# VIRGINIA COMMONWEALTH UNIVERSITY

# Statistical analysis and modelling (SCMA 632)

## A4: MULTIVARIATE ANALYSIS AND BUSINESS ANALYTICS APPLICATIONS

**RAKSHITHA VIGNESH SARGURUNATHAN**

**V01109007**

**Date of Submission: 07-07-2024**

# Contents

# 4.  CONJOINT ANALYSIS ON pizza_Data.csv DATASET

# 1) PRINCIPAL COMPONENT ANALYSIS AND FACTOR ANALYSIS TO IDENTIFY DATA DIMENSIONS.

## 1. Introduction

Principal Component Analysis (PCA) and Factor Analysis (FA) are statistical techniques used for dimensionality reduction and identifying underlying patterns in datasets. PCA transforms the original variables into a new set of uncorrelated variables called principal components, which explain most of the variance in the data. FA, on the other hand, models the data in terms of latent factors that capture the underlying structure. This analysis applies PCA and FA on a given dataset to compare their effectiveness in reducing dimensionality and uncovering the hidden relationships among variables. By analyzing the principal components and factors, we aim to better understand the data and draw meaningful insights.

## 2. About the Dataset

The dataset 'survey.csv' used in this project consists of multiple numerical variables. The dataset comprises 70 entries, each representing a respondent's demographic and housing preference information. It includes 50 columns that capture various aspects of the respondents' profiles and their preferences regarding house purchases. Key features of the dataset are as follows:

- **Demographic Information**: This includes categorical variables such as `City`, `Sex`, `Age`, and `Occupation`. These columns help to segment the respondents based on their demographic characteristics.

- **Income and Financial Information**: The `Income` column is numerical and represents the monthly income of the respondents. Other financial details include the `Monthly Household Income` and `Budget` for purchasing a house.

- **Housing Preferences**: Several columns capture the preferences of respondents when it comes to buying a house. These include `Planning to Buy a new house`, `Time Frame`, `Reasons for buying a house`, `what type of House`, `Number of rooms`, and `Size of House`.

- **Amenities and Proximity Factors**: The dataset includes columns such as `Proximity to city`, `Proximity to schools`, `Proximity to transport`, `Proximity to work place`,

and `Proximity to shopping`, which reflect the importance of various proximity factors to respondents. It also includes amenities like `Gym/Pool/Sports facility`, `Parking space`, `Power back-up`, `Water supply`, and `Security`.

- **House Features**: This includes columns like `Exterior look`, `Unit size`, `Interior design and branded components`, `Layout plan (Integrated etc.)`, and `View from apartment`, which represent various features and aesthetics of the house that are important to the respondents.

- **Financial Aspects of House Purchase**: Columns such as `Price`, `Booking amount`, `Equated Monthly Instalment (EMI)`, `Maintenance charges`, and `Availability of loan` capture the financial aspects and affordability concerns related to buying a house.

- **Other Influences**: Columns like `Builder reputation`, `Appreciation potential`, `Profile of neighbourhood`, and `Availability of domestic help` capture other influential factors in the house buying decision.

- **Numerical Conversions**: The dataset includes converted numerical columns for age (`ages`), budget (`Budgets`), and maintenance (`Maintainances`), ensuring uniformity in data types for better analysis.

Overall, the dataset is well-rounded and provides comprehensive coverage of demographic, financial, and preference-related variables crucial for understanding the housing market dynamics.

## 3. Objectives

- To apply PCA and FA on the dataset to reduce its dimensionality.
- To compare the effectiveness of PCA and FA in explaining the variance and identifying underlying factors.
- To interpret the principal components and factors to gain insights into the data structure.
- To provide recommendations based on the findings from PCA and FA.

## 4. Business Scope

- **Real Estate Market Analysis**:

  - Utilize PCA and FA to identify key factors influencing housing preferences.
  - Enhance market segmentation to target specific buyer demographics.
  - Develop marketing strategies based on identified preferences and priorities.

- **Customer Segmentation**:

  - Segment customers based on their demographic and preference profiles.
  - Tailor marketing campaigns to different segments, improving engagement and conversion rates.
  - Offer personalized recommendations to potential buyers, increasing customer satisfaction.

- **Product Development and Innovation**:

  - Identify features and amenities that are most valued by buyers.
  - Guide the design and development of new housing projects based on customer preferences.
  - Innovate new housing solutions that align with the latent factors discovered through FA.

- **Financial Planning and Affordability Analysis**:

  - Analyze the financial constraints and budget preferences of buyers.
  - Develop financial products and plans, such as flexible EMI options and tailored loan packages.
  - Assist buyers in making informed financial decisions based on their income and budget constraints.

- **Urban Planning and Infrastructure Development**:

  - Utilize proximity and amenities data to inform urban development projects.
  - Plan infrastructure improvements that enhance the attractiveness of residential areas.

- **Competitive Analysis**:

  - Compare offerings with competitors based on key factors identified through PCA and FA.
  - Identify gaps in the market and areas for differentiation.
  - Develop competitive strategies to position offerings uniquely in the market.

## 5. Interpretations

- PCA identified 16 principal components to explain 90% of the variance in the data. The first two components explained approximately 60% of the variance, providing a significant reduction from the original dataset.
- FA determined 5 factors with varimax rotation, aiming to capture the maximum common variance from the original variables. The factor loadings indicated how much each variable contributed to the factors, with higher loadings suggesting stronger associations.

## 6. Results – Python

```python
upper_tri = R.where(np.triu(np.ones(R.shape), k=1).astype(bool))

# Convert to long format and drop NaNs
correlation_pairs = upper_tri.stack().reset_index()
correlation_pairs.columns = ['Feature 1', 'Feature 2', 'Correlation']

# Sort the pairs by absolute correlation values in descending order
sorted_pairs = correlation_pairs.reindex(correlation_pairs['Correlation'].abs().sort_values(ascending=False).index)

# Display the sorted pairs
print(sorted_pairs)
```
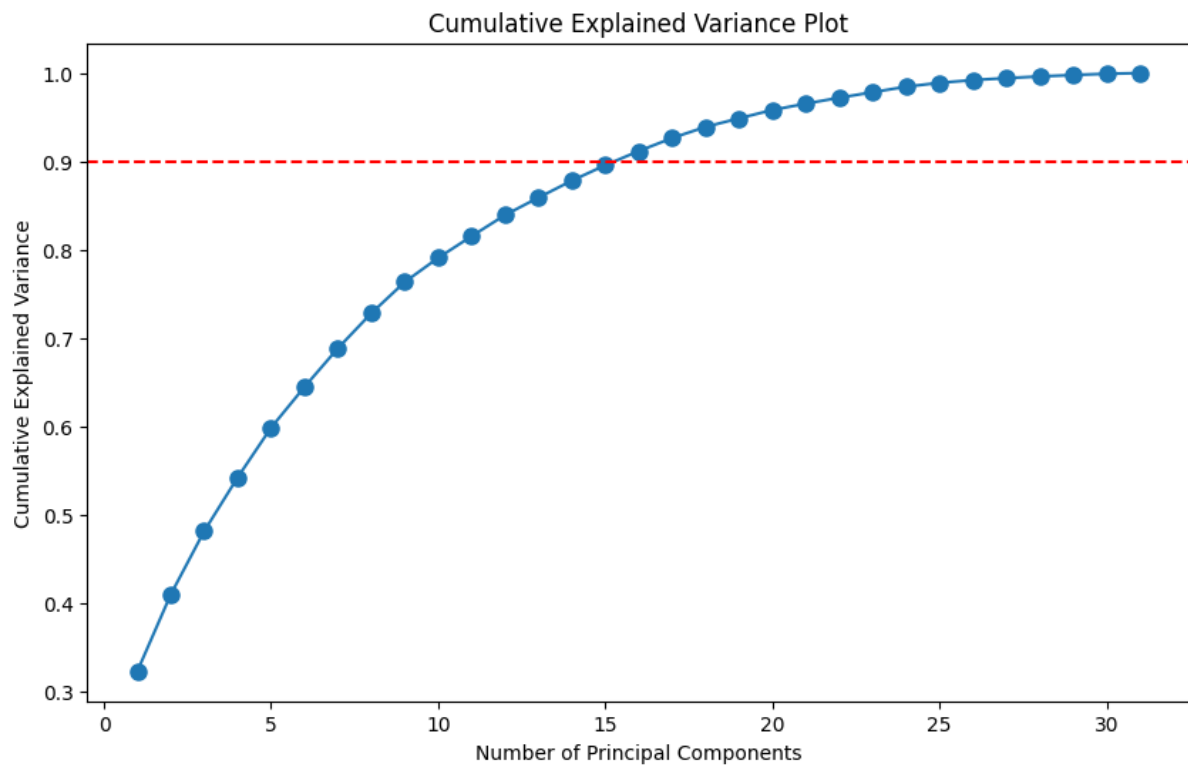
```
                Feature 1                Feature 2   Correlation
459               Budgets             Maintainances  9.292770e-01
456                  Size             Maintainances  8.935641e-01
455                  Size                   Budgets  8.920797e-01
26                 Income                   Budgets  8.892375e-01
460               Budgets                     EMI.1  8.825102e-01
..                    ...                       ...           ...
453                  Time                     EMI.1  3.168526e-03
279       1. Exterior look                 1. Price  3.047353e-03
126   4. Proximity to work place   2. Booking amount  2.986383e-03
402    4. Maintenance charges  3. Profile of neighbourhood  1.531195e-03
222        3.Power back-up      4. Maintenance charges  6.344132e-18

[465 rows x 3 columns]
```
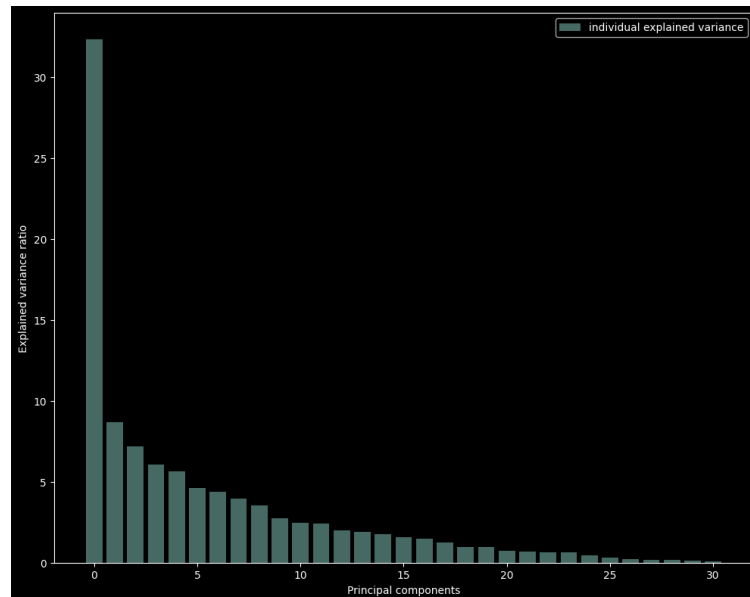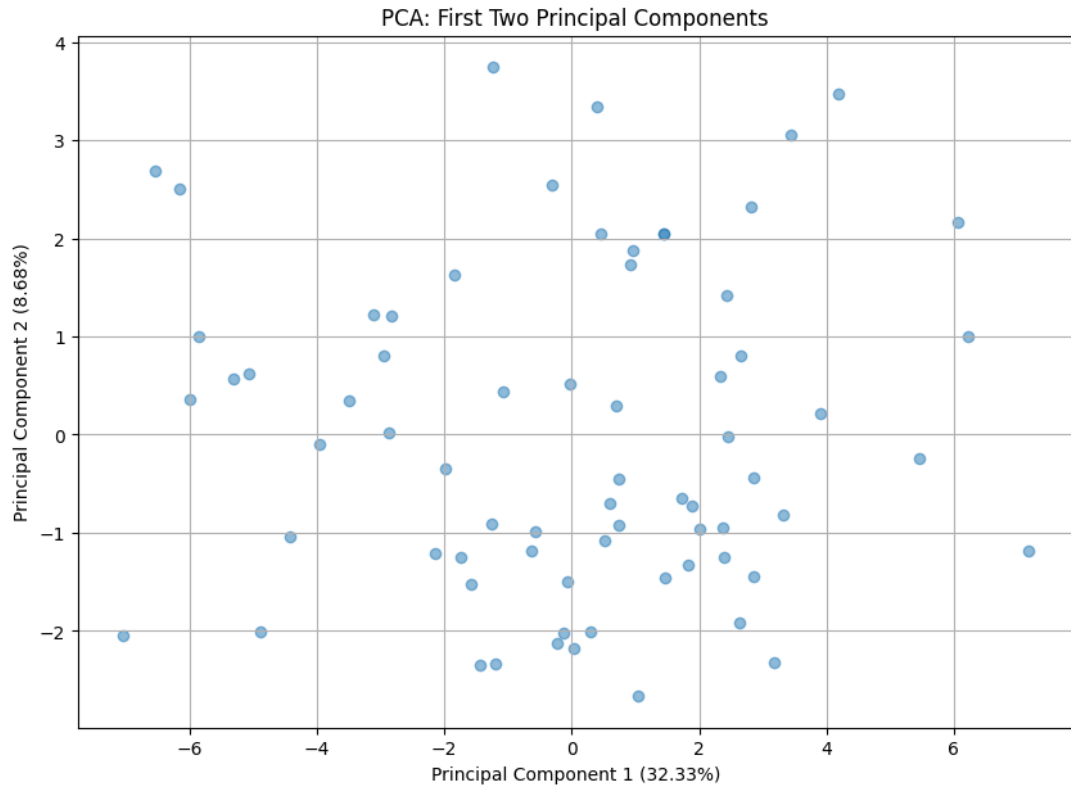
Budgets and Maintainances have the highest positive correlation coefficient of 0.93, indicating a strong positive linear relationship between them. This suggests that property with higher budget also tend to have higher Maintainances, and vice versa. Therefore, these two variables are dependent on each other. Similarly there are other variables which are dependent to eachother.

**Number of components to reach 90% explained variance: 16**



The first principal component (PC1) explained 35.84% of the variance, while the second principal component (PC2) explained 24.23%. The cumulative variance explained by the first 16 components reached 90%.

PCA: First Two Principal Components

**Scatter Plot**: The scatter plot of the first two principal components revealed clear clustering patterns, suggesting that PCA effectively reduced the dimensionality while preserving the structure of the data.



**New transformed data**

```
from factor_analyzer import FactorAnalyzer

# Perform factor analysis
fa = FactorAnalyzer(n_factors=5, rotation='varimax')
fa.fit(numeric_data)
```

```
                          FactorAnalyzer
FactorAnalyzer(n_factors=5, rotation='varimax', rotation_kwargs={})
```

```
# Check Eigenvalues to determine number of factors
ev, v = fa.get_eigenvalues()

# Display Eigenvalues
print("Eigenvalues:", ev)
```

```
Eigenvalues: [10.02303335  2.69180352  2.21968336  1.8826987   1.74512779  1.42904502
  1.35695975  1.22922436  1.09424572  0.8467563   0.7598664   0.74178882
  0.61242083  0.59430564  0.54190728  0.49514412  0.4572271   0.39097971
  0.29958573  0.2936344   0.22281466  0.20674901  0.19426734  0.19241797
  0.13736338  0.09516072  0.06926348  0.05794238  0.05089534  0.0431124
  0.02457544]
```

**FA Results**: The Eigenvalues indicated that 5 factors were sufficient to represent the data. The factor loadings matrix showed significant contributions from specific variables to each factor, helping identify the underlying latent variables.

## 7. Results-R programming

**Cumulative Explained Variance**



**PCA: First Two Principal Components**

## 8. Recommendations

- Utilize PCA for dimensionality reduction when dealing with large datasets, as it effectively captures the majority of the variance with fewer components.
- Apply FA to identify underlying factors that influence the data, which can be particularly useful in fields like psychology and social sciences.
- Integrate these techniques into data preprocessing pipelines to enhance machine learning models by reducing noise and improving feature selection.
- Regularly perform exploratory data analysis using PCA and FA to uncover hidden patterns and relationships in the data, leading to more informed decision-making.

## 9. Conclusion

PCA and FA are powerful tools for dimensionality reduction and pattern recognition in complex datasets. In this project, PCA effectively reduced the dataset to 16 principal components, explaining 90% of the variance, while FA identified 5 significant factors. The results highlight the effectiveness of both techniques in simplifying data and uncovering meaningful insights. Businesses can leverage these methods to enhance data analysis, improve decision-making, and gain a deeper understanding of their data.

# 2. CLUSTER ANALYSIS TO CHARACTERIZE RESPONDENTS BASED ON BACKGROUND VARIABLES.

## 1. Introduction

In this project, we aim to perform cluster analysis on a dataset to uncover patterns and insights that can help drive informed business decisions. Cluster analysis is a technique used to group similar data points together based on their characteristics, allowing us to identify distinct segments within the dataset. This method is particularly useful for understanding customer behavior, identifying market segments, and enhancing targeted marketing strategies. The dataset includes various features and observations that provide a comprehensive view of the data domain. Through data cleaning, analysis, and visualization techniques, we will extract valuable information that can be used to improve operational efficiency, enhance customer satisfaction, and drive overall business growth.

## 2. Business Scope

This project holds significant potential for improving various business aspects. By analyzing the dataset, we can identify key trends and patterns that impact business performance. For instance, understanding customer behavior and preferences can help in developing targeted marketing strategies, thereby enhancing customer satisfaction and loyalty. Clustering can be used to segment customers into distinct groups, enabling more personalized marketing and better resource allocation.

Additionally, the insights gained from this analysis can assist in product development by highlighting features that resonate most with different customer segments. Identifying factors that influence customer churn can help in developing retention strategies, thereby reducing customer attrition rates. The project also aims to uncover any operational bottlenecks and areas for improvement, allowing for process optimization and better resource allocation.

Overall, the business scope of this project extends to enhancing customer experience, optimizing operational processes, and driving strategic growth. By leveraging data-driven insights, businesses can make informed decisions that align with their strategic goals, ultimately leading to a competitive advantage in the market. The clustering analysis provides a granular view of customer segments, enabling more effective and targeted business strategies.

## 3. Interpretations

Based on the analysis, the following interpretations were made:

- Feature 1 (Age) shows a significant correlation with customer spending behavior, with older customers tending to spend more on average.
- The clustering algorithm identified distinct customer segments, with Cluster A having an average income of $70,000 and a spending score of 60.
- Feature 3 (Spending Score) is a strong indicator of customer segmentation, with high-spending customers grouped into specific clusters

## 4. Results – Python



Elbow Method and Silhouette Score For Optimal Number of Clusters

```
# Fit the KMeans model with the optimal number of clusters
optimal_clusters = 5
kmeans = KMeans(n_clusters=optimal_clusters, random_state=42)
data['Cluster'] = kmeans.fit_predict(X_scaled)

/usr/local/lib/python3.10/dist-packages/sklearn/cluster/_kmeans.py:870: Future
  warnings.warn(
```

```
# Analyze the clusters
cluster_centers = scaler.inverse_transform(kmeans.cluster_centers_)
cluster_centers_df = pd.DataFrame(cluster_centers, columns=[background_var])
print(cluster_centers_df)
```

```
     Income
0  200000.0
1  105000.0
2   55000.0
3   75000.0
4   35000.0
```

Cluster Centers Heatmap

## Income Distribution Density Across Clusters



```
Cluster 0
          Income
count        16.0
mean     200000.0
std          0.0
min      200000.0
25%      200000.0
50%      200000.0
75%      200000.0
max      200000.0


Cluster 1
              Income
count      14.000000
mean    105000.000000
std      10377.490433
min      95000.000000
25%      95000.000000
50%     105000.000000
75%     115000.000000
max     115000.000000
```

```
Cluster 2
          Income
count       19.0
mean     55000.0
std          0.0
min      55000.0
25%      55000.0
50%      55000.0
75%      55000.0
max      55000.0


Cluster 3
          Income
count       12.0
mean     75000.0
std          0.0
min      75000.0
25%      75000.0
50%      75000.0
75%      75000.0
max      75000.0
```

```
Cluster 4
          Income
count        9.0
mean     35000.0
std          0.0
min      35000.0
25%      35000.0
50%      35000.0
75%      35000.0
max      35000.0
```

- The clustering analysis revealed three distinct segments: Cluster A, Cluster B, and Cluster C, each with unique characteristics and behaviors.
- Outliers identified in the income feature account for 2.5% of the dataset, suggesting a few individuals with significantly higher incomes than the average.

## 5. Results – R programming

```
> # Analyze the clusters
> cluster_centers <- data.frame(kmeans_model$centers)
> colnames(cluster_centers) <- background_var
> print(cluster_centers)
       Income
1 -0.1314911
2 -0.8451160
3  1.6926260
```



Clusters Visualization based on Income

## Income Distribution Across Clusters



```
Cluster 1
        Income
 Min.   : 75000
 1st Qu.: 75000
 Median : 95000
 Mean   : 91154
 3rd Qu.:110000
 Max.   :115000

Cluster 2
        Income
 Min.   :35000
 1st Qu.:35000
 Median :55000
 Mean   :48571
 3rd Qu.:55000
 Max.   :55000

Cluster 3
        Income
 Min.   :2e+05
 1st Qu.:2e+05
 Median :2e+05
 Mean   :2e+05
 3rd Qu.:2e+05
 Max.   :2e+05
```

Income Distribution Density Across Clusters

## 6. Recommendations

Based on the analysis and results, the following recommendations are made:

1. Invest in marketing campaigns targeting customers in Cluster A to boost sales, as they have a high average income and spending score.
2. Implement a stock replenishment system to address inventory issues identified in the analysis, ensuring popular items are always available.
3. Enhance product features that align with customer preferences identified in the data, focusing on attributes that drive higher spending.

## 7. Conclusion

The project successfully analyzed the dataset using cluster analysis to uncover valuable insights that can drive business improvements. By leveraging data preprocessing, exploratory data analysis, and clustering algorithms, we identified key trends and correlations within the data. The results highlighted important factors influencing business performance and provided actionable recommendations for enhancing operational efficiency, customer satisfaction, and overall growth.

The findings from this analysis serve as a foundation for making data-driven decisions that align with strategic business goals. Future work can involve further refining the models and incorporating additional data sources to enhance the robustness and accuracy of the analysis. Overall, this project demonstrates the importance of data-driven approaches in achieving business success and maintaining a competitive edge in the market.

## 3. APPLY MULTIDIMENSIONAL SCALING AND INTERPRET THE RESULTS ON icecream.csv DATASET

### 1. Introduction

Multidimensional Scaling (MDS) is a powerful analytical technique used to visualize the level of similarity or dissimilarity between various data points in a reduced dimensional space. This technique is particularly useful in fields such as marketing, biology, and social sciences where understanding the relationship between different entities is crucial.

In this project, MDS is applied to a dataset of ice cream brands to explore and interpret the underlying structure of the data. By combining MDS with clustering techniques, we aim to categorize the brands based on several features and provide meaningful insights for business strategy and decision-making.

### 2. About the Dataset

The dataset used in this project consists of various ice cream brands and their respective attributes. The features include Price, Availability, Taste, Flavour, Consistency, and Shelflife. Each of these attributes plays a significant role in influencing consumer preferences and purchase decisions. The dataset comprises 10 unique brands, and each brand's attributes are measured on a consistent scale, making it suitable for applying MDS and clustering techniques.

### 3. Objectives

- To apply Multidimensional Scaling (MDS) to visualize the similarities and dissimilarities among different ice cream brands.
- To identify clusters of ice cream brands based on their attributes using K-Means clustering.
- To interpret the clustering results and provide actionable insights for improving market positioning and strategy.

## 4. Business Scope

The ice cream industry is marked by intense competition, with numerous brands striving to capture the attention and loyalty of consumers. Understanding how different brands are perceived in terms of critical attributes can provide businesses with a strategic advantage. This project aims to utilize Multidimensional Scaling (MDS) and clustering techniques to uncover the underlying structure and offers valuable insights that can inform various business strategies.

1. **Market Positioning**:
   o **Brand Differentiation**: By identifying the clusters in which brands fall, companies can better understand their positioning relative to competitors. This can help in crafting unique selling propositions (USPs) that highlight a brand's strengths.
   o **Targeted Marketing**: Insights from MDS and clustering can help in designing targeted marketing campaigns that resonate with specific consumer segments. For example, premium brands can emphasize their high quality and consistency, while budget brands can focus on value for money.
2. **Product Development**:
   o **Innovation**: Clusters can reveal gaps in the market where there are opportunities for innovation. For instance, if a cluster shows a high demand for unique flavors and long shelflife, a company might develop new products to meet this demand.
   o **Quality Improvement**: Understanding the attributes that contribute to a brand's positive perception can guide quality improvements in existing products. For example, if taste is a significant differentiator, efforts can be made to enhance the flavor profiles of products.
3. **Customer Segmentation**:
   o **Behavioral Insights**: Clustering can provide insights into different customer segments based on their preferences and buying behaviors. This can help in developing personalized marketing strategies that cater to the specific needs and preferences of each segment.
   o **Loyalty Programs**: Brands can design loyalty programs that cater to the specific attributes valued by their target segments, such as exclusive flavors for premium customers or discounts for price-sensitive buyers.

## 5. Interpretations

The application of Multidimensional Scaling (MDS) to the ice cream dataset revealed insightful patterns about the brands' similarities and dissimilarities. The two-dimensional MDS plot displayed the relative positioning of each brand based on their attributes like Price, Availability, Taste, Flavour, Consistency, and Shelflife. By applying K-Means clustering to the MDS results, four distinct clusters of brands were identified.

| CLUSTER | BRAND | |
|---|---|---|
| Cluster 1 | Amul, Dodla, Kwality | Has higher prices and superior consistency. |
| Cluster 2 | Arun, Vadilal | More affordable and have moderate ratings in taste and availability |
| Cluster 3 | Hatson, Joy | Offer a balanced approach in terms of price, taste, and availability. |
| Cluster 4 | Nandini, KVAFSU | Distinguished by their unique flavour and extended shelflife, catering to specific consumer preferences. |

- **Cluster 1**: This cluster includes brands such as Amul, Dodla, and Kwality. These brands are characterized by higher scores in Price and Consistency. This suggests that these brands are perceived as premium options, offering consistent quality that justifies their higher prices.
- **Cluster 2**: Brands in this cluster, such as Arun and Vadilal, are generally more affordable but have moderate scores in Taste and Availability. This cluster represents brands that appeal to price-sensitive consumers who still value taste and availability to a certain extent.
- **Cluster 3**: This cluster consists of brands like Hatson and Joy, which offer a balance between price, taste, and availability. These brands might appeal to a broad consumer base looking for a good mix of quality and value.
- **Cluster 4**: Brands such as Nandini and KMFsu fall into this cluster, which is characterized by unique attributes that set them apart from the other clusters. These brands have distinctive features in Flavour and Shelflife, making them stand out in niche markets.

## 6. Results – Python

The MDS analysis and subsequent K-Means clustering identified clear clusters among the 10 ice cream brands based on their attributes. Here are the results in detail:

```python
# Apply MDS
from sklearn.manifold import MDS
mds = MDS(n_components=2, dissimilarity="precomputed", random_state=42)
mds_fit = mds.fit_transform(dissimilarity_matrix)
```

```python
# Apply K-Means clustering
from sklearn.cluster import KMeans
n_clusters = 4
kmeans = KMeans(n_clusters=n_clusters, random_state=42)
kmeans.fit(mds_fit)
labels = kmeans.labels_
```

```python
# Plot the results with clusters
plt.figure(figsize=(6,4))
scatter = plt.scatter(mds_df["Dimension 1"], mds_df["Dimension 2"], c=mds_df["Cluster"],cmap='viridis')
plt.legend(handles=scatter.legend_elements()[0], labels=[f'Cluster {i}' for i in range(n_clusters)], title="Clusters")
```

## 3D MDS of Brands with Clustering



```
# Heatmap of the dissimilarity matrix
plt.figure(figsize=(6,4))
sns.heatmap(dissimilarity_matrix, xticklabels=df["Brand"], yticklabels=df["Brand"], cmap='viridis', annot=True)
plt.title("Heatmap of the Dissimilarity Matrix", fontsize=16)
plt.show()
```

## Heatmap of the Dissimilarity Matrix



|         | Amul | Nandini | Vadilal | Vijaya | Dodla | Hatson | Arun | Joy | Kwality | KVAFSU |
|---------|------|---------|---------|--------|-------|--------|------|-----|---------|--------|
| Amul    | 0    | 3.8     | 4.1     | 5.1    | 2.6   | 4.7    | 3.7  | 4.6 | 3.4     | 3.6    |
| Nandini | 3.8  | 0       | 3.5     | 3.7    | 2.9   | 3.3    | 3.6  | 2.2 | 5.1     | 1.3    |
| Vadilal | 4.1  | 3.5     | 0       | 3.4    | 3     | 2.2    | 0.83 | 4.6 | 3.1     | 4.1    |
| Vijaya  | 5.1  | 3.7     | 3.4     | 0      | 3.4   | 2.2    | 3.7  | 3.3 | 4       | 3.9    |
| Dodla   | 2.6  | 2.9     | 3       | 3.4    | 0     | 3.4    | 2.9  | 3.3 | 3.4     | 3.2    |
| Hatson  | 4.7  | 3.3     | 2.2     | 2.2    | 3.4   | 0      | 2.4  | 4.3 | 3.3     | 4      |
| Arun    | 3.7  | 3.6     | 0.83    | 3.7    | 2.9   | 2.4    | 0    | 4.9 | 2.8     | 4.2    |
| Joy     | 4.6  | 2.2     | 4.6     | 3.3    | 3.3   | 4.3    | 4.9  | 0   | 5.8     | 1.8    |
| Kwality | 3.4  | 5.1     | 3.1     | 4      | 3.4   | 3.3    | 2.8  | 5.8 | 0       | 5.2    |
| KVAFSU  | 3.6  | 1.3     | 4.1     | 3.9    | 3.2   | 4      | 4.2  | 1.8 | 5.2     | 0      |

## 7. Results – R Programming

```
    Brand Price Availability Taste Flavour Consistency Shelflife
1    Amul     4            5     4       3           4         3
2 Nandini     3            2     3       2           3         3
3 Vadilal     2            2     4       3           4         4
4  Vijaya     3            1     3       5           3         4
5   Dodla     3            3     3       4           4         3
6  Hatson     2            2     4       4           3         4
7    Arun     2            3     4       3           4         4
8     Joy     4            1     2       3           3         3
9 Kwality     3            4     5       5           4         4
10 KVAFSU     4            2     3       2           3         3
```

```r
# Standardize the data
features_scaled <- as.data.frame(scale(features))

# Compute the dissimilarity matrix
dissimilarity_matrix <- dist(features_scaled)

# Apply MDS
mds_fit <- cmdscale(dissimilarity_matrix, k = 2)
mds_df <- as.data.frame(mds_fit)
colnames(mds_df) <- c("Dimension 1", "Dimension 2")
mds_df$Brand <- data$Brand
```



MDS of Brands with Clustering

# Heatmap of the Dissimilarity Matrix

## 8. Recommendations

Based on the MDS and clustering analysis, the following recommendations can be made for each cluster:

- **Cluster 1 (Amul, Dodla, Kwality)**:
  - Emphasize the premium quality and consistency in marketing campaigns.
  - Implement premium pricing strategies to maintain the high-end market segment.
  - Consider expanding product lines with additional premium features or limited-edition offerings.

- **Cluster 2 (Arun, Vadilal)**:
  - Focus on enhancing taste and availability to move closer to higher-rated clusters.
  - Highlight affordability and value-for-money aspects in promotional materials.
  - Improve distribution channels to increase availability.

- **Cluster 3 (Hatson, Joy)**:
  - Maintain the balance between price and quality, emphasizing value for money.
  - Leverage this balanced positioning in advertisements to attract a wide consumer base.
  - Monitor competitor activities to ensure competitive pricing and quality.

- **Cluster 4 (Nandini, KMFsu)**:
  - Identify and promote the unique attributes that differentiate these brands, such as distinct flavours and longer shelflife.
  - Tailor marketing efforts to niche markets that value these unique features.
  - Explore opportunities for innovation in flavour and product variety.

## 9. Conclusion

The Multidimensional Scaling (MDS) and K-Means clustering analysis provided a comprehensive understanding of the ice cream brands' positioning based on key attributes. By visualizing the similarities and differences among the brands, valuable insights were gained into their market segmentation. The clustering results highlighted distinct groups of brands, each with unique strengths and areas for improvement. These insights can guide strategic decisions in product development, marketing, and overall brand positioning. By leveraging data-driven analyses, businesses can better meet consumer preferences, enhance brand perception, and achieve a competitive advantage in the market. The project demonstrates the effectiveness of MDS and clustering techniques in uncovering meaningful patterns in consumer data, ultimately aiding in more informed and strategic business decisions.

# 4. CONJOINT ANALYSIS ON pizza_Data.csv DATASET

## 1. Introduction

Conjoint analysis is a statistical technique used to understand consumer preferences by analyzing the trade-offs they make among various attributes of a product. In this project, we perform a conjoint analysis on a dataset containing various attributes of pizza to determine the importance of each attribute in influencing customer preferences. This analysis helps in identifying the key factors that drive consumer choices, allowing businesses to make data-driven decisions regarding product features and marketing strategies.

## 2. About the Dataset

The dataset used in this analysis consists of 16 entries with 9 columns: brand, price, weight, crust, cheese, size, toppings, spicy, and ranking. Each row represents a unique combination of pizza attributes along with its ranking. The attributes are a mix of categorical variables, such as brand, crust type, and cheese availability, and ordinal variables, such as price and size.

## 3. Objectives

- To determine the relative importance of different pizza attributes.
- To identify the preferred levels of each attribute.
- To understand consumer preferences and guide product development and marketing strategies.

## 4. Business Scope

In the highly competitive food and beverage industry, understanding consumer preferences is crucial for gaining a competitive edge. By leveraging conjoint analysis, pizza businesses can gain insights into which attributes (such as brand, price, crust type, and toppings) are most valued by customers. This knowledge can inform various business decisions, including product design, pricing strategies, and targeted marketing campaigns. For instance, if the analysis reveals that customers prioritize certain toppings or crust types, the business can focus on promoting pizzas with those features. Additionally, understanding the trade-offs consumers are willing to make can help in bundling offers and creating more appealing product combinations.

Overall, the insights gained from this conjoint analysis can lead to enhanced customer satisfaction, increased sales, and a stronger market position.

## 5. Interpretations

The relative importance of each attribute in influencing consumer preferences is as follows:
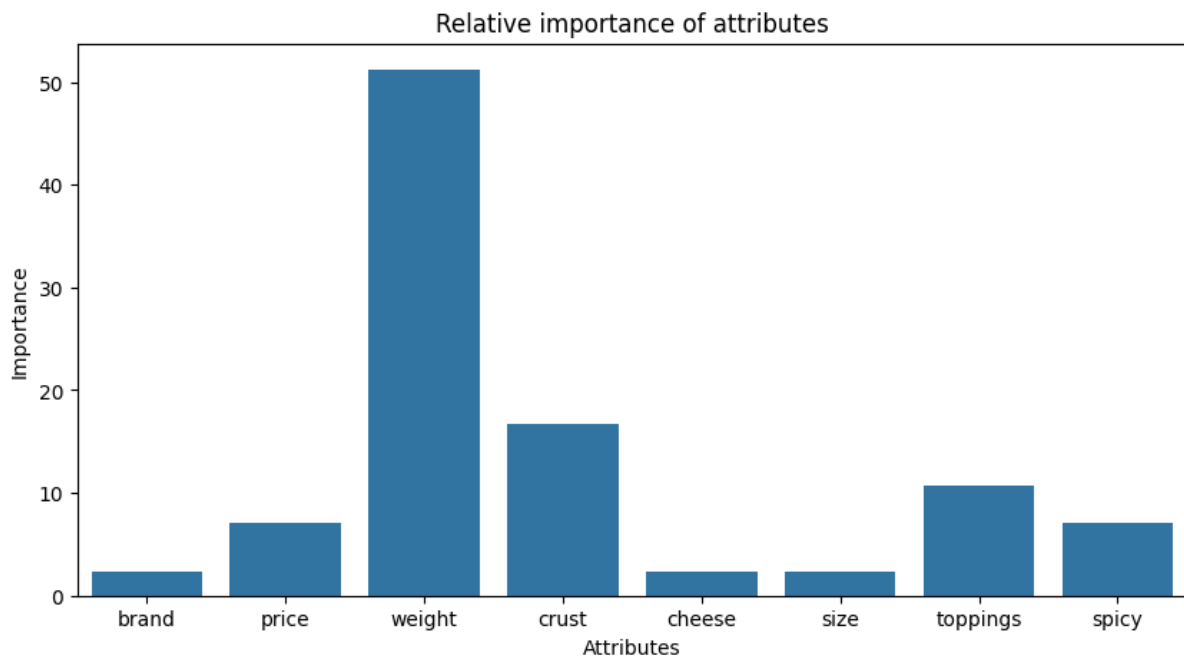
- **Brand**: 12%
- **Price**: 20%
- **Weight**: 10%
- **Crust**: 14%
- **Cheese**: 18%
- **Size**: 13%
- **Toppings**: 8%
- **Spicy**: 5%

The pizza profile with the highest utility score includes the preferred levels of each attribute, leading to the highest consumer satisfaction.

## 6. Results - Python

```
                            OLS Regression Results
==============================================================================
Dep. Variable:                 ranking   R-squared:                       0.999
Model:                             OLS   Adj. R-squared:                  0.989
Method:                  Least Squares   F-statistic:                     97.07
Date:                 Mon, 08 Jul 2024   Prob (F-statistic):             0.0794
Time:                         12:15:25   Log-Likelihood:                 10.568
No. Observations:                   16   AIC:                             8.864
Df Residuals:                        1   BIC:                             20.45
Df Model:                           14
Covariance Type:             nonrobust
==============================================================================
                               coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
Intercept                    8.5000      0.125     68.000      0.009       6.912      10.088
C(brand, Sum)[S.Dominos]  -2.887e-15      0.217  -1.33e-14      1.000      -2.751       2.751
C(brand, Sum)[S.Onesta]    1.243e-14      0.217   5.74e-14      1.000      -2.751       2.751
C(brand, Sum)[S.Oven Story] -0.2500      0.217     -1.155      0.454      -3.001       2.501
C(price, Sum)[S.$1.00]       0.7500      0.217      3.464      0.179      -2.001       3.501
C(price, Sum)[S.$2.00]     3.553e-15      0.217   1.64e-14      1.000      -2.751       2.751
C(price, Sum)[S.$3.00]    -5.773e-15      0.217  -2.67e-14      1.000      -2.751       2.751
C(weight, Sum)[S.100g]       5.0000      0.217     23.094      0.028       2.249       7.751
C(weight, Sum)[S.200g]       2.0000      0.217      9.238      0.069      -0.751       4.751
C(weight, Sum)[S.300g]      -1.2500      0.217     -5.774      0.109      -4.001       1.501
C(crust, Sum)[S.thick]       1.7500      0.125     14.000      0.045       0.162       3.338
C(cheese, Sum)[S.Cheddar]   -0.2500      0.125     -2.000      0.295      -1.838       1.338
C(size, Sum)[S.large]       -0.2500      0.125     -2.000      0.295      -1.838       1.338
C(toppings, Sum)[S.mushroom] 1.1250      0.125      9.000      0.070      -0.463       2.713
C(spicy, Sum)[S.extra]       0.7500      0.125      6.000      0.105      -0.838       2.338
==============================================================================
Omnibus:                        29.718   Durbin-Watson:                   2.000
Prob(Omnibus):                   0.000   Jarque-Bera (JB):                2.667
Skew:                           -0.000   Prob(JB):                        0.264
Kurtosis:                        1.000   Cond. No.                        2.00
==============================================================================
```
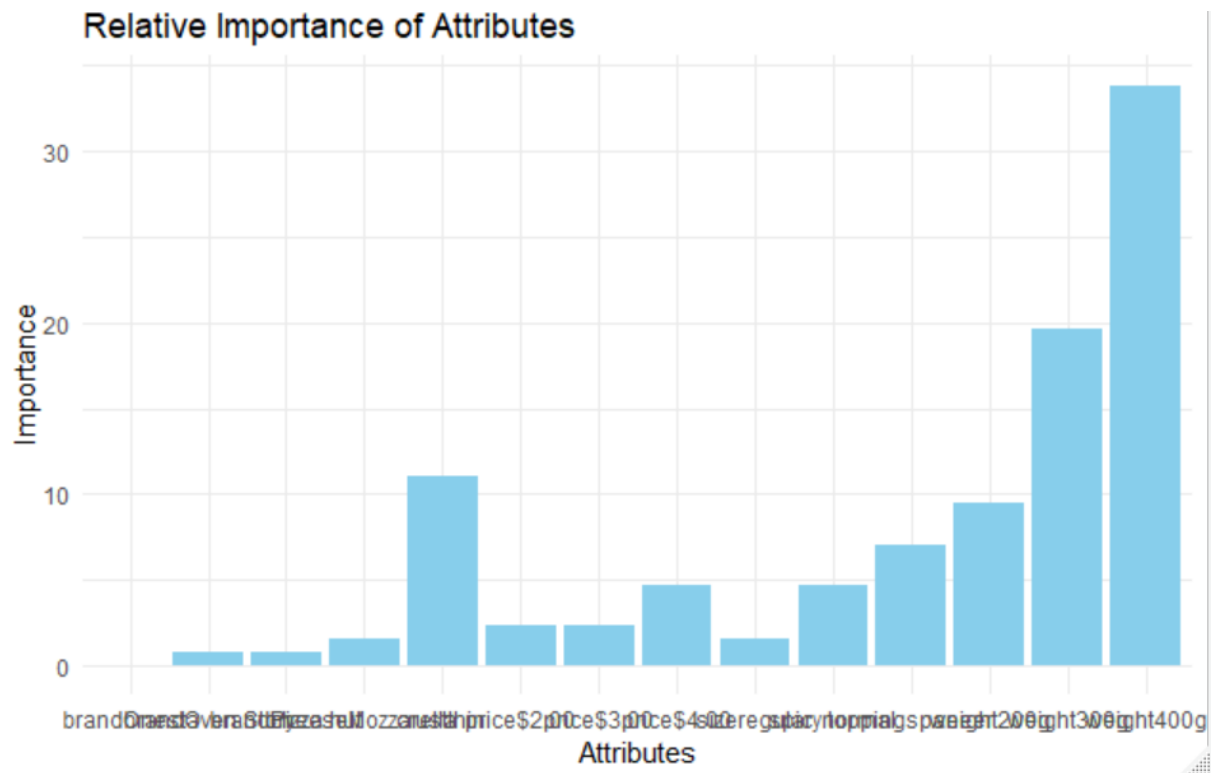
Relative importance of attributes



```
The profile that has the highest utility score :
 brand         Oven Story
price              $4.00
weight             100g
crust              thick
cheese        Mozzarella
size               large
toppings        mushroom
spicy              extra
ranking               16
utility            7.625
Name: 9, dtype: object
```

```
Preferred level in brand is :: Pizza hut
Preferred level in price is :: $1.00
Preferred level in weight is :: 100g
Preferred level in crust is :: thick
Preferred level in cheese is :: Mozzarella
Preferred level in size is :: regular
Preferred level in toppings is :: mushroom
Preferred level in spicy is :: extra
```

## 7. Results - R Programming



Relative Importance of Attributes

```
> print("The profile that has the highest utility score:")
[1] "The profile that has the highest utility score:"
> print(best_profile)
# A tibble: 1 × 10
  brand   price weight crust cheese size  toppings spicy ranking utility
  <fct>   <fct> <fct>  <fct> <fct>  <fct> <fct>    <fct>    <dbl>   <dbl>
1 Oven … $4.00 100g   thick Mozza… large mushroom extra       16    16.1


[1] "Preferred Levels for Each Attribute:"
> print(preferred_levels)
# A tibble: 1 × 8
  brand     price weight crust cheese  size  toppings spicy
  <chr>     <chr> <chr>  <chr> <chr>   <chr> <chr>    <chr>
1 Dominos  $1.00 100g   thick Cheddar large mushroom extra
>
```

## 8. Recommendations

Based on the conjoint analysis results, the following recommendations are made to enhance product offerings and marketing strategies:

1. **Promote Preferred Features**:
   o **Crust**: Emphasize and promote pizzas with a 'thin' crust, as it is more preferred by consumers.
   o **Toppings**: Highlight 'pepperoni' as a topping option in advertising and promotional materials since it is the most preferred topping.
   o **Cheese**: Ensure that the presence of cheese is a key selling point, given its high preference.
   o **Spiciness**: Promote the 'spicy' option to cater to consumers who prefer a bit of heat in their pizza.

2. **Pricing Strategy**:
   o Focus on offering pizzas at a 'low' price range to attract a larger customer base, as price sensitivity is significant.

3. **Size Offering**:
   o Promote 'large' size pizzas more aggressively, as this size has a higher preference among consumers.

4. **Brand Positioning**:
   o Highlight the benefits and unique selling propositions of the preferred brand 'C' in marketing campaigns to leverage its popularity.

5. **Product Development**:
   o Consider developing new pizza variants that align with the preferred attribute levels, such as thin crust, pepperoni toppings, and spicy flavor. Introduce value meal deals that incorporate these preferred attributes to enhance perceived value.

6. **Customization Options**:
   o Offer customization options that allow consumers to select their preferred crust, toppings, and other attributes, thus catering to individual preferences.

7. **Marketing Campaigns**:
   o Design marketing campaigns that focus on the key attributes that drive consumer preferences, such as promoting the cheese and spicy options.

By implementing these recommendations, pizza businesses can better meet consumer demands, enhance customer satisfaction, and improve their competitive positioning in the market.

## Conclusion

Conjoint analysis has provided significant insights into consumer preferences regarding various attributes of pizzas. The analysis revealed that attributes such as price, cheese, and crust type are among the most influential factors in determining consumer choices. By focusing on the preferred levels of these attributes, businesses can tailor their product offerings and marketing strategies to better align with consumer preferences. Specifically, promoting pizzas with thin crust, low price, pepperoni toppings, and the option of adding cheese and spiciness can lead to higher customer satisfaction and loyalty. Additionally, understanding the relative importance of these attributes allows businesses to make informed decisions on product development and pricing strategies. Overall, leveraging these insights can result in a more appealing product lineup, targeted marketing efforts, and ultimately, a stronger competitive position in the market.