# Unsupervised Learning Wholesale data

Bernard Adeboye

# Project/Goals

- Apply unsupervised learning techniques to a real-world data set and use data visualization tools to communicate the insights gained from the analysis

- Perform exploratory data analysis,

- Preprocessing and feature engineering,
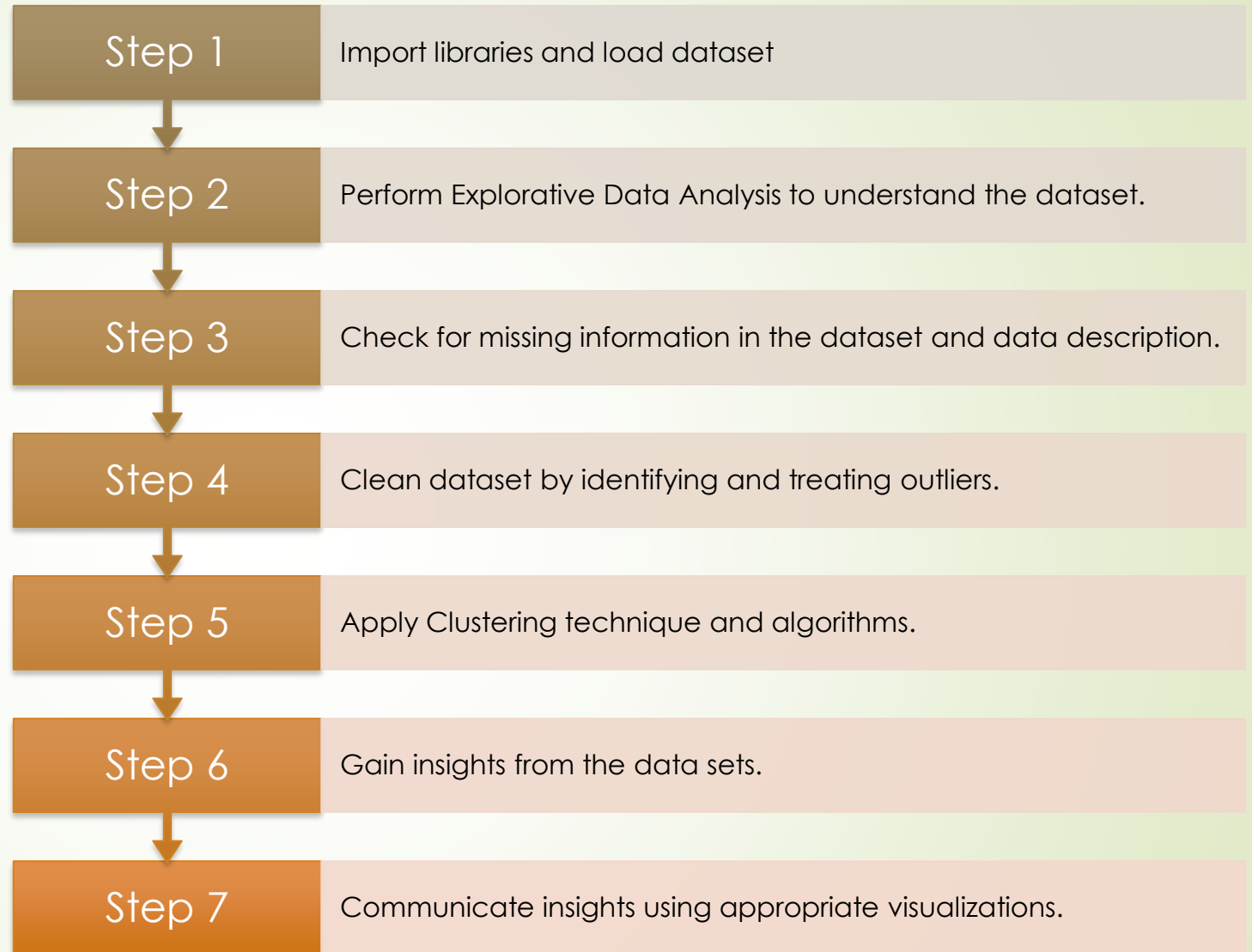
- Perform KMeans clustering,
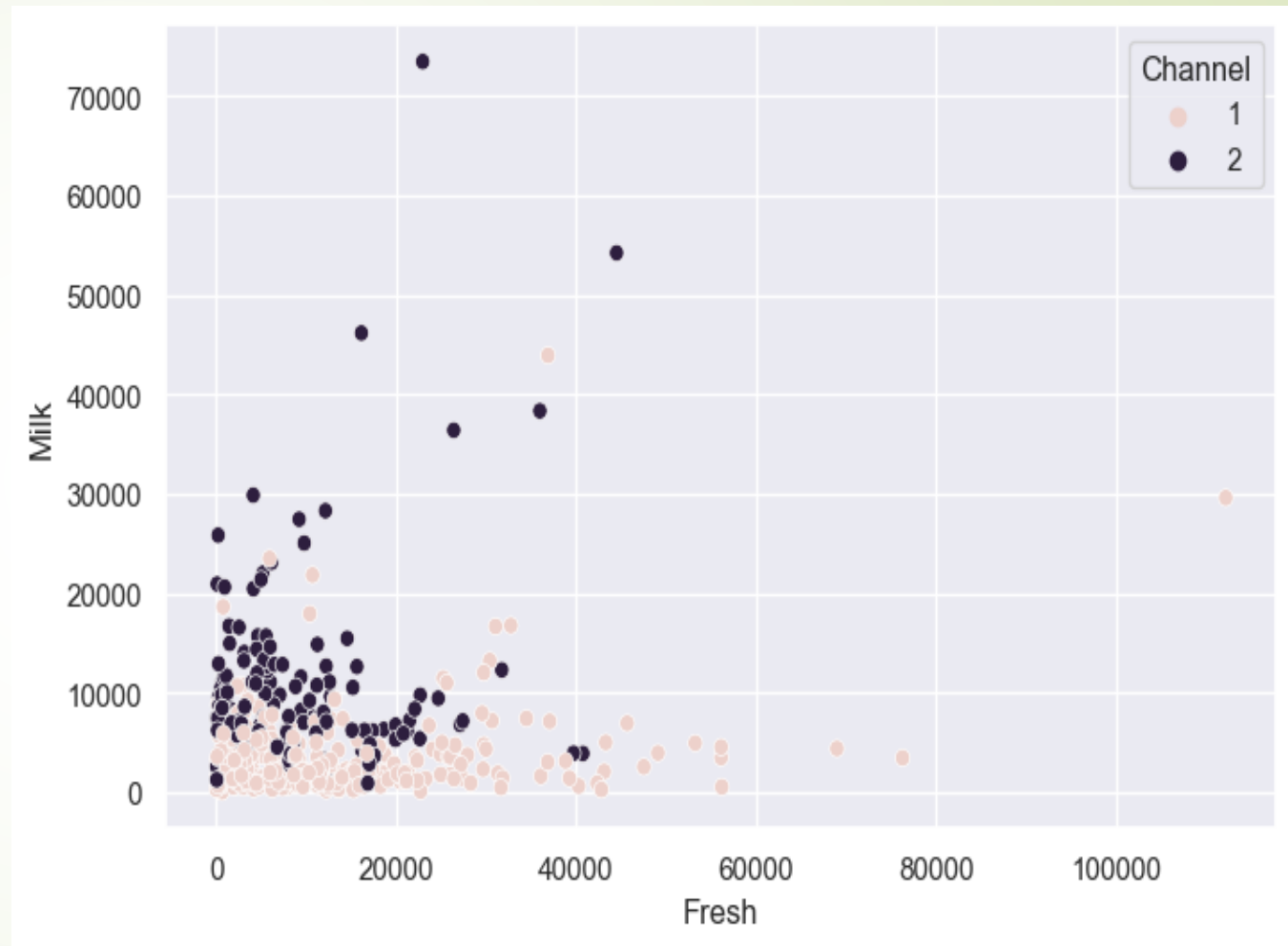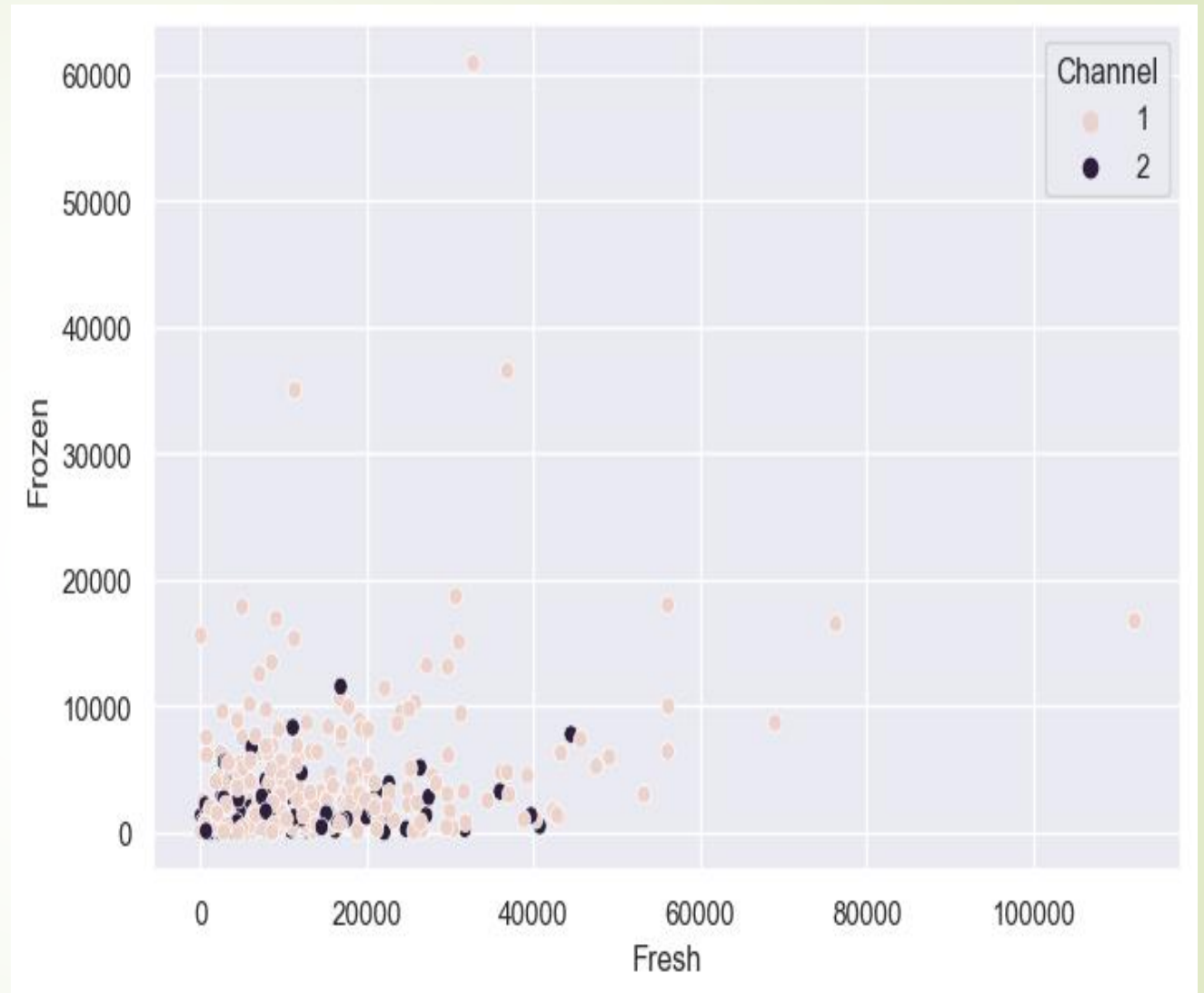
- Perform hierarchical clustering,

- Perform PCA

# Process

**Step 1** — Import libraries and load dataset

**Step 2** — Perform Explorative Data Analysis to understand the dataset.

**Step 3** — Check for missing information in the dataset and data description.

**Step 4** — Clean dataset by identifying and treating outliers.

**Step 5** — Apply Clustering technique and algorithms.

**Step 6** — Gain insights from the data sets.

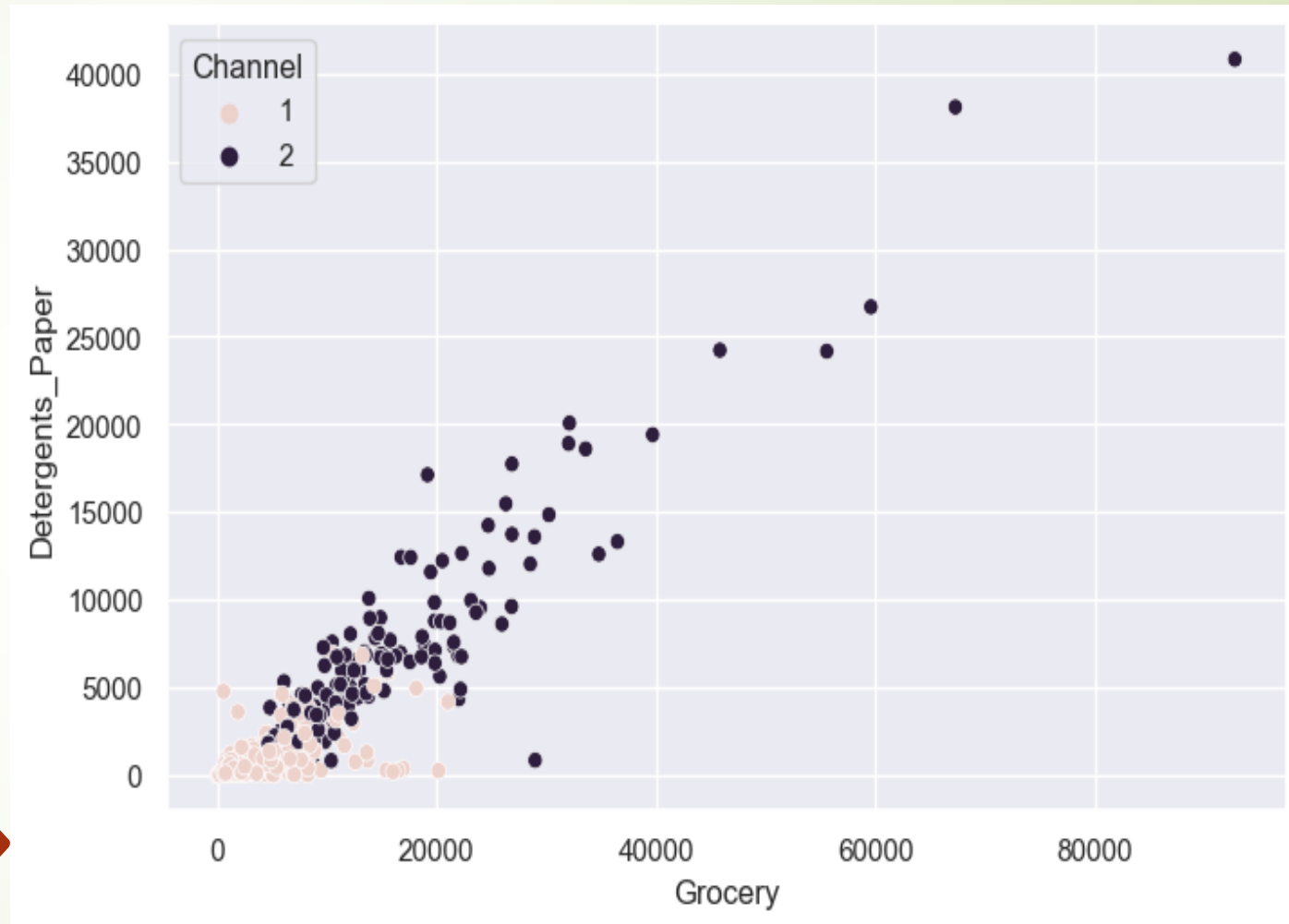**Step 7** — Communicate insights using appropriate visualizations.

# Explorative Data Analysis Relationship between the Fresh, Milk and Channel variables

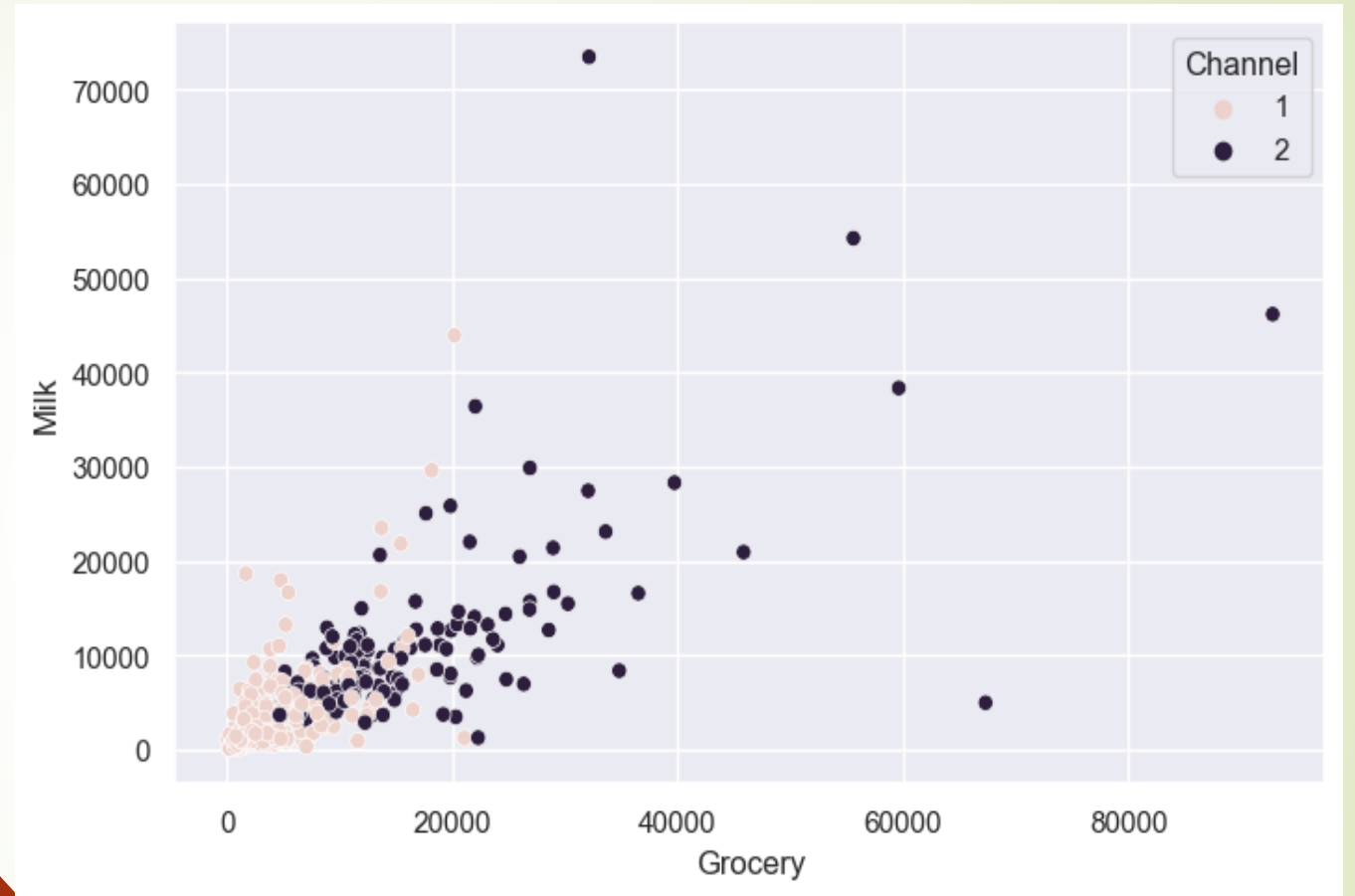# Relationship between the Fresh, Frozen and Channel variables

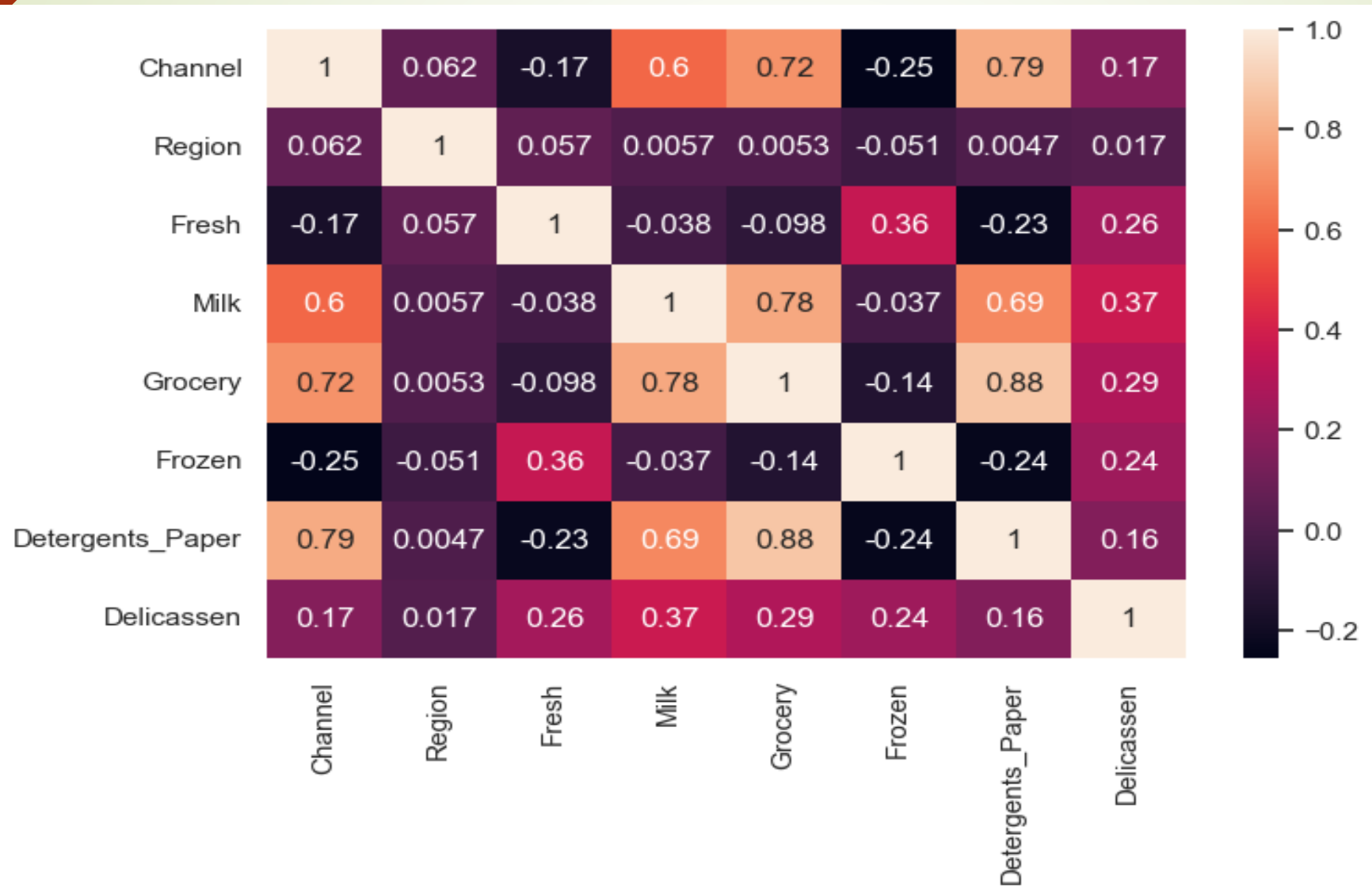# Relationship between the Grocery, Detergents_paper and Channel variables

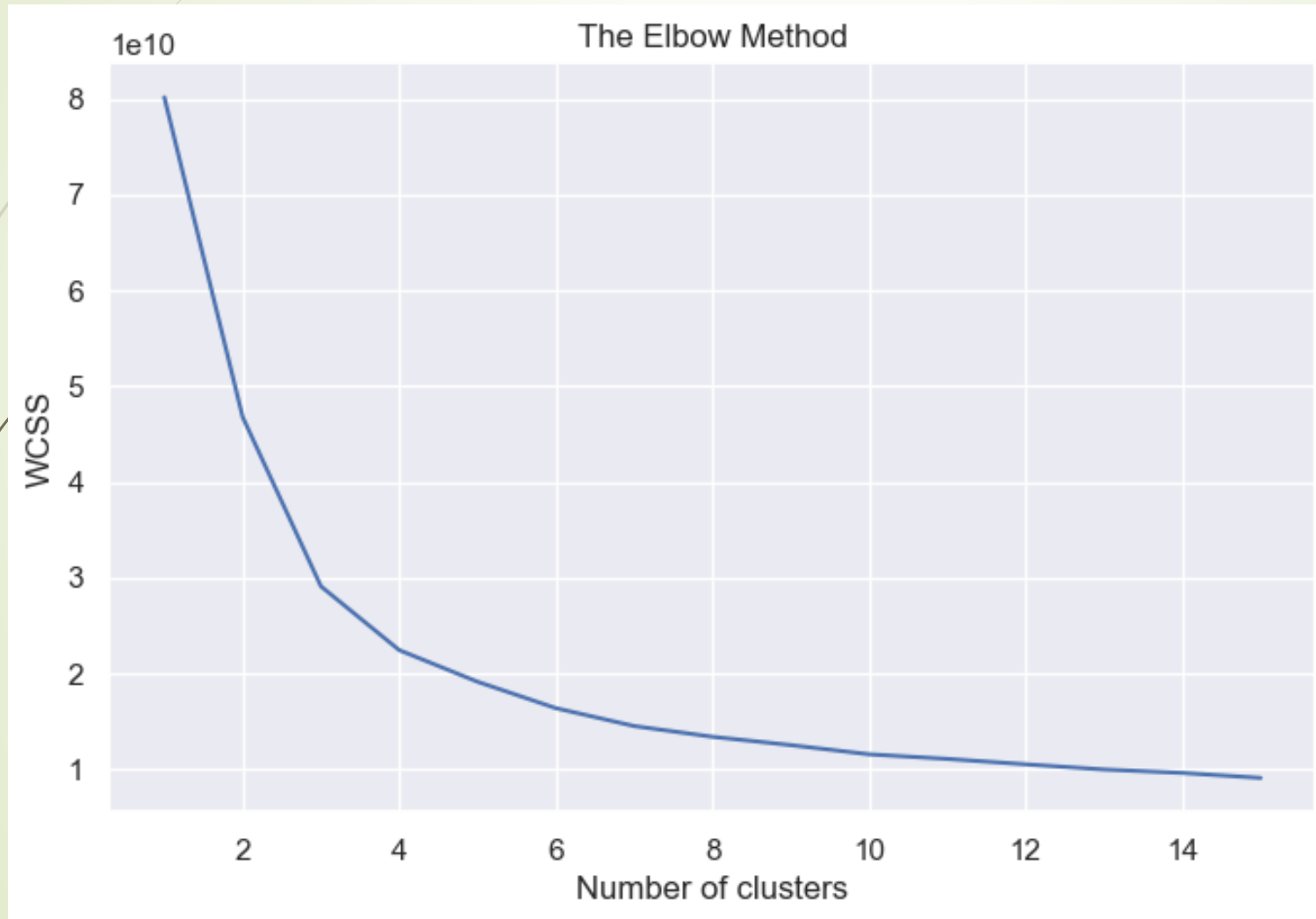# Relationship between the Grocery, Milk and region variables

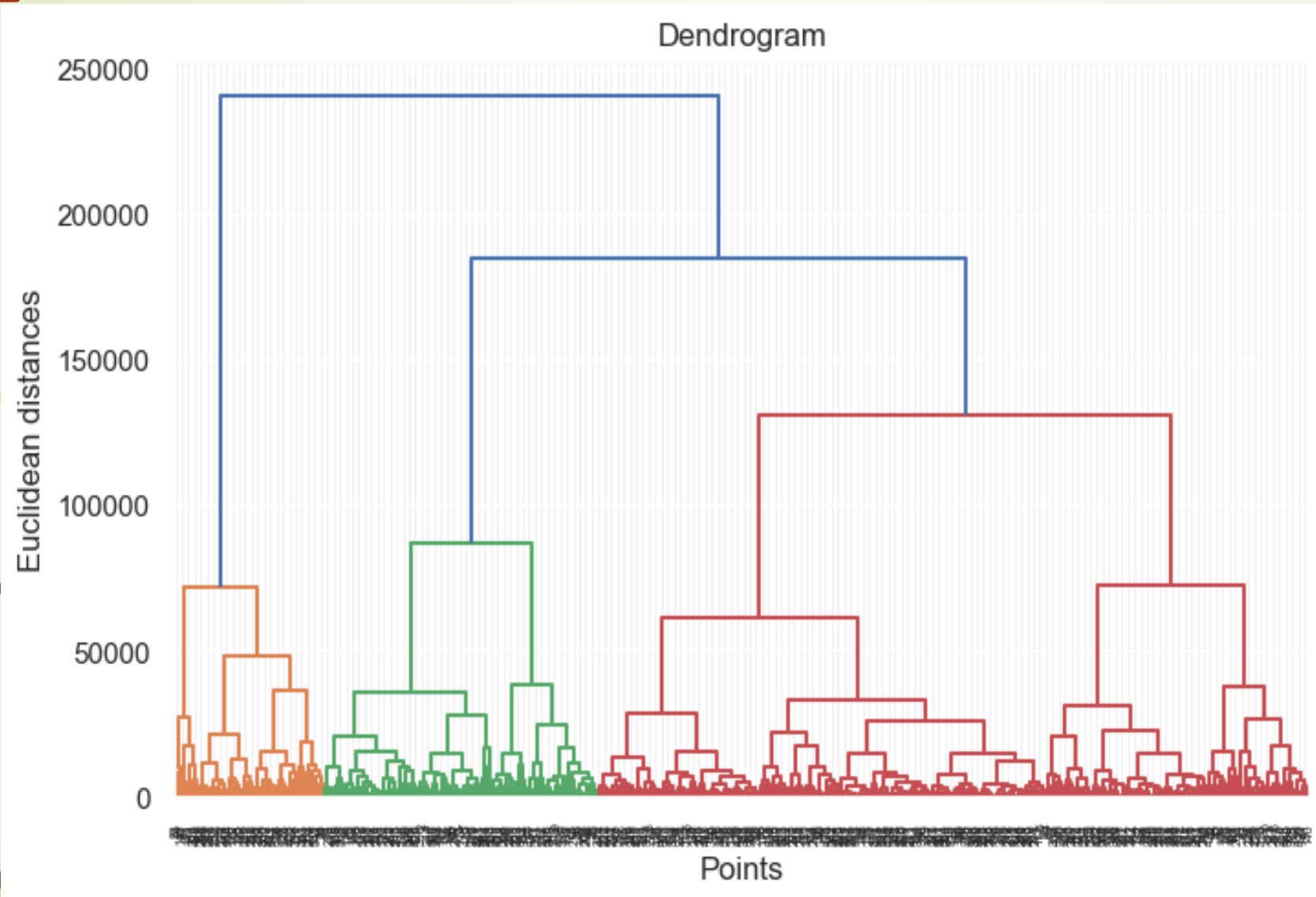# Relationship between the Grocery, Milk and Channel variables
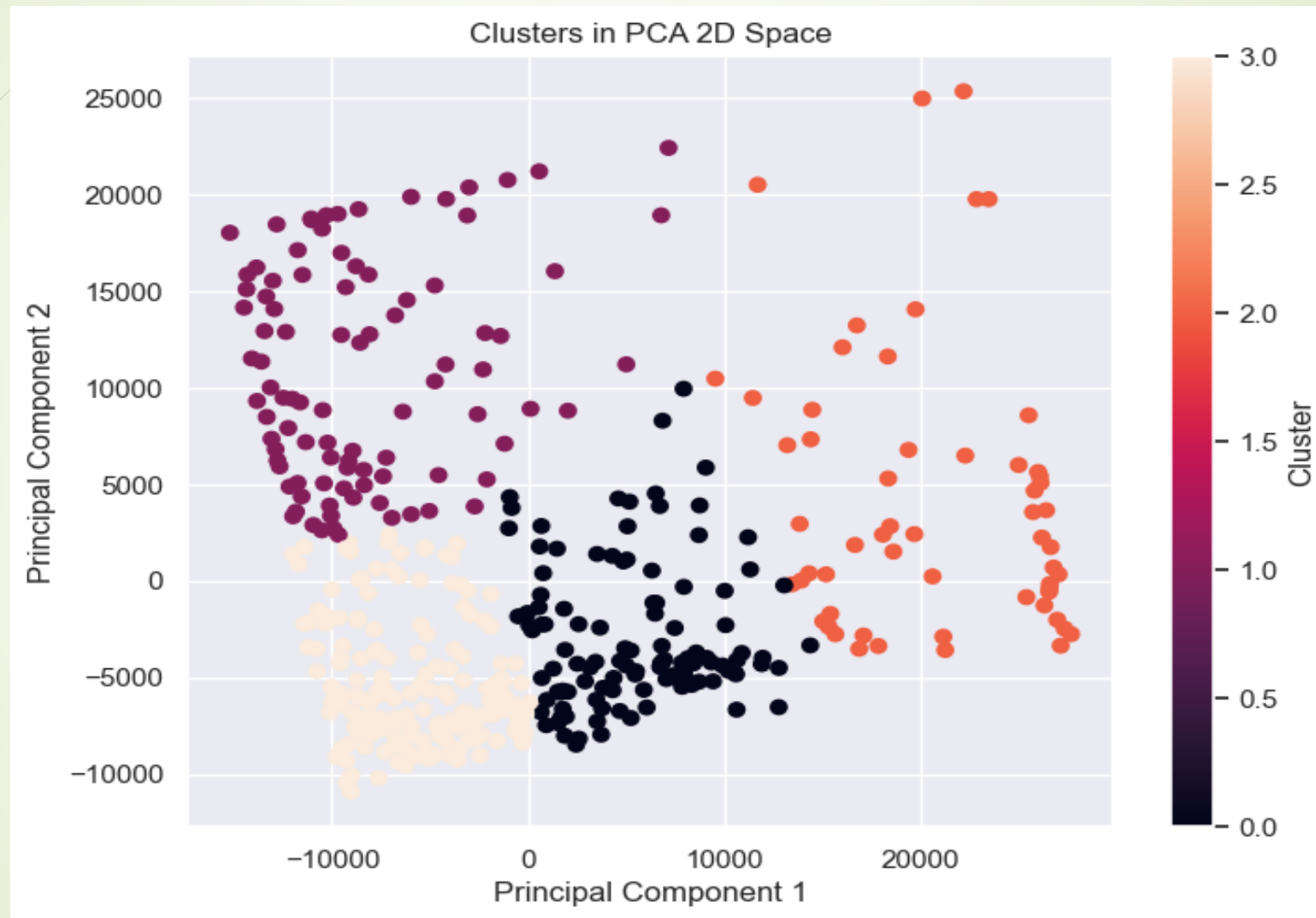
# Correlation heatmap

# The Elbow Method

# Hierarchical Clustering - Dendrogram

# Principal Component Analysis

# Conclusion

▶ Performing PCA on dataset for feature reduction yielded four clusters;

   1. Cluster 3 is the combination of products that falls around zero and negative. Stakeholders should consider promo or discount to encourage customers to engage in the purchase of those items.

   2. Cluster 2 is the combination of products that are of high interest to customers, therefore, stakeholder should provide     measure to retain this group and provide ways to increase their    purchasing activities.

▶ Performing K-means on the dataset grouped the products into 4 clusters.

▶ Right Skewness: Features such as 'Fresh', 'Milk', 'Grocery', 'Frozen', 'Detergents Paper', and 'Delicassen' show a right skew (mean > median).This could impact the performance of some machine learning algorithms.

▶ Combination of Milk and detergents paper/Grocery shows a positive linear relationship with channel 2 and region 3 playing a significant part.

▶ There is a very strong Relationship between the Grocery, Detergents_paper and Channel variables

▶ The Hierarchical Cluster shows two clusters based on agglomerative clustering.