outliers from the dataset
- I checked the data after removing the outliers and it looked better.



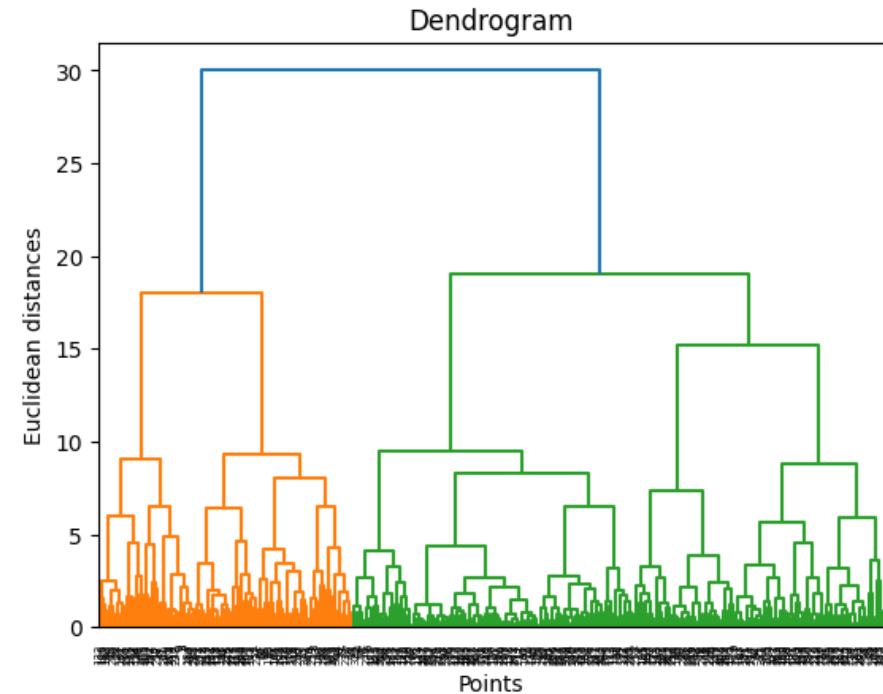
KMeans Clustering:

- Using KMeans, I started with three clustering groups however using the inertia score showed that it is best to make the clustering into 6 groups.
- To confirm the findings, I used the silhouette score which gave the same results.
- From the result, it means that all the shoppers can be divided into six groups and that depends on the type of items they purchase for each shopping trip.



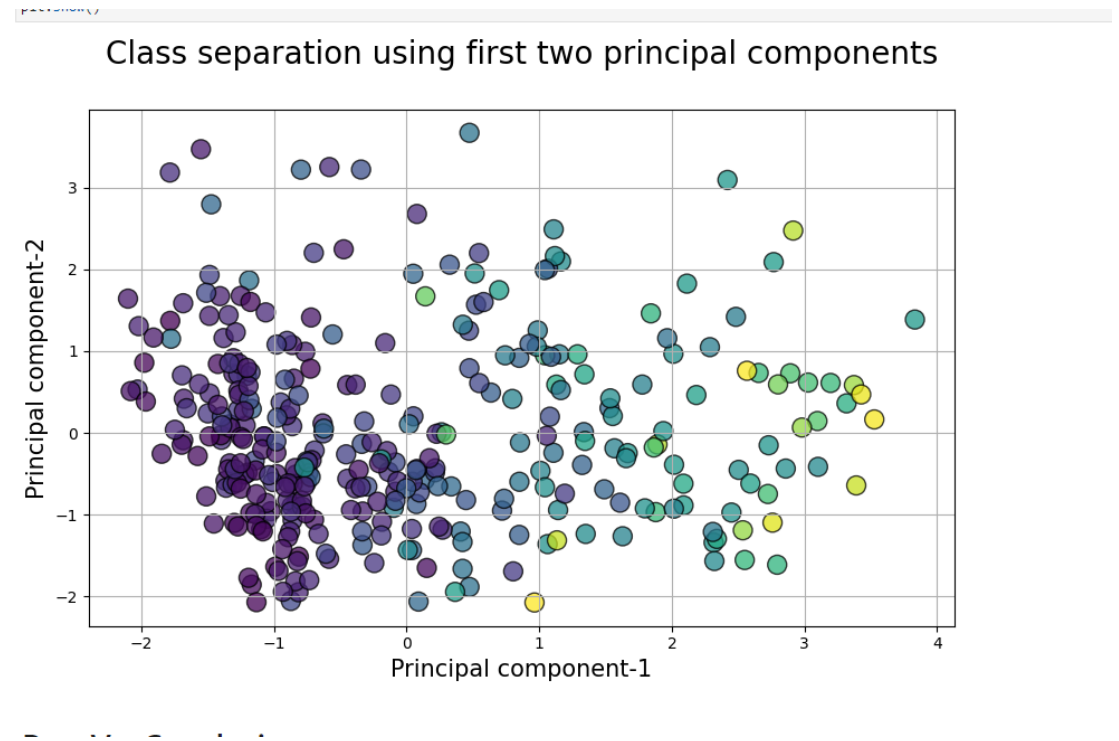
Hierarchical Clustering:

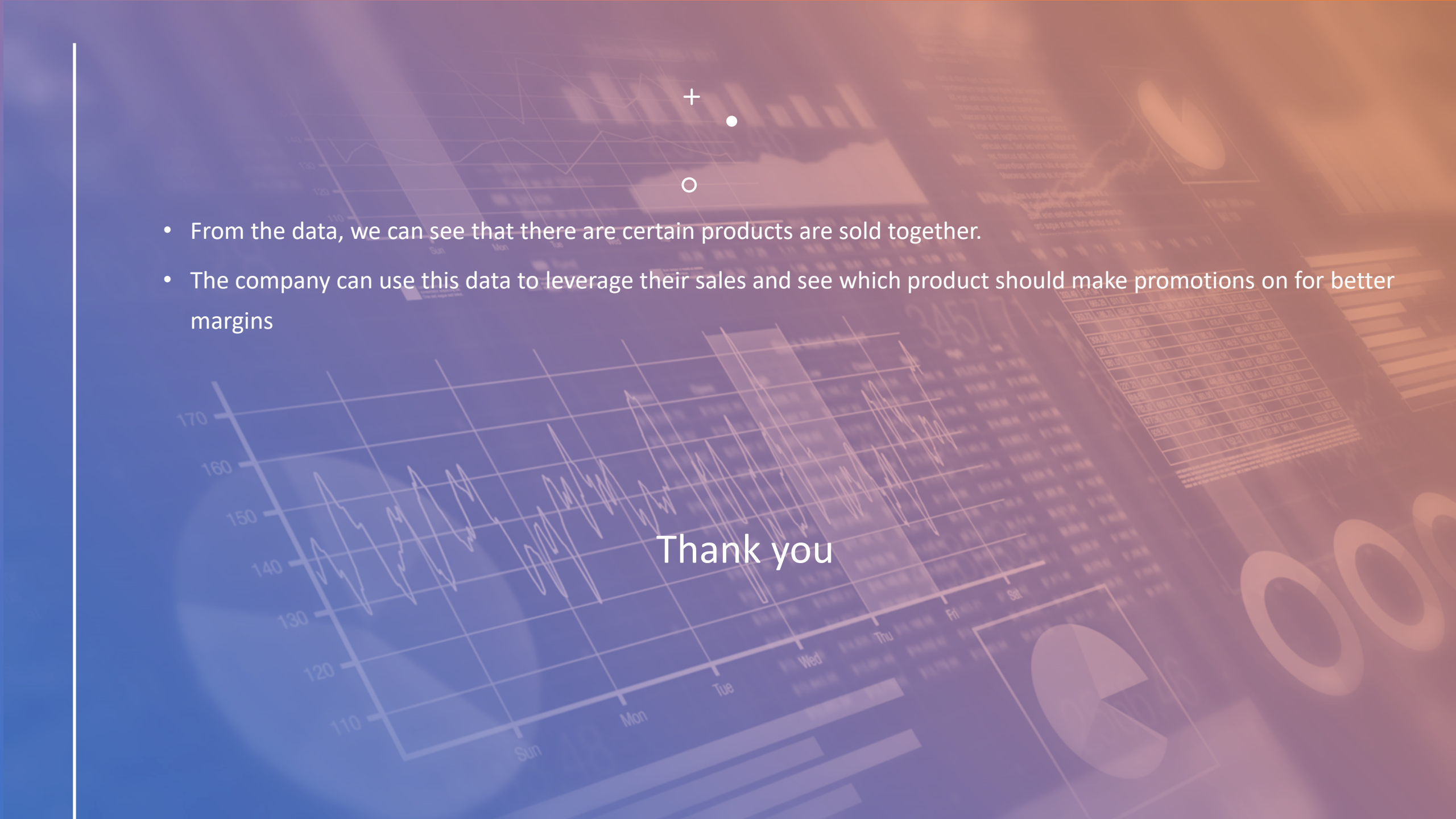
- After multiple tries, using different numbers of clusters, Dendrogram showed the best number of clusters is three



PCA:

- I used five products, which are 'Fresh', 'Milk', 'Frozen', 'Detergents_Paper', and 'Delicassen'. From processing the data, it is apparent that the first two components combined make 64% of the variance ratio.
- Using the previous information, we used only the first two components in our training.



- 
- The background features a collage of financial data visualizations. At the top, there is a line chart with a y-axis ranging from 110 to 140 and an x-axis with labels for Sun, Mon, Tue, and Wed. Below this, a larger line chart shows a y-axis from 110 to 170 and an x-axis with labels for Sun, Mon, Tue, Wed, Thu, Fri, and Sat. To the right of the charts, there is a table of numerical data. In the bottom right corner, a pie chart is visible. The entire background is overlaid with a semi-transparent orange and blue gradient.
- From the data, we can see that there are certain products are sold together.
 - The company can use this data to leverage their sales and see which product should make promotions on for better margins

Thank you