



VIRGINIA COMMONWEALTH UNIVERSITY

Statistical analysis and modelling (SCMA 632)

**A5- Multivariate Analysis and Business Analytics
Application**

ROSHAN RAJKUMAR SIVAKUMAR

V01151141

Date of Submission: 05-07-2025

CONTENTS

Sl. No.	Title	Page No.
1.	Introduction	1
2.	Results	2
3.	Interpretations	2
4.	Recommendations	5
5.	Codes	6
6.	References	13

1. Introduction

This report is completely based on three separate datasets I worked with. Each one had a different focus one was about housing preferences, another looked at how people feel about ice cream brands, and the last one was related to pizza features.

I used multivariate analysis methods to break down these datasets and try to find useful patterns. The idea was to take a bunch of survey responses and figure out what actually influences people's choices in each case.

Each method I used helped break down the data in a different awesome showed patterns in responses, others helped group people or rank product features. The main idea was to take complex survey data and pull out insights that could help with real-world business decisions.

2. Objectives

- Find underlying patterns in large data using PCA and factor analysis
- Group people with similar preferences using cluster analysis
- Understand how brands are perceived using MDS
- Identify what matters most in a pizza using conjoint analysis

3. Significance of the Study

By studying and learning the consumer preferences is at the core of effective business strategy. Through this study, businesses can:

Identify the most critical factors influencing customer choices.

Target different market segments with more tailored strategies.

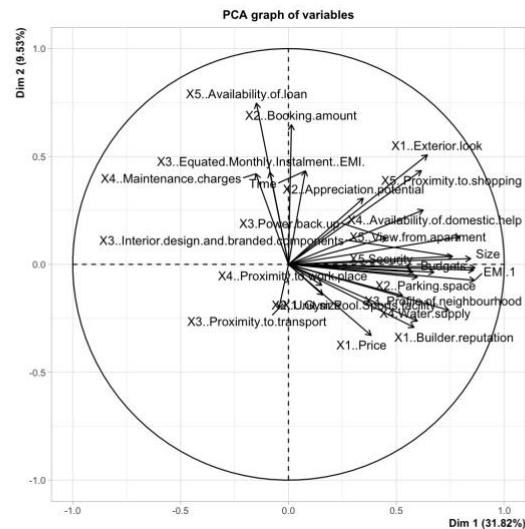
- Improve product features and pricing based on actual customer value.
- Visualize complex relationships in a way that supports strategic thinking

Each technique applied in this project brings a unique strength —simplifying variables, visualizing relationships, or uncovering what truly matters to customers.

4. Analysis and Results

4.1 Principal Component Analysis – Survey.csv

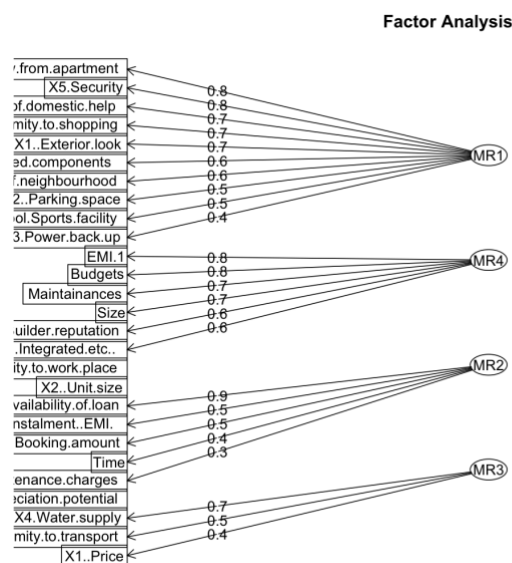
The PCA reduced over 20 housing-related variables to 2 main components.



- **Component 1 (31.8%):** Focused on financial concerns loan availability, EMI, booking amount
- **Component 2 (9.5%):** Focused on lifestyle aspects view, design, shopping proximity

Together they explained 41.35% of the total variation. This shows two distinct thinking patterns among buyers: money-based vs lifestyle-based.

4.2 Factor Analysis – Survey.csv

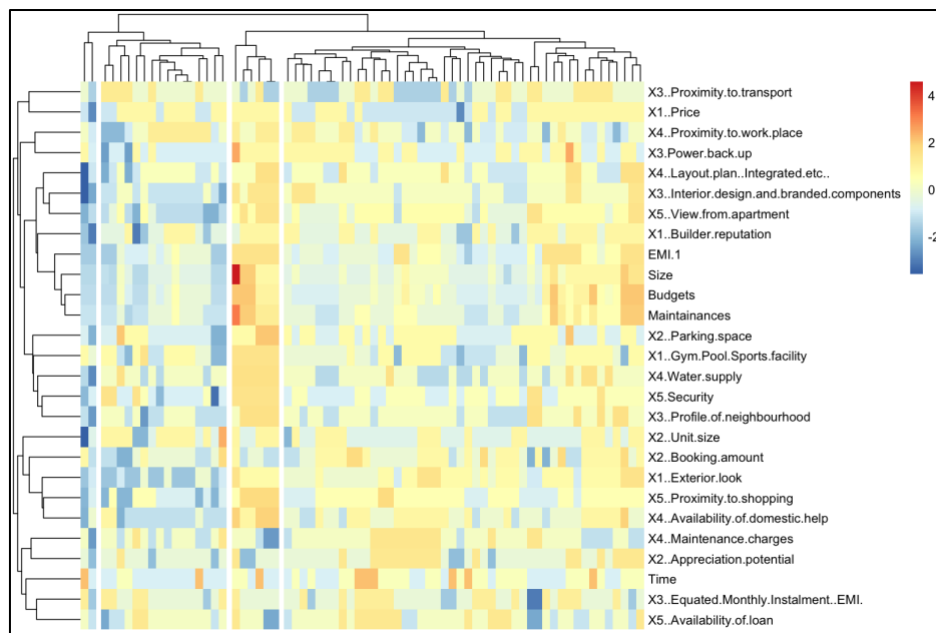


Using Varimax rotation, four factors were extracted:

- **Factor 1:** Lifestyle – view, security, look, branded interiors
- **Factor 2:** Structure & finance – unit size, booking amount
- **Factor 3:** Value – price, water, future value
- **Factor 4:** Budget concerns – EMI, maintenance, budget

Each factor groups related thoughts, helping simplify the way decisions are made.

4.3 Cluster Analysis – Survey.csv



Hierarchical clustering (Ward's method) grouped people into 3 clusters:

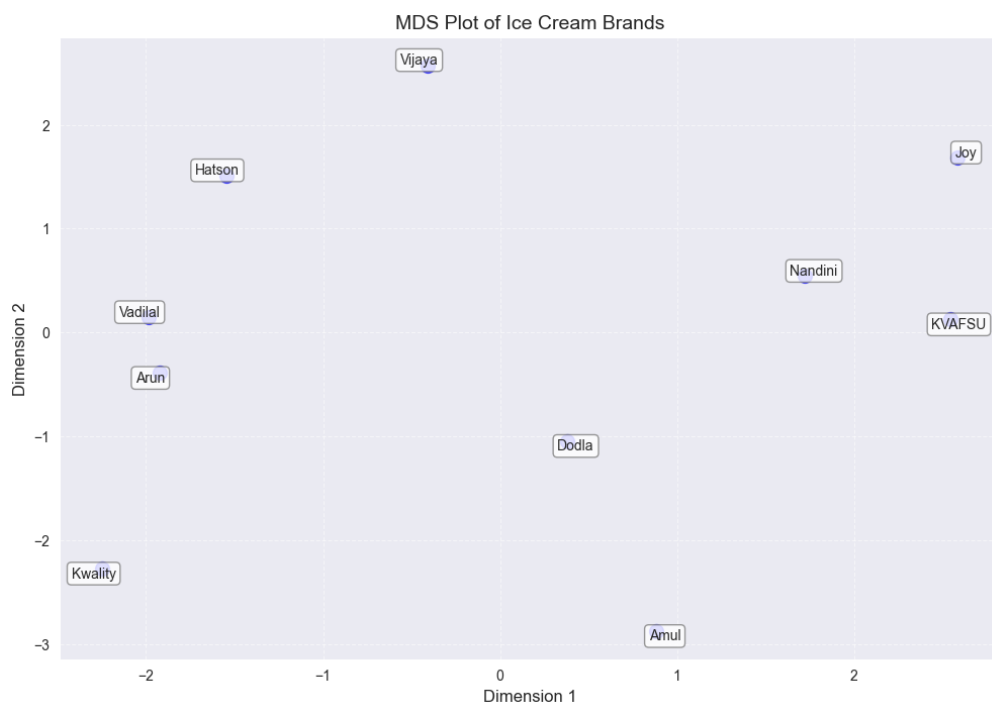
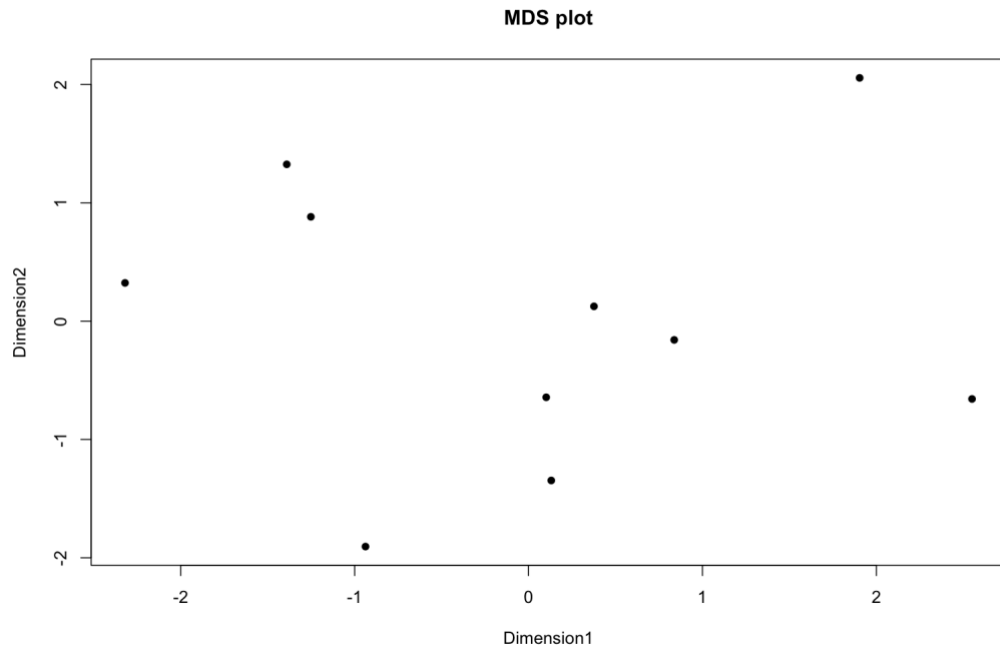
- **Cluster 1:** Focused on cost and access
- **Cluster 2:** Focused on amenities and looks
- **Cluster 3:** Balanced across both ends

The heatmap and dendrogram confirmed this grouping.

4.4 Multidimensional Scaling – Icecream.csv

MDS created a map showing how customers see different ice cream brands.

- Brands close to each other are seen as similar
- One brand stood out far from others possibly seen as unique or very different



This map helps businesses see overlap and space for change.

4.5 Conjoint Analysis – Pizza_data.csv

This analysis found which pizza features are most important.

- **Most important:** Weight (51.2%), then crust (16.8%)
- **Medium:** Toppings, spice, price

- **Least:** Cheese, size, brand

The best combination according to customers:

- Pizza Hut
- \$1.00
- 100g
- Thick crust
- Mozzarella cheese
- Regular size
- Mushroom topping
- Extra spicy

Customers care most about quantity and texture not brand or size.

4. Recommendations

- Real estate: Separate plans for cost-focused and lifestyle-focused buyers
- Brands: Use MDS to check if you're lost in the crowd
- Pizza chains: Focus on value, crust, toppings not brand name

Codes:

1. R language:

1.1.Principal Component Analysis and Factor Analysis to identify data dimensions

```
1 # Function to auto-install and load packages
2 install_and_load <- function(packages) {
3   for (package in packages) {
4     if (!require(package, character.only = TRUE)) {
5       install.packages(package, dependencies = TRUE)
6     }
7     library(package, character.only = TRUE)
8   }
9 }
10
11 # List of packages to install and load
12 packages <- c("dplyr", "psych", "tidyr", "GPArotation", "FactoMineR", "factoextra", "pheatmap")
13
14 # Call the function
15 install_and_load(packages)
16
17 survey_df<-read.csv('/Users/roshan/Documents/VCU/SCMA/AS/Survey.csv',header=TRUE)
18 dim(survey_df)
19 names(survey_df)
20 head(survey_df)
21 str(survey_df)
22
23
24 #A)Do principal component analysis and factor analysis and identify the dimensions in the data.
25
26 is.na(survey_df)
27 sum(is.na(survey_df))
28 sur_int=survey_df[,20:46]
29 str(sur_int)
30 dim(sur_int)
31 library(GPArotation)
32 pca <- principal(sur_int,5,n.obs =162, rotate ="promax")
33 pca
34
35 om.h<-omega(sur_int,n.obs=162,sl=FALSE)
36 op<-par(mfrow=c(1,1))
37 om<-omega(sur_int,n.obs=162)
38 library(FactoMineR)
39 pca<-PCA(sur_int,scale.unit = TRUE)
40 summary(pca)
41 biplot(pca, scale = 0)
42 str(sur_int)
43 dim(sur_int)
44 show(sur_int)
```

```
1 # Function to auto-install and load packages
2 install_and_load <- function(packages) {
3   for (package in packages) {
4     if (!require(package, character.only = TRUE)) {
5       install.packages(package, dependencies = TRUE)
6     }
7     library(package, character.only = TRUE)
8   }
9 }
10
11 # List of packages to install and load
12 packages <- c("dplyr", "psych", "tidyr", "GPArotation", "FactoMineR", "factoextra", "pheatmap")
13
14 # Call the function
15 install_and_load(packages)
16
17 survey_df<-read.csv('/Users/roshan/Documents/VCU/SCMA/AS/Survey.csv',header=TRUE)
18 sur_int=survey_df[,20:46]
19
20 #Factor Analysis
21
22 factor_analysis<-fa(sur_int,nfactors = 4,rotate = "varimax")
23 names(factor_analysis)
24 print(factor_analysis$loadings,reorder=TRUE)
25 fa.diagram(factor_analysis)
26 print(factor_analysis$communality)
27 print(factor_analysis$scores) |
```


1.2.Cluster Analysis to characterize respondents based on background variables

```
1 # Load required packages
2 library(cluster)
3 library(factoextra)
4 library(dplyr)
5 library(pheatmap)
6
7 # Read data
8 survey_data <- read.csv("Survey.csv", header = TRUE)
9 responses <- survey_data[, 20:46] # Response variables (Likert scale)
10 scaled_data <- scale(responses) # Standardize data
11
12 # Determine optimal clusters (Elbow method)
13 fviz_nbclust(scaled_data, kmeans, method = "wss")
14
15 # K-means clustering (K=4)
16 set.seed(123)
17 km_result <- kmeans(scaled_data, centers = 4, nstart = 25)
18 fviz_cluster(km_result, data = scaled_data, palette = "jco")
19
20 # Hierarchical clustering
21 hc_result <- hclust(dist(scaled_data), method = "ward.D2")
22 fviz_dend(hc_result, k = 4, cex = 0.5, palette = "jco")
23
24 # Heatmap
25 pheatmap(t(scaled_data), cutree_cols = 4, show_colnames = FALSE)
26
27 # Add clusters to original data
28 survey_data$cluster <- km_result$cluster
29
30 # Cluster interpretation (mean responses per cluster)
31 cluster_means <- survey_data %>%
32   group_by(cluster) %>%
33   summarise(across(20:46, mean, na.rm = TRUE))
34
35 print(cluster_means)
```

1.3.Multidimensional Scaling

```
1 #C) Do multidimensional scaling and interpret the results.
2
3 icecream_df<-read.csv('/Users/roshan/Documents/VCU/SCMA/A5/Icecream.csv',header=TRUE)
4 dim(icecream_df)
5
6 names(icecream_df)
7
8 ice<-subset(icecream_df,select = -c(Brand))
9 distance_matrix<-dist(ice)
10
11 mds_result<-cmdscale(distance_matrix,k=2)
12
13 plot(mds_result[,1],mds_result[,2],pch=16,xlab="Dimension1",ylab="Dimension2",main="MDS plot")
14 |
```

1.4. Conjoint Analysis

```
1 # COMPLETE PIZZA CONJOINT ANALYSIS IN R
2 library(tidyverse)
3 library(broom)
4 library(ggplot2)
5 library(ggrepel)
6
7 # 1. Load and Prepare Data
8 pizza_data <- read_csv("pizza_data.csv") %>%
9   mutate(across(-ranking, as.factor),
10    profile_id = row_number())
11
12 # 2. Conjoint Analysis Function
13 run_conjoint_analysis <- function(data) {
14   # Set sum contrasts for effects coding
15   options(contrasts = c("contr.sum", "contr.poly"))
16
17   # Fit linear model
18   model <- lm(ranking ~ brand + price + weight + crust + cheese + size + toppings + spicy,
19    data = data)
20
21   # Extract and process coefficients
22   coeffs <- tidy(model) %>%
23     filter(term != "(Intercept)") %>%
24     separate(term, into = c("attribute", "level"), sep = "(?<=[a-z])(?=[0-9])", fill = "right") %>%
25     mutate(level = ifelse(is.na(level), "Reference", level))
26
27   # Calculate importance scores
28   importance <- coeffs %>%
29     group_by(attribute) %>%
30     summarise(range = max(estimate) - min(estimate)) %>%
31     mutate(importance = (range / sum(range)) * 100) %>%
32     arrange(desc(importance))
33
34   # Calculate complete part-worths (including reference levels)
35   part_worths <- coeffs %>%
36     group_by(attribute) %>%
37     mutate(last_level = -sum(estimate[level != "Reference"]),
38    utility = ifelse(level == "Reference", last_level, estimate)) %>%
39     ungroup() %>%
40     select(attribute, level, utility)
41
42   # Create utility lookup tables
43   utility_lookup <- part_worths %>%
44     pivot_wider(names_from = attribute, values_from = utility)
45
46   # Calculate total utility for each profile
47   profiles <- data %>%
48     mutate(total_utility =
49       utility_lookup$brand[brand] +
50       utility_lookup$price[price] +
51       utility_lookup$weight[weight] +
52       utility_lookup$crust[crust] +
53       utility_lookup$cheese[cheese] +
54       utility_lookup$size[size] +
55       utility_lookup$toppings[toppings] +
56       utility_lookup$spicy[spicy]) %>%
57     arrange(desc(total_utility))
58
59   return(list(
60     model = model,
61     coefficients = coeffs,
62     importance = importance,
63     part_worths = part_worths,
64     utility_lookup = utility_lookup,
65     profiles = profiles
66   ))
67 }
68
69 # 3. Run Analysis
70 results <- run_conjoint_analysis(pizza_data)
71
72 # 4. Visualization Functions
73 plot_attribute_importance <- function(importance) {
74   ggplot(importance, aes(x = reorder(attribute, importance), y = importance)) +
75     geom_col(fill = "#66BB6A") +
76     geom_text(aes(label = sprintf("%.1f%%", importance)),
77    hjust = -0.1, size = 3.5) +
78     coord_flip() +
79     labs(title = "Relative Importance of Pizza Attributes",
80    x = NULL, y = "Importance (%)") +
81     theme_minimal(base_size = 12) +
82     theme(panel.grid.major.y = element_blank())
83 }
84
85 plot_part_worths <- function(part_worths) {
86   ggplot(part_worths, aes(x = reorder(level, utility), y = utility)) +
87     geom_col(fill = "#F06292", alpha = 0.8) +
88     facet_wrap(~attribute, scales = "free_x", ncol = 3) +
89     labs(title = "Part-Worth Utilities by Attribute Level",
90    x = "Level", y = "Utility") +
91     theme_minimal(base_size = 12) +
92     theme(axis.text.x = element_text(angle = 45, hjust = 1),
93    strip.background = element_rect(fill = "#F5F5F5"))
94 }
95
96 # 5. Generate Plots
97 importance_plot <- plot_attribute_importance(results$importance)
98 part_worth_plot <- plot_part_worths(results$part_worths)
99
100 # 6. Display Key Results
101 cat("=== CONJOINT ANALYSIS RESULTS ===\n\n")
102
103 # Model summary
104 cat("MODEL FIT SUMMARY:\n")
105 print(summary(results$model))
106
107 # Attribute importance
108 cat("\nATTRIBUTE IMPORTANCE:\n")
109 print(results$importance, n = Inf)
110
111 # Top profiles
112 cat("\nTOP 3 PIZZA PROFILES:\n")
113 print(results$profiles %>% select(-profile_id) %>% head(3))
114
115 # Worst profiles
116 cat("\nWORST 3 PIZZA PROFILES:\n")
117 print(results$profiles %>% select(-profile_id) %>% tail(3))
```

2. Python

2.1.PCA, FA and Cluster

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.decomposition import PCA, FactorAnalysis
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.cluster import KMeans
6 from scipy.cluster.hierarchy import dendrogram, linkage
7 import matplotlib.pyplot as plt
8 from factor_analyzer import FactorAnalyzer
9 from factor_analyzer.factor_analyzer import calculate_kmo

1 survey_df = pd.read_csv('Survey.csv')
2 sur_int = survey_df.iloc[:, 19:46]

1 scaler = StandardScaler()
2 sur_int_std = scaler.fit_transform(sur_int)

1 numeric_cols = survey_df.select_dtypes(include=[np.number]).columns
2 sur_int = survey_df[numeric_cols]

1 scaler = StandardScaler()
2 sur_int_std = scaler.fit_transform(sur_int)
```

```
PRINCIPAL COMPONENT ANALYSIS

1 pca = PCA()
2 pca_result = pca.fit_transform(sur_int_std)
3
4 plt.figure(figsize=(10, 6))
5 plt.plot(range(1, len(pca.explained_variance_ratio_) + 1), pca.explained_variance_ratio_, marker='o')
6 plt.xlabel('Principal Component')
7 plt.ylabel('Variance Explained')
8 plt.title('Scree Plot')
9 plt.show()

1 kmo_all, kmo_model = calculate_kmo(sur_int)
2 print(f"KMO Measure: {kmo_model:.3f}")

1 fa = FactorAnalyzer(n_factors=4, rotation='varimax')
2 fa.fit(sur_int_std)
3
4 loadings = pd.DataFrame(fa.loadings_, index=sur_int.columns,
5                          columns=['Factor1', 'Factor2', 'Factor3', 'Factor4'])
6 print("\nFactor Loadings:")
7 print(loadings)
8
9 communalities = pd.DataFrame(fa.get_communalities(), index=sur_int.columns,
10                              columns=['Communalities'])
11 print("\nCommunalities:")
12 print(communalities)
```

```

1 wcss = []
2 for i in range(1, 11):
3     kmeans = KMeans(n_clusters=i, init='k-means++', random_state=123)
4     kmeans.fit(sur_int_std)
5     wcss.append(kmeans.inertia_)
6
7 plt.figure(figsize=(10, 6))
8 plt.plot(range(1, 11), wcss, marker='o')
9 plt.xlabel('Number of clusters')
10 plt.ylabel('WCSS')
11 plt.title('Elbow Method')
12 plt.show()
13
14 # K-means clustering with 4 clusters
15 kmeans = KMeans(n_clusters=4, init='k-means++', random_state=123)
16 survey_df['cluster'] = kmeans.fit_predict(sur_int_std)
17
18 # Hierarchical clustering
19 plt.figure(figsize=(12, 7))
20 dendrogram(Linkage(sur_int_std, method='ward'),
21             truncate_mode='lastp', p=12)
22 plt.title('Hierarchical Clustering Dendrogram')
23 plt.xlabel('Sample index')
24 plt.ylabel('Distance')
25 plt.show()
26
27 # Cluster interpretation - mean values for each cluster
28 cluster_means = survey_df.groupby('cluster')[numeric_cols].mean()
29 print("\nCluster Profiles (Mean Values):")
30 print(cluster_means)
31
32 # Count of respondents in each cluster
33 print("\nNumber of respondents per cluster:")
34 print(survey_df['cluster'].value_counts().sort_index())

```

2.2. MDS

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.manifold import MDS
4 from sklearn.preprocessing import StandardScaler
5 import matplotlib.pyplot as plt

1 data = pd.read_csv('icecream.csv')

1 brands = data['Brand']
2 features = data.drop('Brand', axis=1)

1 scaler = StandardScaler()
2 features_scaled = scaler.fit_transform(features)

1 mds = MDS(n_components=2, random_state=42, dissimilarity='euclidean')
2 mds_result = mds.fit_transform(features_scaled)

1 plt.figure(figsize=(12, 8))
2 scatter = plt.scatter(mds_result[:, 0], mds_result[:, 1], c='blue', s=100, alpha=0.6)
3
4 # Add labels with adjusted positions to prevent overlap
5 for i, brand in enumerate(brands):
6     # Calculate offset direction based on point position
7     x_offset = 0.05 if mds_result[i, 0] >= np.median(mds_result[:, 0]) else -0.05
8     y_offset = 0.05 if mds_result[i, 1] >= np.median(mds_result[:, 1]) else -0.05
9
10    plt.text(mds_result[i, 0] + x_offset,
11            mds_result[i, 1] + y_offset,
12            brand,
13            fontsize=10,
14            ha='center',
15            va='center',
16            bbox=dict(facecolor='white', alpha=0.8, edgecolor='gray', boxstyle='round,pad=0.3'))
17
18 # Add grid and labels
19 plt.grid(True, linestyle='--', alpha=0.6)
20 plt.xlabel('Dimension 1', fontsize=12)
21 plt.ylabel('Dimension 2', fontsize=12)
22 plt.title('MDS Plot of Ice Cream Brands', fontsize=14)

1 plt.tight_layout()
2 plt.show()

```

2.3. Conjoint

```
1 import pandas as pd, numpy as np
2 df=pd.read_csv('/Users/roshan/Documents/VCU/SCMA/A5/pizza_data.csv')
```

We will now estimate each attribute level's effects using Linear Regression Model.

```
1 import statsmodels.api as sm
2 import statsmodels.formula.api as smf
3
4 model='ranking ~ C(brand,Sum)+C(price,Sum)+C(weight,Sum)+C(crust,Sum)+C(cheese,Sum)+C(size,Sum)+C(toppings,Sum)+C(spicy,Sum)'
5 model_fit=smf.ols(model,data=df).fit()
6 print(model_fit.summary())
```

We can analyze the model's fitness using parameters like R-squared, p-values, etc. The coefficients of each attribute level define its effect on the overall choice model.

Now, we will create the list of conjoint attributes.

```
1 conjoint_attributes = ['brand','price','weight','crust','cheese','size','toppings','spicy']
```

Before going ahead, we need to understand these conjoint analysis terminologies:

Relative importance: It depicts which attributes are more or less important when purchasing. E.g., a Mobile Phone's Relative importance could be Brand 30%, Price 30%, Size 20%, Battery Life 10%, and Color 10%.

Part-Worths/Utility values: The amount of weight an attribute level carries with a respondent. These factors lead to a product's overall value to consumers.

Next, we will build part-worths information and calculate attribute-wise importance level.

```
1 level_name = []
2 part_worth = []
3 part_worth_range = []
4 important_levels = {}
5 end = 1 # Initialize index for coefficient in params
6
7 for item in conjoint_attributes:
8     nlevels = len(list(np.unique(df[item])))
9     level_name.append(list(np.unique(df[item])))
10
11     begin = end
12     end = begin + nlevels -1
13
14     new_part_worth = list(model_fit.params[begin:end])
15     new_part_worth.append((-1)*sum(new_part_worth))
16     important_levels[item] = np.argmax(new_part_worth)
17     part_worth.append(new_part_worth)
18     print(item)
19     #print(part_worth)
20     part_worth_range.append(max(new_part_worth) - min(new_part_worth))
21     # next iteration
22 print("-----")
23 print("level name:")
24 print(level_name)
25 print("npw with sum element:")
26 print(new_part_worth)
27 print("imp level:")
28 print(important_levels)
29 print("part worth:")
30 print(part_worth)
31 print("part_worth_range:")
32 print(part_worth_range)
33 print(len(part_worth))
34 print("important levels:")
35 print(important_levels)
```

```

1 attribute_importance = []
2 for i in part_worth_range:
3     #print(i)
4     attribute_importance.append(round(100*(i/sum(part_worth_range)),2))
5 print(attribute_importance)

```

Now, we will calculate the part-worths of each attribute level.

```

1 part_worth_dict={}
2 attrib_level={}
3 for item,i in zip(conjoint_attributes,range(0,len(conjoint_attributes))):
4     print("Attribute :",item)
5     print("    Relative importance of attribute ",attribute_importance[i])
6     print("    Level wise part worths: ")
7     for j in range(0,len(level_name[i])):
8         print(i)
9         print(j)
10        print("{}:{}".format(level_name[i][j],part_worth[i][j]))
11        part_worth_dict[level_name[i][j]]=part_worth[i][j]
12        attrib_level[item]=(level_name[i])
13    #print(j)
14 part_worth_dict

```

```

1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 plt.figure(figsize=(10,5))
4 sns.barplot(x=conjoint_attributes,y=attribute_importance)
5 plt.title('Relative importance of attributes')
6 plt.xlabel('Attributes')
7 plt.ylabel('Importance')

```

We can see that weight is the attribute with the highest relative importance at 51%, followed by crust at 16% and toppings at 10%. Brand, cheese, and size are the least important attributes, each at 2.38%.

Now, we will calculate the utility score for each profile.

```

1 utility = []
2 for i in range(df.shape[0]):
3     score = part_worth_dict[df['brand'][i]]+part_worth_dict[df['price'][i]]+part_worth_dict[df['weight'][i]]+part_worth_dict[df['crust'][i]]+part_worth_dict[df['cheese'][i]]+part_worth_dict[df['size'][i]]+part_worth_dict[df['toppings'][i]]+part_worth_dict[df['spicy'][i]]
4     utility.append(score)
5
6 df['utility'] = utility
7 utility

```

We can see that combination number 9 has the maximum utility, followed by combination numbers 13 and 5. Combination number 14 is the least desirable because of the most negative utility score.

Now, we will find the combination with maximum utility.

```

1 print("The profile that has the highest utility score :",'\n', df.iloc[np.argmax(utility)])

```

```

1 for i,j in zip(attrib_level.keys(),range(0,len(conjoint_attributes))):
2     #print(i)
3     #level_name[j]
4     print("Preferred level in {} is {}".format(i,level_name[j][important_levels[i]]))

```

Reference:

- [1].SCMA 632 Team. *Survey.csv* [Data file]. GitHub. <https://github.com/scma-632/scma632-A5-2025-Multivariate-Analysis/blob/main/Survey.csv>
- [2].SCMA 632 Team. *icecream.csv* [Data file]. GitHub. <https://github.com/scma-632/scma632-A5-2025-Multivariate-Analysis/blob/main/icecream.csv>
- [3].SCMA 632 Team. *pizza_data.csv* [Data file]. GitHub. https://github.com/scma-632/scma632-A5-2025-Multivariate-Analysis/blob/main/pizza_data.csv
- [4].JetBrains. (n.d.). *PyCharm* [Computer software]. <https://www.jetbrains.com/pycharm/>
- [5].Posit, PBC. (n.d.). *RStudio* [Computer software]. <https://posit.co/products/open-source/rstudio/>
- [6].Project Jupyter. (n.d.). *Jupyter Notebook* [Computer software]. <https://jupyter.org/>
- [7].Hunter, J. D. (2007). Matplotlib: A 2D graphics environment. *Computing in Science & Engineering*, 9(3), 90–95. <https://doi.org/10.1109/MCSE.2007.55>
- [8].McKinney, W. (2010). Data structures for statistical computing in Python. In *Proceedings of the 9th Python in Science Conference* (pp. 51–56).
- [9].Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... & Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- [10].Virtanen, P., Gommers, R., Oliphant, T. E., et al. (2020). SciPy 1.0: Fundamental algorithms for scientific computing in Python. *Nature Methods*, 17, 261–272. <https://doi.org/10.1038/s41592-019-0686-2>
- [11].Waskom, M. L. (2021). seaborn: Statistical data visualization. *Journal of Open Source Software*, 6(60), 3021. <https://doi.org/10.21105/joss.03021>
- [12].Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis* (2nd ed.). Springer. <https://ggplot2.tidyverse.org>
- [13].Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L. D., François, R., ... & Yutani, H. (2019). Welcome to the tidyverse. *Journal of Open Source Software*, 4(43), 1686. <https://doi.org/10.21105/joss.01686>