# Market Segmentation using Machine Learning:

Kanishka Arora

June 25, 2024

Indian Institute of Space Science and Technology

# Contents

# List of Tables

# List of Figures

<small>CHAPTER</small> **1**

# Market Segmentation

---

## 1.1 What is Marketing?

```
┌─────────────────────────────────────┐
│   Identify Consumer Needs and Desires │
└─────────────────────────────────────┘
                    │
                    ▼
          ┌──────────────────┐
          │  Develop Offers   │
          └──────────────────┘
                    │
                    ▼
          ┌──────────────────┐
          │ Match Offers to Needs │
          └──────────────────┘
                    │
                    ▼
        ┌────────────────────────────┐
        │ Benefit Consumers and Suppliers │
        └────────────────────────────┘
                    │
                    ▼
        ┌────────────────────────────┐
        │ Drive Marketing Planning Process │
        └────────────────────────────┘
```

The above is a basic flowchart that explains the meaning of Marketing. It is broadly divided into two categories, which are,

1. **Strategic Marketing**
   It focuses largely on the long-term goal and direction of the organization, without paying much attention to the roadmap of short-term marketing, which leads to long-term goals. It is typically associated with consumer needs and desires, strengths and weaknesses within the organization, external threats, and opportunities for the organization.
   Once the above-listed factors have been studied and accounted for, we make two key decisions in the strategic marketing process i.e. which consumers to focus on (segmentation and targeting), and which image of the organization to create in the market (positioning).

**Strategic Marketing**

Bad | Good

|  | Bad | Good |



**Fig. 1.1:** SWOT Analysis

2. **Tactical Marketing**
   The tactical marketing plan is invested in the short-term goals and actions required to bring forth the long-term desired results of the organization. Only when segmentation targeting and positioning have been done can we begin the tactical marketing process.
   Tactical marketing broadly covers four areas, namely

   - Product
   - Price
   - Place
   - Promotion

In an organization, the combination of both good strategic marketing and good tactical marketing leads to the best possible outcome. But one must remember,

So the foundation of growth for the organization's success lies in proper strategic marketing planning.

## 1.2   What is Market Segmentation?

Market segmentation is essential for marketing success, as it allows companies to target specific groups within a larger market. Smith first proposed this strategy in 1956, who defined it as dividing a heterogeneous market into smaller, homogeneous markets. Effective market segmentation means that consumers within a segment share similar characteristics important to management, while those in different segments are distinct in those characteristics.

Segmentation criteria can range from single attributes like age or gender to multiple factors like benefits sought or values held. This tailored approach helps organizations better meet the needs of each segment, leading to higher sales and stronger market positioning.

A **concentrated market strategy** focuses on one segment, making it ideal for resource-limited organizations facing competition. However, it carries the risk of relying on a single segment. A **differentiated market strategy** targets multiple segments with customized marketing efforts, suitable for mature markets. An undifferentiated market strategy uses the same product and marketing mix for the entire market, which can work for resource-rich organizations or new products.

## 1.3   What are the benefits of Market Segmentation?

These are the following benefits of Market Segmentation,

- Market segmentation encourages organizations to assess their current position and future goals.
- It helps organizations identify their strengths relative to competitors.
- Forces reflection on consumer preferences and needs.
- Leads to critical insights and perspectives.
- Enhances understanding of consumer differences.
- Improves alignment between organizational strengths and consumer needs.
- Forms the basis for long-term competitive advantage.
- Facilitates market dominance in niche segments.
- Enables customization of products/services for specific consumer groups.
- Supports micro marketing and hyper-segmentation strategies.
- Enhances ROI by focusing marketing efforts more effectively.
- Essential for small organizations focusing on distinct consumer needs.
- Effective in targeting sales efforts towards specific consumer groups.
- Contributes to team building within organizations.
- Improves communication and information sharing across organizational units.

## 1.4   What are the costs of Execution of Market Segmentation?

Implementing market segmentation requires significant investment in terms of time, human resources, and finances. It involves thorough analysis, custom marketing strategies, and ongoing monitoring. Success promises a competitive advantage, but failure can lead to wasted resources and demoralization among the staff involved.

Hence the organization has to make a wise decision to move with the market segmentation analysis and strategy.

CHAPTER **2**

# Market Segmentation Analysis

> The process of grouping consumers into naturally existing or artificially created segments of consumers who share similar product preferences or characteristics.

## 2.1 Layers of Market Segmentation Analysis

Market segmentation involves three critical layers,

- **Technical Process**
  Led by data analysts, involves data collection, exploration, and statistical segment extraction to group consumers effectively. The grouping of consumers can always only be as good as the data provided by the segment extraction method.

- **Analytical Tasks**
  Includes profiling and describing each segment to guide strategic marketing decisions and customize the marketing mix.

- **Implementation and Strategic Decision-making**
  Requires organizational commitment to long-term segmentation strategies based on market opportunities identified. User involvement is crucial throughout, from data collection to selecting target segments and developing tailored marketing plans.

## 2.2 Approaches in Market Segmentation Analysis

Here are two approaches to Market Segmentation Analysis,

- **Based on Organizational Constraints**

```
┌─────────────────────────────────────────┐
│      Market Segmentation Approaches       │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│ Segment Revolution(Quantitative Survey-Based) │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│ Segment Evolution(Refining Existing Segments) │
└─────────────────────────────────────────┘
                    │
                    ▼
┌─────────────────────────────────────────┐
│  Segment Mutation(Exploratory Research)   │
└─────────────────────────────────────────┘
```

- **Based on choice of Segmentation Variables**



## 2.3   More on Data Driven Market Segmentation Approaches

Over here data analysis created the solution for the market segmentation.

- **Rarity of Natural Segments:** Naturally occurring, distinct, and well-separated market segments are rare in real consumer data.

- **Conceptual Approaches:** There are three approaches to data-driven market segmentation: natural, reproducible, and constructive segmentation.

- **Natural Segmentation:** Assumes pre-existing segments in the data, aiming to uncover these distinct market segments.

- **Reproducible Segmentation:** Acknowledges some structure in the data, allowing for consistent and reliable segmentation solutions across different analyses.

- **Constructive Segmentation:** Involves creating artificial segments in the absence of natural data structure, which can still provide strategic marketing value.

- **Data Structure Analysis:** Essential to analyze data structure before segmentation to avoid methodological errors and misinterpretations. This can involve repeated segmentation with different algorithms.

- **Practical Implications:** Most data sets contain some exploitable structure, even if not in clear cluster form, making collaborative constructive segmentation useful for developing effective marketing strategies.

In the proceeding chapters we will discuss the step by step analysis of market segmentation.

# Steps in Market Segmentation

In this chapter we will discuss the steps in Market Segmentation one by one.

## 3.1 Step 1 - Deciding Not to Segment

Market segmentation, a fundamental strategy in marketing, demands careful consideration due to its profound implications for organizations. It requires a steadfast, long-term commitment akin to a marriage, involving substantial investments in research, surveys, product development, and tailored marketing communications. Before embarking on segmentation, organizations must weigh these costs against potential sales increases to ensure profitability.

Implementing segmentation often necessitates organizational adjustments, such as restructuring around market segments rather than products, and potentially developing new products or adjusting pricing and distribution strategies. Strategic business units focused on segments can facilitate ongoing adaptation to evolving market needs.

Below are the following key implementation barriers,

- **Senior Management Involvement**
  Lack of leadership, proactive championing, commitment, and involvement in the segmentation process by senior executives.

- **Resource Allocation**
  Insufficient allocation of resources (financial and human) for conducting the segmentation analysis and implementing the strategy.

- **Organizational Culture**
  Lack of market or consumer orientation, resistance to change, inadequate communication and sharing of information across departments, short-term thinking, and office politics.

- **Training and Expertise**
  Lack of understanding and expertise in market segmentation among senior management and the segmentation team.

- **Formal Marketing Function**
  Absence of a formal marketing function or qualified marketing experts within the organization.

- **Data Management**
  Lack of qualified data managers and analysts to handle segmentation data effectively.

- **Objective Restrictions**
  Financial constraints or organizational structures that hinder the execution of segmentation strategies.

- **Process-related Issues**
  Unclear objectives, inadequate planning, lack of structured processes, unclear allocation of responsibilities, and time constraints during the segmentation process.

- **Acceptance of Management Techniques**
  Resistance to using management techniques (like segmentation analysis) that senior management does not fully understand.

Ultimately, the decision to pursue market segmentation should be made at the highest executive level and consistently communicated throughout the organization to foster sustained commitment and alignment across all operational units.

## 3.2   Step 2 - Identifying the Ideal Target Segment

### 3.2.1   Segment Evaluation Criteria

In this step we focus on defining segment evaluation criteria, comprising essential knock-out criteria and flexible attractiveness criteria. While knock-out criteria are non-negotiable and automatically disqualify segments lacking essential features, attractiveness criteria are selectively applied by the segmentation team to gauge the relative desirability of compliant segments. This structured approach, supported by a diverse array of proposed criteria detailed in the literature, facilitates a thorough evaluation process essential for identifying and prioritizing market segments effectively.

### 3.2.2   Knock out Criteria

Knock-out criteria are used to determine if market segments resulting from the market segmentation analysis qualify to be assessed using segment attractiveness criteria, which are the following,

- The segment must be homogeneous.
- The segment must be distinct.
- The segment must be large enough.
- The segment must match the strengths of the organization.
- Members of the segment must be identifiable.
- The segment must be reachable.

### 3.2.3   Attractiveness Criteria

Segments also require attractiveness criteria pertaining to any given situation, and it is not binary in nature. Attractiveness criteria are not assessed as either complying or non-complying, rather it is rated in a spectrum.

### 3.2.4   Implementing a Structured Process

- Structured processes are widely acknowledged as beneficial for assessing market segments.
- The most popular method is the segment evaluation plot, which plots segment attractiveness against organizational competitiveness.
- Segment attractiveness and organizational competitiveness values are determined by the segmentation team, with no universal criteria applicable to all organizations.

- Criteria for both segment attractiveness and organizational competitiveness must be negotiated and agreed upon, with a recommendation to use no more than six factors.

- A team approach is optimal, where a core team proposes initial solutions and an advisory committee, representing all organizational units, discusses and possibly modifies these proposals.

- Including representatives from all organizational units is crucial as each unit has different business perspectives and all units are stakeholders in the segmentation strategy.

- Segment attractiveness criteria should be selected early to ensure relevant data is captured during collection and to simplify target segment selection later.

- The segmentation team should finalize approximately six segment attractiveness criteria, each with a weighted importance relative to the others.

- Team members should distribute 100 points across the criteria to determine weights, negotiating until agreement is reached, with approval ideally sought from the advisory committee for a comprehensive perspective.

## 3.3   Step 3 - Data Collection

In this section we will discuss the collection of data for Segmentation.

- Empirical data is the foundation for both commonsense and data-driven market segmentation, used to identify, create, and describe market segments.

- Commonsense segmentation typically uses one single characteristic, such as gender, as the segmentation variable, with other characteristics serving as descriptor variables.

- Descriptor variables, such as socio-demographics and media behavior, are used to describe segments in detail, aiding in the development of effective marketing strategies.

- Data-driven market segmentation uses multiple segmentation variables to identify or create market segments that are useful to the organization.

- The quality of empirical data is critical for assigning individuals to the correct market segments and for accurately describing these segments.

- Empirical data for segmentation can come from surveys, observations (such as scanner data), and experimental studies, with data reflecting actual consumer behavior being preferable.

- Survey data, while common, can be unreliable for behaviors that are socially desirable, so exploring a range of data sources is recommended.

### 3.3.1   Segmentation Criteria

> The term segmentation criterion relates to the nature of the information used for market segmentation.The most common segmentation criteria include geographic, socio-demographic, psychographic and bhevioural.

| Geographic Segmentation | | Socio-demographic Segmentation |
| --- | --- | --- |
| | Segmentation Criteria | |
| Psychographic Segmentation | | Behavioral Segmentation |

Below we will discuss the criteria individually,

### 3.3.2   Geographic Segmentation

- Geographic segmentation uses the consumer's location of residence as the primary criterion for forming market segments, making it straightforward for targeting communication messages and channels.

- This method is particularly useful for companies needing to cater to language differences and regional preferences.

- The key advantage is the ease of assigning each consumer to a geographic unit, allowing precise targeting through local media.

- The main disadvantage is that geographic location alone often fails to capture other relevant consumer characteristics, such as the benefits sought in products or services.

- Geographic segmentation has seen a revival in international market studies, although it faces challenges like ensuring meaningful segmentation variables across regions and avoiding cultural bias in survey responses.

### 3.3.3   Socio-Demographic Segmentation

- Typical socio-demographic segmentation criteria include age, gender, income, and education.

- Socio-demographic segments are particularly useful in industries like luxury goods (high income), cosmetics (gender-specific marketing), and baby products (gender-specific needs).
- In industries such as retirement villages and tourism resorts, age and family status (having children) play crucial roles in segmenting consumers.
- Like geographic segmentation, socio-demographic criteria allow easy determination of segment membership for each consumer.
- Socio-demographic factors can explain specific product preferences in some cases (e.g., family vacation choices), but often they do not fully account for consumer behavior variations.
- Studies suggest that demographics explain a small percentage (about 5%) of consumer behavior variance, with values, tastes, and preferences being more influential in buying decisions.

### 3.3.4   Psychographic Segmentation

- Psychographic segmentation categorizes people based on psychological criteria such as beliefs, interests, preferences, aspirations, and benefits sought in products.
- Benefit segmentation, focuses on identifying consumer benefits as a key psychographic approach.
- Lifestyle segmentation, another popular method, segments based on activities, opinions, and interests.
- Psychographic criteria are complex compared to geographic or socio-demographic criteria, often requiring multiple variables like travel motives or perceived risks in segmentation studies.
- Psychographic segmentation provides deeper insights into consumer behavior by focusing on underlying motivations. For example, tourists motivated by cultural exploration are likely to choose destinations rich in cultural experiences.

### 3.3.5   Behavioral Segmentation

- Behavioral segmentation extracts segments based on similarities in actual or reported behaviors, such as prior product experience, purchase frequency, amount spent per purchase occasion, and information search habits.
- Research indicates that behaviors reported by tourists are often more effective segmentation criteria than geographic variables.
- The primary advantage of behavioral approaches is their basis in actual behavior, rather than stated or intended behavior, which allows for segmentation by the most relevant consumer similarities.

- Examples include using actual consumer expenses and purchase data across product categories as segmentation variables and longitudinal brand choice behavior.

- Using behavioral data eliminates the need for developing valid measures of psychological constructs but accessing such data can be challenging, particularly when including potential customers who haven't yet purchased the product.

### 3.3.6   Data from Survey Studies

Most of the analyses in market segmentation are based on survey data as it is cheap and easy to collect but can contain a significant number of biases. A few things need to be discussed before we use the survey data for analysis,

Choice of
Variables

Response
Options

Aspects which
have to be
considered
for Data from
Survey Studies

Response
Styles

Sample
Size

### 3.3.7   Choice of Variables

The below points summarize information about the choice of variables for analysis of survey data.

- Careful selection of segmentation variables is crucial in both commonsense and data-driven segmentation to ensure the quality of market segmentation solutions.

- In data-driven segmentation, relevant variables capturing the segmentation criterion should be included, while unnecessary variables must be avoided to prevent respondent fatigue and maintain response quality.

- Unnecessary variables increase the complexity of segmentation problems without adding relevant information, complicating the extraction of optimal market segments for analytical techniques.

- Noisy or masking variables, which do not contribute relevant information, can hinder algorithms from identifying correct segmentation solutions.

- Developing a high-quality questionnaire involves both exploratory qualitative research to understand consumer beliefs and quantitative survey research to ensure all critical variables are included without redundancy.

### 3.3.8 Response Options

The below points summarize information about the Response Options of data.

- Survey response options determine the data scale for segmentation analysis, crucial for distance-based techniques.

- Binary responses (e.g., yes/no) are ideal for segmentation, represented as 0s and 1s with clear distances between options.

- Nominal variables (e.g., occupation choices) can be converted to binary for segmentation, simplifying analysis.

- Metric data (e.g., age, nights stayed) allow for precise measurement and are well-suited for segmentation.

- Ordinal data (e.g., Likert scales) lack defined distances between responses, complicating distance-based segmentation methods.

- Visual analogue scales, like slider scales, provide metric data suitable for nuanced responses in online surveys.

### 3.3.9 Response Styles

The below points summarize information about the Response Styles of data.

- Response bias in surveys occurs when respondents consistently answer based on factors unrelated to the specific item content, known as response styles.

- Common response styles include using extreme options (e.g., STRONGLY AGREE/DISAGREE) or agreeing with all statements, affecting segmentation accuracy.

- Segmentation algorithms may misinterpret data due to response styles, leading to incorrect market segment conclusions.

- Minimizing response styles is critical in market segmentation to accurately identify and target genuine consumer segments.
- Additional analysis or exclusion of respondents affected by response styles is necessary to mitigate biased segment interpretations.

### 3.3.10 Sample Size

The below points summarize information about the Sample Size of data.

- Market segmentation analysis requires sufficient sample sizes to accurately identify segments; inadequate samples lead to ambiguous results.
- Recommended sample sizes vary: Formann suggests at least 2p (or five times 2p) where p is the number of segmentation variables, while Qiu and Joe (2015) propose $10 \cdot p \cdot k$ for accurate segment identification.
- Empirical studies by Dolnicar et al. (2014) and others show that increasing sample size improves segment extraction algorithm accuracy, especially beneficial for smaller initial samples.
- Challenges such as unequal segment sizes, overlap, and data quality issues (e.g., response biases) complicate segment recovery, necessitating sufficient sample sizes.
- Dolnicar et al. (2016) recommend a guideline of at least 100 respondents per segmentation variable to ensure robust segmentation outcomes, emphasizing high-quality data collection.
- The impact of sample size on segment recovery varies with data characteristics; while some challenges can be mitigated with larger samples, others, like high correlation between variables, remain problematic even with increased sample size.

### 3.3.11 Data from Internal sources & Experimental studies

- Organizations increasingly rely on internal data sources such as scanner data from grocery stores, booking data from airline loyalty programs, and online purchase data for market segmentation.
- Internal data strengths include its reflection of actual consumer behavior rather than self-reported data, which is prone to memory imperfections and response biases (Niemi 1993; Fisher 1993; Paulhus 1991; Dolnicar and Grün 2007a,b, 2009).
- Such data is automatically generated and easily accessible if stored in accessible formats, requiring minimal additional effort for collection and processing.
- However, internal data may be biased towards existing customers, potentially missing insights into different consumption patterns of future customers.

- Experimental data, derived from field or laboratory experiments, also informs market segmentation, such as studies on consumer responses to advertisements or choice experiments evaluating preferences based on product attributes.

- Conjoint analyses and choice experiments provide insights into how specific product attributes influence consumer choices, which can be used as segmentation criteria.

- Experimental data offers controlled settings for studying consumer behavior, although translating experimental findings to real-world market segments requires careful consideration.

## 3.4   Step 4 - Data Exploration

### 3.4.1   Loading and having a glance at the data

After data collection, exploratory data analysis (EDA) is crucial for cleaning and pre-processing the data and selecting appropriate segmentation algorithms. EDA helps to identify measurement levels of variables, investigate univariate distributions, and assess dependencies between variables. This pre-processing ensures the data is compatible with segmentation algorithms and guides the selection of the most suitable segmentation methods based on EDA results.

This analysis can be done on languages such as python or R. In the textbook, Market Segmentation Analysis they have used R to explore the data of the Australian Vacation dataset, which can be found over here. Through this they have extracted important information from the dataset. Some of the following R commands were used,

```
R> vaccsv <- system.file("csv/vacation.csv",+   package = "MSA")
R> file.copy(vaccsv, ".")
R> vac <- read.csv("vacation.csv", check.names = FALSE)
R> colnames(vac)
R> dim(vac)
R> summary(vac[, c(1, 2, 4, 5)])
```

implementing these lines of code in R will give you the desired outputs with the dataset. Please explore!

### 3.4.2   Data Cleaning

- The first step in data analysis is to clean the data, ensuring all values are recorded correctly and categorical variables have consistent labels.
- Check the range of plausible values for metric variables, such as age, which should lie between 0 and 110 years, to identify any data entry errors.
- Verify that categorical variables contain only permissible values; for example, gender should typically have only two values: female and male.
- In the Australian travel motives data set, the variables Gender and Age required no cleaning, but the Income2 variable categories were not sorted.
- This sorting issue is due to R's `read.csv()` or `read.table()` functions converting non-numeric columns into factors, which are sorted alphabetically by default.
- To re-order categories in R, copy the column to a helper variable (`inc2`), store its levels (`lev`), find the correct order, and convert it into an ordered factor.

- Copy the `Income2` column to a helper variable and check its levels:

```
R> inc2 <- vac$Income2
R> levels(inc2)
[1] "<30k"    ">120k"    "30-60k"
```

- Store the levels in a variable and reorder them correctly:

```
R> lev <- levels(inc2)
R> lev[c(1, 3, 4, 5, 2)]
[1] "<30k"    "30-60k"  "60-90k"  "90-120k" ">120k"
```

- Transform the variable into an ordered factor:

```
R> inc2 <- factor(inc2, levels = lev[c(1, 3, 4, 5, 2)], ordered = TRUE)
```

- Double-check the transformation by cross-tabulating the original and new variables:

```
R> table(orig = vac$Income2, new = inc2)
```

- Overwrite the original column with the correctly ordered version:

```
R> vac$Income2 <- inc2
```

- Ensure reproducibility by keeping all R code for data transformations and saving the cleaned data set using `save()` and `load()` functions.

### 3.4.3   Descriptive Analysis

- Being familiar with the data avoids misinterpretation of results from complex analyses.
- Descriptive numeric and graphic representations provide insights into the data.
- In R, the `summary()` command returns ranges, quartiles, and means for numeric variables, as well as frequency counts and missing values for categorical variables.
- Useful graphical methods for numeric data include histograms, boxplots, and scatter plots.
- Bar plots visualize frequency counts for categorical variables, while mosaic plots illustrate the association of multiple categorical variables.
- Histograms, created through binning, visualize the distribution of numeric variables by plotting the frequency of observations within specified value ranges.

### 3.4.3.1  Histograms

- Use the `lattice` package in R for creating histograms by segments:

    ```
    R> library("lattice")
    ```

- Construct a histogram for variable `Age` in the dataset `vac`:

    ```
    R> histogram(~ Age, data = vac)
    ```

- Specify the number of bins (`breaks`) to gain deeper insights:

    ```
    R> histogram(~ Age, data = vac, breaks = 50, type = "density")
    ```

- The finer bins provide more detailed information, revealing the distribution characteristics such as bi-modality.
- Use `type = "density"` to scale the y-axis by density estimates, allowing comparison with parametric distributions.
- This representation is typically preferred for histograms to visualize the probability density functions.

This is what the histogram looks like,



**Fig. 3.1:** Histogram from *Market Segmentation Analysis* textbook

**Fig. 3.2:** Box-and-Whisker plot of tourist age from *Market Segmentation Analysis* textbook

#### 3.4.3.2 Box and Whisker plot

- **Box-and-Whisker Plot (Boxplot):**
    - Summarizes data using the minimum, first quartile, median, third quartile, and maximum.
    - Provides insights into the distributional properties assuming unimodality.
    - Highlights outliers beyond 1.5 times the size of the box as circles to prevent loss of outlier information.

- **Use in R:**
    - In R, create a horizontal boxplot for variable `Age` with `boxplot(vac$Age, horizontal = TRUE, xlab = "Age")`.

- **Dot Chart for Percentages:**
    - The dot chart visualizes the percentages of agreement with travel motives from columns 13 to 32 in the dataset.
    - Each dot on the chart represents the percentage of respondents indicating that a specific travel motive was important to them on their last vacation.
    - To compute these percentages in R, we use `100 * colMeans(vac[, 13:32] == "yes")`. This calculates the mean percentage of "yes" responses across all columns.

– The dot chart is sorted to display these percentages from lowest to highest, providing a clear view of the distribution of agreement levels across different motives.

– It serves as a powerful tool to quickly grasp which travel motives are universally popular among respondents and which ones are less universally appealing.

– Insights gained from the dot chart highlight the heterogeneity in the importance attributed to various travel motives among survey participants.

– This heterogeneity underscores the potential of these motives as effective segmentation variables in market analysis, as they reveal distinct preferences and behaviors among different segments of the population.

- **Insights from Dot Chart:**
    – Illustrates varying agreement levels with different travel motives.
    – Confirms heterogeneity in importance attributed to travel motives, suggesting suitability as segmentation variables.



**Fig. 3.3:** Dot Chart of Yes percentages *Market Segmentation Analysis* textbook

So in short we have the following plotting devices to help us explain and portray the data set in a thorough way. I will be mentioning other plots as well, which have not been used in the Market Segmentation Analysis, but I have knowledge of them.

```
Box-and-Whisker Plot          Correlation Heat Map


Bar Plot ←  Graphical Plotting Techniques  → Histogram Plot


    Pair Plot                        Mosaic Plot


              Scatter Plot
```

### 3.4.4   Pre-Processing

#### 3.4.4.1   Categorical Variables

- Merge levels of categorical variables to reduce complexity.
- Convert categorical variables to numeric where scale assumptions hold.
- Consider binary options over multi-category scales for simplicity.
- Example: Convert survey responses to 0/1 matrix in R (`vacmot <- (vac[, 13:32] == "yes") + 0`).
- Preprocessing alters data for compatibility with statistical methods.
- Reflect on survey response options and cultural influences.
- Assume numeric variables with comparable scales for analysis methods.
- Use R packages like `flexclust` for handling categorical data effectively.

#### 3.4.4.2   Numerical Variables

- Standardize numerical variables to put them on a common scale:

$$z_i = \frac{x_i - \bar{x}}{s}$$

where $\bar{x}$ is the mean and $s$ is the standard deviation.

- Use R function `scale()` for standardization:

$$\texttt{vacmot.scaled <- scale(vacmot)}$$

- Consider robust methods for location and spread if data has outliers:

$$\text{Median: } \tilde{x} \quad \text{Interquartile Range (IQR): IQR} = Q3 - Q1$$

- Distance-based methods assume variables are numeric and comparable.

- Ordinal data can be treated as numeric if distances between scale points are assumed equal.

- Likert scales may not have equal distances between categories due to response biases.

- Convert dichotomous ordinal or nominal variables to binary (0/1) for simplicity:

$$\texttt{vacmot <- (vac[, 13:32] == "yes") + 0}$$

- Use appropriate statistical methods after conversion to ensure validity of results.

### 3.4.5  Principal Components Analysis

- PCA transforms multivariate data into uncorrelated principal components that explain variance in descending order.

- It works from the covariance or correlation matrix of numeric variables, sensitive to data scaling.

- Principal components are used for dimensionality reduction and visualizing high-dimensional data.

- Interpret PCA results using standard deviations, explained variances, and cumulative proportions.

- Visualization typically involves the first few principal components, e.g., PC1 and PC2 in scatter plots.

- PCA assists in identifying patterns and relationships among variables, aiding in market segmentation and data reduction tasks.

## 3.5 Step 5 - Selecting the Targeting Segments

We had talked about this somewhat briefly in the first few chapters, and now we revisit again to discuss in more detail

### 3.5.1 Target Decision

- **Critical Decision**
  Selecting specific market segments for targeting is a crucial, long-term decision affecting the organization's future performance.

- **Segmentation Process**
  After conducting a global market segmentation analysis, segments are available for detailed inspection and evaluation.

- **Segment Profiling and Description**
  Segments are profiled by inspecting key characteristics and described in detail to ensure they are identifiable, reachable, and their needs align with the organization's capabilities.

- **Double-Check Knock-Out Criteria**
  Ensure all segments meet essential criteria: sufficient size, homogeneity, distinctiveness, identifiability, reachability, and alignment with the organization's capabilities.

- **Evaluating Segment Attractiveness**
  Assess the attractiveness of each segment based on factors such as potential profitability, growth potential, and strategic fit with the organization.

- **Assessing Organizational Competitiveness** Evaluate the organization's competitiveness in serving each segment, considering factors like market position, resources, and ability to meet segment needs effectively.

- **Key Targeting Questions**
  Determine which segments the organization most wants to target and commit to, and which segments are most likely to choose the organization over competitors.

- **Basis for Decision**
  The answers to the above questions form the basis for making the final target segment selection.

### 3.5.2 Market Segment Evaluation

- **Decision Matrix Usage**
  Decision matrices help visualize segment attractiveness and organizational competitiveness, aiding in target market selection. Various versions exist, including the Boston matrix, GE/McKinsey matrix, and directional policy matrix.

- **Segment Attractiveness and Competitiveness**
  The x-axis represents segment attractiveness ("How attractive is the segment to us?") and the y-axis represents relative organizational competitiveness ("How attractive are we to the segment?").

- **Criteria Weighting and Rating**
  Criteria for segment attractiveness and organizational competitiveness are weighted and rated. For example, Criterion 1 (weight 25%) for Segment 1 is rated 5 for attractiveness and 2 for competitiveness.

- **Weighted Calculation**
  Total values for each segment are calculated by multiplying the weights with the ratings and summing them up. Segment 1's attractiveness is calculated as $0.25 \times 5 + 0.35 \times 2 + 0.20 \times 10 + 0.10 \times 8 + 0.10 \times 9 = 5.65$.

- **Bubble Size Representation**
  The size of the bubbles in the plot represents another criterion, such as profit potential. For example, Segment 1 has a bubble size of 2.25.

- **Segment Evaluation Plot**
  The plot helps identify the most promising segments. Segment 8 is highly attractive to the organization and vice versa, but has low profit potential (bubble size 1.50). Segment 5 is highly attractive with high profit potential but less organizational appeal.

- **Selecting Target Segments**
  Based on the plot, segments 3 and 7 might be excluded despite their high profit potential due to low attractiveness. Segment 5 is highly attractive but has low compatibility. Segment 8 is highly compatible but has low profit potential.

- **Practical Application**
  Organizations can use the decision matrix to prioritize segments that align well with their strategic goals, balancing attractiveness and competitiveness.

- **Balancing Act**
  The goal is to find segments that balance attractiveness and competitiveness. A highly attractive segment may not always be the best choice if the organization cannot compete effectively in that segment.

- **Profit Potential**
  Bubble size in the plot often represents profit potential, which combines segment size and spending behavior, crucial for selecting target segments.

- **Strategy Alignment**
  The selected segments should align with the organization's strategic goals and capabilities, ensuring that resources are effectively utilized.

- **Visual Decision Aid**
  The segment evaluation plot serves as a visual aid to facilitate discussion and decision-making within the segmentation team, making it easier to compare and contrast different segments.

- **Hypothetical Example**
  In the example, Segment 8 is highly attractive to both the organization and the segment itself, making it a good target despite lower profit potential. Segment 5, while highly attractive, may require more effort to improve organizational competitiveness.

**Table 3.1:** Data underlying the segment evaluation plot(Taken from the textbook)

| | Weight | Seg 1 | Seg 2 | Seg 3 | Seg 4 | Seg 5 | Seg 6 | Seg 7 | Seg 8 |
|---|---|---|---|---|---|---|---|---|---|
| **How attractive is the segment to us? (segment attractiveness)** | | | | | | | | | |
| **Criterion 1** | 25% | 5 | 10 | 1 | 5 | 10 | 3 | 1 | 10 |
| **Criterion 2** | 35% | 2 | 1 | 2 | 6 | 9 | 4 | 2 | 10 |
| **Criterion 3** | 20% | 10 | 6 | 4 | 4 | 8 | 2 | 1 | 9 |
| **Criterion 4** | 10% | 8 | 4 | 2 | 7 | 10 | 8 | 3 | 10 |
| **Criterion 5** | 10% | 9 | 6 | 1 | 4 | 7 | 9 | 7 | 8 |
| **Total** | 100% | 5.65 | 5.05 | 2.05 | 5.25 | 8.95 | 4.25 | 2.15 | 9.60 |
| **How attractive are we to the segment? (relative organisational competitiveness)** | | | | | | | | | |
| **Criterion 1** | 25% | 2 | 10 | 10 | 10 | 1 | 5 | 2 | 9 |
| **Criterion 2** | 25% | 3 | 10 | 4 | 6 | 2 | 4 | 3 | 8 |
| **Criterion 3** | 25% | 4 | 10 | 8 | 7 | 3 | 3 | 1 | 10 |
| **Criterion 4** | 15% | 9 | 8 | 3 | 9 | 4 | 5 | 3 | 9 |
| **Criterion 5** | 10% | 1 | 8 | 6 | 2 | 1 | 4 | 4 | 8 |
| **Total** | 100% | 3.70 | 9.50 | 6.55 | 7.30 | 2.20 | 4.15 | 2.35 | 8.90 |
| **Size** | | 2.25 | 5.25 | 6.00 | 3.75 | 5.25 | 2.25 | 4.50 | 1.50 |

## 3.6   Step 6 - Customising the Marketing Mix

**Implementations for Marketing Mix Decisions**

- **Evolution of Marketing**
  Initially viewed as a toolbox to assist in selling products, marketing combined various elements to achieve optimal sales results (Dolnicar and Ring 2014). Borden (1964) identified 12 marketing ingredients, but the common understanding now revolves around the 4Ps: Product, Price, Promotion, and Place.

- **Market Segmentation**
  It is integral to strategic marketing, closely linked with positioning and competition. The segmentation-targeting-positioning (STP) approach involves segment extraction, profiling, description, targeting, and positioning.
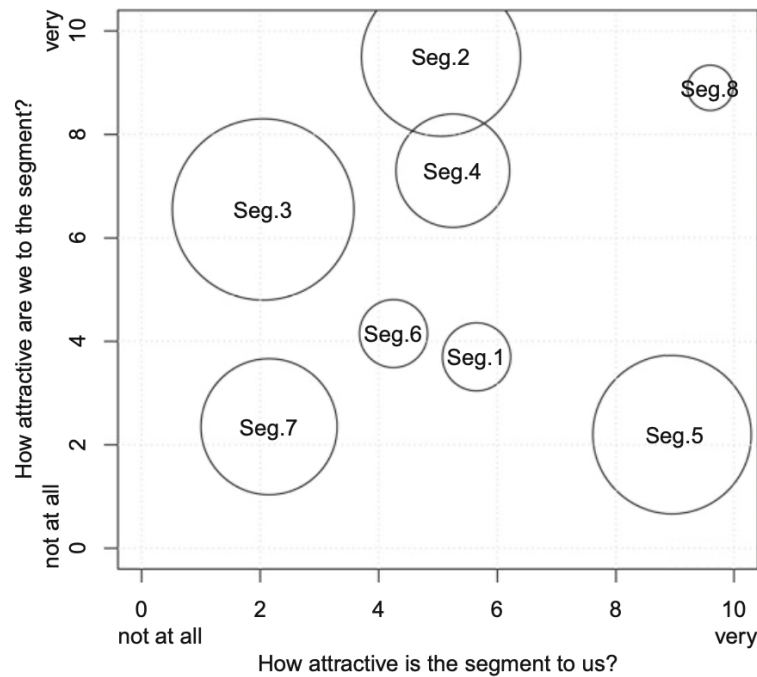
**Fig. 3.4:** Segment Evaluation Plot(Taken from the Textbook)

- **STP Approach**
  Ensures segmentation is integrated with other strategic decisions. The process may not always be linear, requiring adjustments between segmentation and targeting.

- **Impact on Marketing Mix**
  The selection of target segments influences the development of the marketing mix, traditionally comprising the 4Ps. Each aspect must be customized to align with the selected target segments.

- **Customizing Marketing Mix**
  To maximize benefits, the marketing mix should be tailored to the target segment. This may involve designing new products, adjusting prices, selecting appropriate distribution channels, and crafting targeted promotional strategies.

- **Segmentation Variables**
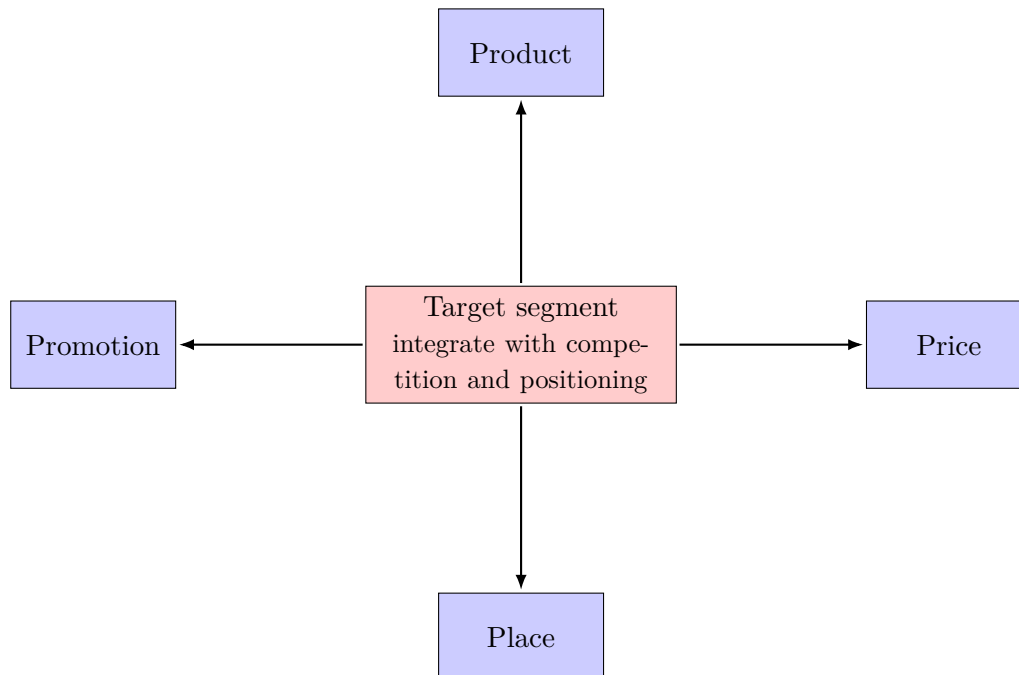  The focus can be structured around one of the 4Ps, influencing the choice of segmentation variables. For pricing, variables like price sensitivity are relevant; for advertising, lifestyle and psychographic variables are useful; for distribution, store loyalty and patronage are valuable.

- **Holistic Approach**
  Typically, segmentation analysis is not limited to one of the 4Ps. Instead,

insights from a detailed target segment description guide the organization in developing or adjusting the entire marketing mix to cater to the chosen segment.

```
                        ┌───────────┐
                        │  Product  │
                        └───────────┘
                              ▲
                              │
┌──────────┐    ┌────────────────────┐    ┌─────────┐
│Promotion │◄───│   Target segment   │───►│  Price  │
└──────────┘    │ integrate with com-│    └─────────┘
                │pe-tition and       │
                │positioning         │
                └────────────────────┘
                              │
                              ▼
                        ┌───────────┐
                        │   Place   │
                        └───────────┘
```

### 3.6.1   Product

- The product dimension of the marketing mix involves decisions related to specifying the product based on customer needs.
- This typically means modifying an existing product rather than designing a completely new one.
- Decisions include naming the product, packaging, offering warranties, and providing after-sales support services.
- An example of a product measure could be developing a "MUSEUMS, MONUMENTS  MUCH, MUCH MORE" product with an activities pass.
- Another example could be making gardens at the destination an attraction in their own right.

### 3.6.2   Price

- The price dimension of the marketing mix involves setting the price for a product and deciding on discounts to offer.
- It is important to understand the spending behavior of the target segment to make informed pricing decisions.

- Analysis of expenditure data can reveal if a segment has higher vacation expenditures per person per day.

- Higher expenditure segments may allow for premium pricing strategies rather than discounted prices.

- For instance, if a target segment shows higher spending, a premium price can be attached to a specialized product like "MUSEUMS, MONUMENTS MUCH, MUCH MORE".

### 3.6.3   Place

- The place dimension of the marketing mix involves decisions on how to distribute the product to customers.

- Key considerations include whether the product should be available online, offline, or both.

- Decisions need to be made on selling directly to customers or through wholesalers or retailers.

- Understanding the booking preferences of the target segment can help ensure that the product is available through preferred distribution channels.

- Visualizing booking behavior can provide insights into the preferred booking methods of the target segment.

### 3.6.4   Promotion

- The promotion dimension involves developing an advertising message and identifying effective communication methods.

- Tools in this category include public relations, personal selling, and sponsorship.

- It is important to determine the best information sources to reach the target segment.

- Comparing information sources and preferred TV stations of the target segment can help tailor the promotional strategy.

- Insights into preferred information sources can guide the design of specific information packs available both in tourist centers and online.

- Understanding TV channel preferences can help develop a media plan for maximum exposure to the target segment.

# Code

In this section, I have written the converted Python code for the Fast Food Case
Study from the R code given in the Market Segmentation Analysis textbook.

## A.1   Exploring the Data

The dataset contains responses from 1453 Australian adults about their perceptions
of McDonald's based on these attributes: YUMMY, CONVENIENT, SPICY, FAT-
TENING, GREASY, FAST, CHEAP, TASTY, EXPENSIVE, HEALTHY, and
DISGUSTING. Each attribute has a binary response: YES (McDonald's has this
attribute) or NO (McDonald's does not).

Below is the Python code for performing Principal Component Analysis (PCA) on
the McDonald's dataset:

```python
import pandas as pd
import numpy as np
from sklearn.cluster import KMeans
from sklearn.mixture import GaussianMixture
from sklearn.decomposition import PCA
from sklearn.preprocessing import StandardScaler
import matplotlib.pyplot as plt
import seaborn as sns
from scipy.stats import entropy
from patsy import dmatrices
from scipy.stats import chi2_contingency
import matplotlib.patches as mpatches
from statsmodels.graphics.mosaicplot import mosaic
from sklearn.tree import DecisionTreeClassifier
from sklearn.tree import plot_tree
from adjustText import adjust_text

mcdonalds_data =
    pd.read_csv('/Users/kanishka.arora/Downloads/mcdonalds.csv')

# Convert YES/NO to binary
MD_x = mcdonalds_data.iloc[:, :11].applymap(lambda x: 1 if x == 'Yes'
    else 0)

# Principal Components Analysis (PCA)
MD_pca = PCA()
```

```
25   MD_pca_fit = MD_pca.fit(MD_x)
26
27   # Transform data to principal components
28   MD_pca_proj = MD_pca.transform(MD_x)
29
30   # Plot PCA results with annotations
31   plt.figure(figsize=(12, 10))
32   plt.scatter(MD_pca_proj[:, 0], MD_pca_proj[:, 1], color='grey',
     ↪  alpha=0.6, s=100, edgecolor='k', label='Data Points')
33
34   # Collect all texts to adjust their positions
35   texts = []
36   for i, (pc1, pc2) in enumerate(zip(MD_pca.components_[0, :],
     ↪  MD_pca.components_[1, :])):
37       plt.arrow(0, 0, pc1, pc2, color='red', alpha=0.75, head_width=0.05,
         ↪  head_length=0.1)
38       text = plt.text(pc1 * 1.2, pc2 * 1.2, MD_x.columns[i], color='red',
         ↪  ha='center', va='center', fontsize=12, fontweight='bold')
39       texts.append(text)
40
41   # Adjust text positions to avoid overlap
42   adjust_text(texts, arrowprops=dict(arrowstyle='->', color='blue'))
43
44   # Customize plot appearance
45   plt.xlabel('Principal Component 1', fontsize=14, fontweight='bold')
46   plt.ylabel('Principal Component 2', fontsize=14, fontweight='bold')
47   plt.title('PCA of McDonald\'s Attributes', fontsize=16,
     ↪  fontweight='bold')
48   plt.grid(True, linestyle='--', alpha=0.7)
49   plt.axhline(0, color='black', linewidth=0.8)
50   plt.axvline(0, color='black', linewidth=0.8)
51   plt.legend(loc='upper right')
52   plt.tight_layout()
53
54   plt.show()
```
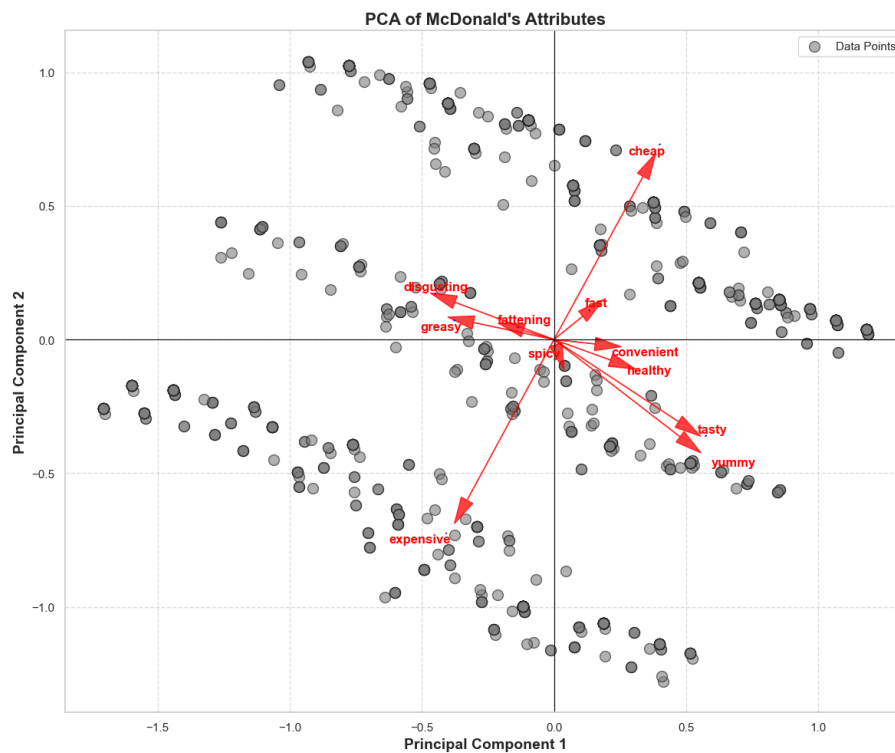
**Fig. A.1:** PCA of McDonald's Attributes

## A.2 Extracting Segments

### A.2.1 Using K-Means

```python
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
range_n_clusters = range(2, 9)
kmeans_models = {}
sum_of_distances = []
for n_clusters in range_n_clusters:
    kmeans = KMeans(n_clusters=n_clusters, n_init=10, random_state=1234)
    kmeans.fit(MD_x)
    kmeans_models[n_clusters] = kmeans
    sum_of_distances.append(kmeans.inertia_)
# Scree Plot
plt.figure(figsize=(12, 8))
bars = plt.bar(range(2, 9), sum_of_distances, color='skyblue',
    edgecolor='black')
for bar in bars:
    yval = bar.get_height()
```

```
16      plt.text(bar.get_x() + bar.get_width()/2, yval + 0.05 *
        ↪  max(sum_of_distances), round(yval, 2), ha='center', va='bottom',
        ↪  fontsize=10, fontweight='bold')
17  plt.plot(range(2, 9), sum_of_distances, color='red', marker='o',
    ↪  linestyle='-', linewidth=2, markersize=6)
18  plt.xlabel('Number of Segments', fontsize=14, fontweight='bold')
19  plt.ylabel('Sum of Distances within Segments', fontsize=14,
    ↪  fontweight='bold')
20  plt.title('Scree Plot for K-means Clustering', fontsize=16,
    ↪  fontweight='bold')
21  plt.grid(axis='y', linestyle='--', alpha=0.7)
22  plt.xticks(range(2, 9), fontsize=12)
23  plt.yticks(fontsize=12)
24  plt.tight_layout()
25  plt.show()
```
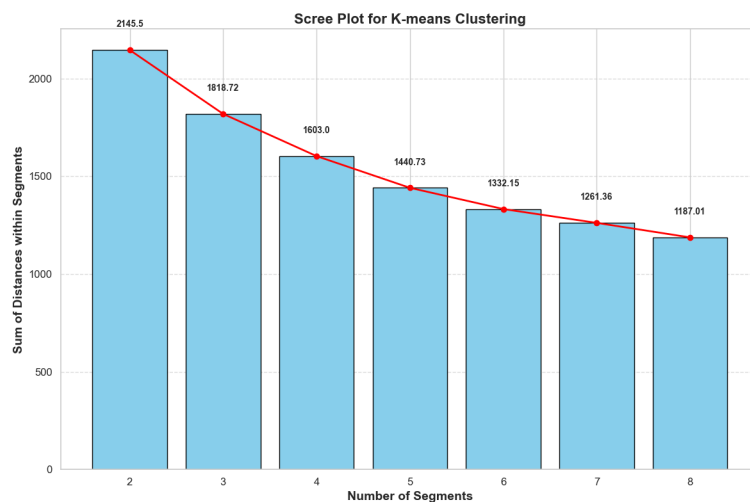


**Fig. A.2:** Scree plot for K-Means clustering.png

```
1  # Convert YES/NO to binary (1 for 'Yes', 0 for 'No')
2  MD_x = mcdonalds_data.iloc[:, :11].applymap(lambda x: 1 if x == 'Yes'
   ↪  else 0)
3
4  # Function to perform k-means clustering with random restarts
5  def kmeans_with_random_restarts(data, n_clusters, n_init=10):
6      kmeans = KMeans(n_clusters=n_clusters, n_init=n_init,
       ↪  random_state=1234)
7      kmeans.fit(data)
8      return kmeans
9
10 # Function to calculate adjusted Rand index for stability
11 def calculate_stability(data, kmeans_model, n_boot=100, n_init=10):
```

```
12      stability_scores = []
13      for _ in range(n_boot):
14          boot_sample = resample(data, n_samples=len(data), random_state=_)
15          kmeans_boot = kmeans_with_random_restarts(boot_sample,
        ↪  kmeans_model.n_clusters, n_init)
16          labels_true = kmeans_model.predict(data)
17          labels_boot = kmeans_boot.predict(data)
18          ari_score = adjusted_rand_score(labels_true, labels_boot)
19          stability_scores.append(ari_score)
20      return stability_scores
21
22  # Perform global stability analysis for each number of segments (2 to 8)
23  range_n_clusters = range(2, 9)
24  stability_results = {}
25  for n_clusters in range_n_clusters:
26      kmeans_model = kmeans_with_random_restarts(MD_x, n_clusters)
27      stability_scores = calculate_stability(MD_x, kmeans_model)
28      stability_results[n_clusters] = stability_scores
29
30  # BoxPlot
31  plt.figure(figsize=(10, 6))
32  boxprops = dict(linestyle='-', linewidth=2, color='black')
33  medianprops = dict(linestyle='-', linewidth=2, color='orange')
34  whiskerprops = dict(linestyle='-', linewidth=1.5, color='black')
35  capprops = dict(linestyle='-', linewidth=1.5, color='black')
36  plt.boxplot(stability_results.values(), notch=False, patch_artist=False,
        ↪  boxprops=boxprops, medianprops=medianprops,
        ↪  whiskerprops=whiskerprops, capprops=capprops)
37  plt.xticks(range(1, len(range_n_clusters) + 1), range_n_clusters,
        ↪  fontsize=12)
38  plt.xlabel('Number of Segments', fontsize=12)
39  plt.ylabel('Adjusted Rand Index', fontsize=12)
40  plt.title('Global Stability Boxplot', fontsize=14)
41  plt.ylim(0.4, 1.0)
42  plt.grid(axis='y', linestyle='--', linewidth=0.5)
43  plt.show()
```
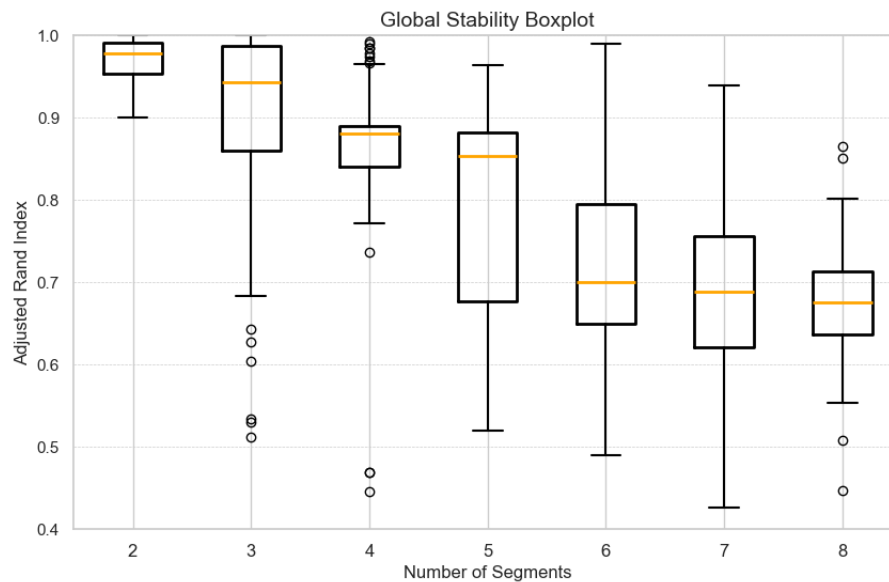
**Fig. A.3:** Global Stability Boxplot

```python
# One-hot encode categorical variables
MD_x_numeric = pd.get_dummies(mcdonalds_data, drop_first=True)
scaler = StandardScaler()
MD_x_numeric_scaled = scaler.fit_transform(MD_x_numeric)
n_clusters = 4
kmeans_model = KMeans(n_clusters=n_clusters, n_init=10,
↪    random_state=1234)
kmeans_model.fit(MD_x_numeric_scaled)
labels = kmeans_model.labels_

# Calculate pairwise similarity within each cluster
def calculate_similarity_within_cluster(data, labels, cluster_label):
    cluster_data = data[labels == cluster_label]
    similarity = []
    for i in range(len(cluster_data)):
        for j in range(i + 1, len(cluster_data)):
            similarity.append(np.linalg.norm(cluster_data[i] -
            ↪    cluster_data[j]))
    return similarity
# Plot histograms for each cluster
fig, axes = plt.subplots(2, 2, figsize=(12, 10))
axes = axes.flatten()
fig.suptitle('Gorge plot of the four-segment k-means solution for the
↪    fast food dataset', fontsize=16)
for cluster_label in range(n_clusters):
    similarity = calculate_similarity_within_cluster(MD_x_numeric_scaled,
    ↪    labels, cluster_label)
```
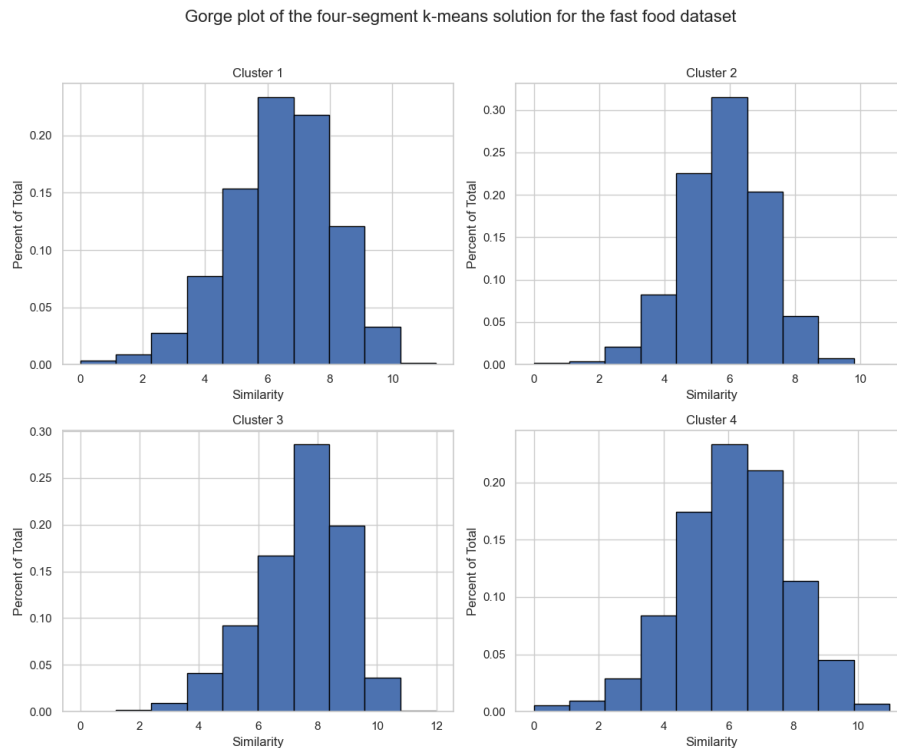
```
24      axes[cluster_label].hist(similarity, bins=10, edgecolor='black',
        ↪   density=True)
25      axes[cluster_label].set_title(f'Cluster {cluster_label + 1}')
26      axes[cluster_label].set_xlabel('Similarity')
27      axes[cluster_label].set_ylabel('Percent of Total')
28  plt.tight_layout(rect=[0, 0, 1, 0.96])  # Adjust layout to make room for
    ↪   the suptitle
29  plt.show()
```



**Fig. A.4:** Gorge plot of the four-segment k-means solution for the fast food dataset

```
1   MD_x_numeric = pd.get_dummies(mcdonalds_data, drop_first=True)
2
3   # Standardize the data
4   scaler = StandardScaler()
5   MD_x_numeric_scaled = scaler.fit_transform(MD_x_numeric)
6   num_segments = range(2, 9)
7
8   # Fit KMeans models for each segment
9   kmeans_models = {}
10  for segment in num_segments:
11      kmeans_models[segment] = KMeans(n_clusters=segment, n_init=10,
        ↪   random_state=1234)
12      kmeans_models[segment].fit(MD_x_numeric_scaled)
```

```
13
14  # Calculate segment stability
15  segment_stability = []
16  for segment in num_segments:
17      labels_segment = kmeans_models[segment].labels_
18      segment_stability.append(labels_segment)
19
20  # SLSA Plot
21  plt.figure(figsize=(12, 8))
22  for i, segment in enumerate(num_segments):
23      stability_scores = [np.mean(segment_stability[i] == labels) for
        ↪   labels in segment_stability]
24      plt.plot(num_segments, stability_scores, marker='o', label=f'Segment
        ↪   {segment}')
25  plt.xlabel('Number of Segments', fontsize=14, fontweight='bold')
26  plt.ylabel('Segment Level Stability', fontsize=14, fontweight='bold')
27  plt.title('Segment Level Stability Across Solutions (SLSA) Plot',
    ↪   fontsize=16, fontweight='bold')
28  plt.xticks(num_segments, fontsize=12)
29  plt.yticks(fontsize=12)
30  plt.legend(fontsize=12)
31  plt.grid(True, linestyle='--', alpha=0.7)
32  plt.tight_layout()
33  plt.show()
```
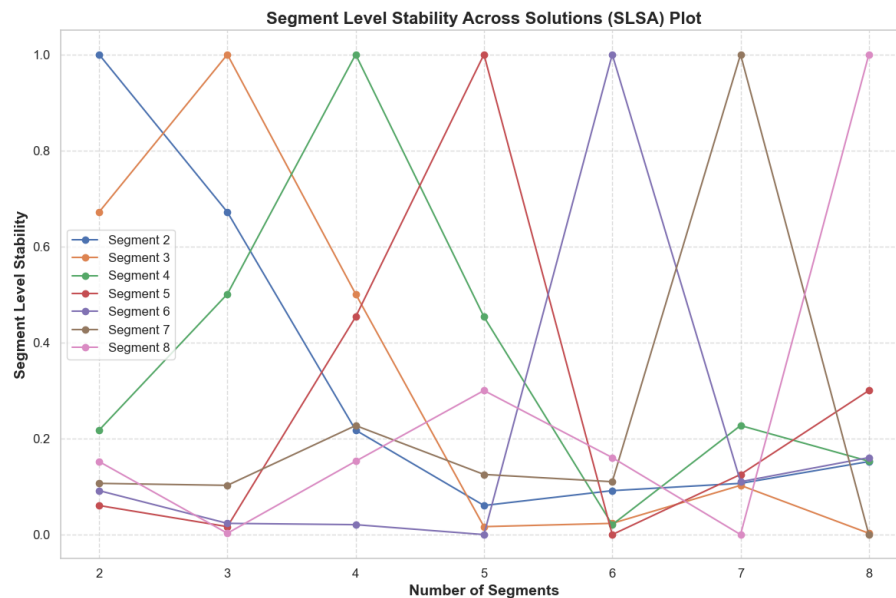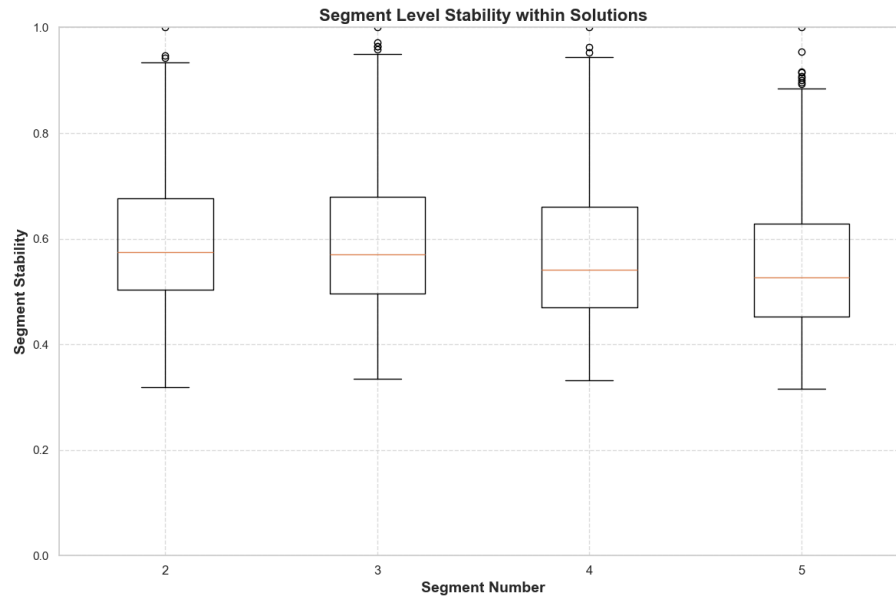


**Fig. A.5:** Segment Level Stability Across Solutions (SLSA) Plot

```
1  MD_x_numeric = pd.get_dummies(mcdonalds_data, drop_first=True)
2  scaler = StandardScaler()
```

```python
MD_x_numeric_scaled = scaler.fit_transform(MD_x_numeric)
segment_solutions = ["2", "3", "4", "5"]

# Fit KMeans models for each segment and store labels and similarities
segment_labels = {}
segment_similarities = {}
for segment in segment_solutions:
    kmeans = KMeans(n_clusters=int(segment), n_init=10,
    ↪   random_state=1234)
    kmeans.fit(MD_x_numeric_scaled)
    segment_labels[segment] = kmeans.labels_
    segment_similarities[segment] =
    ↪   kmeans.transform(MD_x_numeric_scaled).min(axis=1)

# Calculate segment stability values
segment_stability_values = []
for segment in segment_solutions:
    similarities = segment_similarities[segment]
    normalized_similarities = similarities / np.max(similarities)
    segment_stability_values.append(normalized_similarities)

# Plot Segment Level Stability within Solutions
plt.figure(figsize=(12, 8))
plt.boxplot(segment_stability_values, whis=1.5)
plt.xlabel("Segment Number", fontsize=14, fontweight='bold')
plt.ylabel("Segment Stability", fontsize=14, fontweight='bold')
plt.xticks(range(1, len(segment_solutions) + 1), segment_solutions,
↪   fontsize=12)
plt.ylim(0, 1)
plt.title("Segment Level Stability within Solutions", fontsize=16,
↪   fontweight='bold')
plt.grid(True, linestyle='--', alpha=0.7)
plt.tight_layout()
plt.show()
```

**Fig. A.6:** Segment Level Stability within solutions

## A.2.2  Using mixtures of distributions

```python
# Convert non-numeric columns to numeric using one-hot encoding
MD_x_numeric = pd.get_dummies(MD_x)
np.random.seed(1234)
k_values = range(2, 9)
MD_m28 = []

for k in k_values:
    model = KMeans(n_clusters=k, random_state=1234)
    model.fit(MD_x_numeric.values)
    iter_val = model.n_iter_
    converged = model.n_iter_ < model.max_iter
    log_likelihood = -model.inertia_
    n_samples, _ = MD_x_numeric.shape
    aic = -2 * log_likelihood + 2 * k
    bic = -2 * log_likelihood + np.log(n_samples) * k
    labels = model.labels_
    counts = np.bincount(labels)
    probs = counts / float(counts.sum())
    class_entropy = entropy(probs)
    icl = bic - class_entropy

    MD_m28.append((iter_val, converged, k, k, log_likelihood, aic, bic,
        icl))

```

```
24   MD_m28 = pd.DataFrame(MD_m28, columns=['iter', 'converged', 'k', 'k0',
     ↪  'logLik', 'AIC', 'BIC', 'ICL'])
25
26   print(MD_m28)
27   num_segments = MD_m28["k"]
28   AIC_values = MD_m28["AIC"]
29   BIC_values = MD_m28["BIC"]
30   ICL_values = MD_m28["ICL"]
31
32   # Set Seaborn style for better aesthetics
33   sns.set(style="whitegrid")
34   plt.figure(figsize=(12, 8))
35   sns.lineplot(x=num_segments, y=AIC_values, marker='o', label='AIC (Akaike
     ↪  Information Criterion)', linestyle='-', color='b')
36   sns.lineplot(x=num_segments, y=BIC_values, marker='o', label='BIC
     ↪  (Bayesian Information Criterion)', linestyle='-', color='r')
37   sns.lineplot(x=num_segments, y=ICL_values, marker='o', label='ICL
     ↪  (Integrated Complete Likelihood)', linestyle='-', color='g')
38   plt.xlabel('Number of Segments', fontsize=14)
39   plt.ylabel('Value of Information Criteria', fontsize=14)
40   plt.title('Information Criteria (AIC, BIC, ICL)', fontsize=16,
     ↪  fontweight='bold')
41   plt.legend(title='Criteria', fontsize=12, title_fontsize='13')
42   plt.xticks(num_segments)
43   plt.grid(True)
44   plt.tight_layout()
45   plt.show()
```

**Fig. A.7:** Information Criteria (AIC, BIC, ICL)

```python
# Perform KMeans clustering
k_values = range(2, 9)
kmeans_models = {}
np.random.seed(1234)
for k in k_values:
    model = KMeans(n_clusters=k, random_state=1234)
    model.fit(MD_x_numeric)
    kmeans_models[k] = model

# Select the four-segment solution from k-means
MD_k4 = kmeans_models[4]

# Create a function to generate a mixture model
def fit_gmm(X, n_components, random_state=None):
    gmm = GaussianMixture(n_components=n_components,
    ↪  random_state=random_state)
    gmm.fit(X)
    return gmm

# Fit a Gaussian Mixture Model (GMM) with 4 components
MD_m4 = fit_gmm(MD_x_numeric, 4, random_state=1234)

# Cross-tabulation of k-means and GMM clusters
kmeans_labels = MD_k4.labels_
gmm_labels = MD_m4.predict(MD_x_numeric)
cross_tab = pd.crosstab(kmeans_labels, gmm_labels, rownames=['kmeans'],
↪  colnames=['mixture'])
```

```
26  print(cross_tab)
27
28  # Initialize GMM with k-means labels
29  MD_m4a = GaussianMixture(n_components=4, random_state=1234)
30  MD_m4a.fit(MD_x_numeric.values, y=kmeans_labels)
31  gmm_labels_a = MD_m4a.predict(MD_x_numeric)
32
33  # Cross-tabulation of k-means and GMM clusters with initialized labels
34  cross_tab_a = pd.crosstab(kmeans_labels, gmm_labels_a,
    ↪  rownames=['kmeans'], colnames=['mixture'])
35  print(cross_tab_a)
36
37  # Comparing the log-likelihood values
38  logLik_m4 = MD_m4.score(MD_x_numeric) * len(MD_x_numeric)
39  logLik_m4a = MD_m4a.score(MD_x_numeric) * len(MD_x_numeric)
40  print(f'Log Likelihood (random init): {logLik_m4}')
41  print(f'Log Likelihood (k-means init): {logLik_m4a}')
42
```

```
mixture    0    1    2    3
kmeans
0        508    0    4   37
1          0  215   11    4
2         44    3  267    8
3          0   36   15  301
mixture    0    1    2    3
kmeans
0        508    0    4   37
1          0  215   11    4
2         44    3  267    8
3          0   36   15  301
Log Likelihood (random init): 16082.886182864011
Log Likelihood (k-means init): 16082.886182864011
```

**Fig. A.8:** Output

## A.2.3   Using Regression Models

```
1  like_counts = pd.value_counts(mcdonalds_data['Like'])
2  reversed_counts = like_counts.iloc[::-1]
3  print(reversed_counts)
4
```

```
5    # Define a mapping of string values to numeric codes
6    like_mapping = {
7        'I HATE IT!-5': -5,
8        '-4': -4,
9        '-3': -3,
10       '-2': -2,
11       '-1': -1,
12       '0': 0,
13       '1': 1,
14       '2': 2,
15       '3': 3,
16       '4': 4,
17       'I LOVE IT!+5': 5
18   }
19
20   # Map 'Like' to numeric values using the defined mapping
21   mcdonalds_data['Like.n'] = mcdonalds_data['Like'].map(like_mapping)
22
23   # Count the occurrences of each numeric value in 'Like.n'
24   like_n_counts = mcdonalds_data['Like.n'].value_counts()
25   print(like_n_counts)
```

```
Like
-1              58
-2              59
-4              71
-3              73
I love it!+5    143
I hate it!-5    152
+1              152
+4              160
0               169
+2              187
+3              229
Name: count, dtype: int64
Like.n
 0.0     169
-3.0      73
-4.0      71
-2.0      59
-1.0      58
Name: count, dtype: int64
```

**Fig. A.9:** Output

```
1    # Applying the mapping to create 'Like_n' column
2    mcdonalds_data['Like_n'] = mcdonalds_data['Like'].map(like_mapping)
```

```
3
4   # Handle missing values (if any) - remove rows with missing values
5   mcdonalds_data.dropna(subset=['Like_n'] +
    ↪  mcdonalds_data.columns[1:-1].tolist(), inplace=True)
6
7   print(mcdonalds_data.head())
```

|    | yummy | convenient | spicy | fattening | greasy | fast | cheap | tasty | expensive | healthy | \ |
|----|-------|------------|-------|-----------|--------|------|-------|-------|-----------|---------|---|
| 0  | No    | Yes        | No    | Yes       | No     | Yes  | Yes   | No    | Yes       | No      |   |
| 10 | No    | Yes        | No    | Yes       | No     | Yes  | Yes   | No    | No        | No      |   |
| 12 | No    | Yes        | No    | Yes       | No     | Yes  | Yes   | No    | No        | No      |   |
| 14 | No    | Yes        | No    | Yes       | No     | Yes  | No    | No    | Yes       | No      |   |
| 16 | Yes   | Yes        | No    | Yes       | Yes    | Yes  | Yes   | Yes   | No        | No      |   |

|    | disgusting | Like | Age | VisitFrequency     | Gender | Like_n | Cluster |
|----|------------|------|-----|--------------------|--------|--------|---------|
| 0  | No         | -3   | 61  | Every three months | Female | -3     | 2       |
| 10 | Yes        | -2   | 53  | Every three months | Female | -2     | 2       |
| 12 | No         | 0    | 65  | Every three months | Male   | 0      | 2       |
| 14 | No         | -3   | 67  | Once a month       | Male   | -3     | 2       |
| 16 | No         | 0    | 34  | Once a month       | Female | 0      | 2       |

**Fig. A.10:** Output

```
1   import seaborn as sns
2   from patsy import dmatrices
3   from sklearn.mixture import GaussianMixture
4   import statsmodels.api as sm
5
6   # Cleaning column names
7   print("Original column names:", mcdonalds_data.columns)
8   mcdonalds_data.columns = [col.replace(' ', '_').replace('.', '_') for col
    ↪  in mcdonalds_data.columns]
9   mcdonalds_data = mcdonalds_data.loc[:,
    ↪  ~mcdonalds_data.columns.duplicated()]
10  print("Updated column names:", mcdonalds_data.columns)
11  independent_vars = mcdonalds_data.columns[1:-1].tolist()
12  dependent_var = 'Like_n'
13  formula_str = dependent_var + ' ~ ' + ' + '.join(independent_vars)
14
15  y, X = dmatrices(formula_str, data=mcdonalds_data,
    ↪  return_type='dataframe')
16
17  # Fit Gaussian Mixture Model
18  np.random.seed(1234)
19  n_components = 2  # Adjust number of components if needed
```

```python
20  model = GaussianMixture(n_components=n_components, n_init=10,
    ↪   random_state=0)
21  model.fit(X)
22
23  # Predict the cluster labels
24  labels = model.predict(X)
25  mcdonalds_data['Cluster'] = labels
26
27  # Fit linear regression models for each cluster and store coefficients
    ↪   and standard errors
28  coefficients = []
29  standard_errors = []
30  for cluster in range(n_components):
31      cluster_data = mcdonalds_data[mcdonalds_data['Cluster'] == cluster]
32      y_cluster, X_cluster = dmatrices(formula_str, data=cluster_data,
        ↪   return_type='dataframe')
33      model = sm.OLS(y_cluster, X_cluster).fit()
34      coefficients.append(model.params)
35      standard_errors.append(model.bse)
36
37  # Convert coefficients to DataFrame for plotting
38  coefficients_df = pd.DataFrame(coefficients).T
39  coefficients_df.columns = [f'Comp. {i + 1}' for i in range(n_components)]
40  coefficients_df = coefficients_df.reset_index()
41
42  # Convert standard errors to DataFrame for error bars
43  errors_df = pd.DataFrame(standard_errors).T
44  errors_df.columns = [f'Comp. {i + 1}' for i in range(n_components)]
45  errors_df = errors_df.reset_index()
46
47  # Plotting the coefficients with error bars
48  fig, ax = plt.subplots(figsize=(10, 8))
49
50  # Plot each component's coefficients with error bars
51  colors = sns.color_palette("Paired", n_components)
52  for i in range(n_components):
53      ax.barh(coefficients_df['index'], coefficients_df[f'Comp. {i + 1}'],
        ↪   xerr=errors_df[f'Comp. {i + 1}'],
54              label=f'Comp. {i + 1}', alpha=0.6, color=colors[i],
                ↪   edgecolor='black')
55
56  # Beautifying the plot
57  ax.set_xlabel('Coefficient Value', fontsize=14)
58  ax.set_ylabel('Variable', fontsize=14)
59  ax.set_title('Regression Coefficients for Each Component', fontsize=16,
    ↪   fontweight='bold')
60  ax.legend(title='Components', fontsize=12)
61  ax.grid(True, which='both', linestyle='--', linewidth=0.5)
62  ax.axvline(x=0, color='grey', linestyle='--', linewidth=0.8)
```

```
63  plt.xticks(fontsize=12)
64  plt.yticks(fontsize=12)
65
66  plt.show()
```
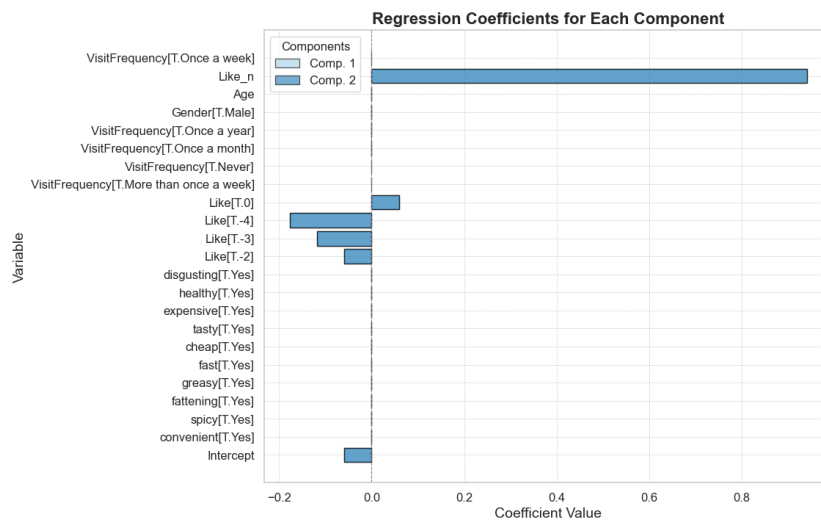


**Fig. A.11:** Regression coefficients for each component

## A.3 Profiling Segments

```python
1  from scipy.cluster.hierarchy import linkage, leaves_list
2
3
4  # Setting seed for reproducibility
5  np.random.seed(42)
6
7  # Define attributes and create simulated data
8  attributes = ['disgusting', 'expensive', 'greasy', 'healthy', 'spicy',
   ↪ 'fast', 'convenient', 'fattening', 'cheap', 'tasty', 'yummy']
9  data = np.random.rand(100, len(attributes))  # Replace with your actual
   ↪ data
10 df = pd.DataFrame(data, columns=attributes)
11
12 # Simulated k-means results (replace with actual segment percentages)
13 segment_data = {
14     'Cluster 1': np.random.rand(len(attributes)),
15     'Cluster 2': np.random.rand(len(attributes)),
16     'Cluster 3': np.random.rand(len(attributes)),
17     'Cluster 4': np.random.rand(len(attributes))
18 }
```

```python
MD_k4 = pd.DataFrame(segment_data, index=attributes)
linked = linkage(df.T, method='ward')
order = leaves_list(linked)
ordered_attributes = [attributes[i] for i in order]

ordered_MD_k4 = MD_k4.loc[ordered_attributes]
fig, axs = plt.subplots(2, 2, figsize=(14, 10), sharex=True, sharey=True)
axs = axs.flatten()

# Sample size and percentage of each cluster (replace with actual
    numbers)
cluster_info = {'Cluster 1': (470, 32), 'Cluster 2': (257, 18), 'Cluster
    3': (324, 22), 'Cluster 4': (402, 28)}
colors = ['blue', 'green', 'red', 'pink']  # Different colors for each
    segment

for i, (cluster, (size, percent)) in enumerate(cluster_info.items()):
    bars = ordered_MD_k4[cluster]
    marker_threshold = 0.25  # Highlighting threshold (change according
        to the data)
    ax = axs[i]

    # Plot the bars and highlight marker variables
    for j, (attribute, value) in enumerate(bars.items()):
        ax.barh(j, value, color=colors[i], edgecolor='none', alpha=0.6)
        # Highlight the marker variables
        if abs(value - ordered_MD_k4.mean(axis=1)[attribute]) >
            marker_threshold:
            ax.plot(value, j, 'ro')

    # Add horizontal lines indicating the overall sample percentage
    sample_means = df.mean(axis=0)
    for k, mean in enumerate(sample_means[ordered_attributes]):
        ax.plot([mean, mean], [k - 0.4, k + 0.4], 'k--', lw=0.7)
    ax.set_yticks(range(len(ordered_attributes)))
    ax.set_yticklabels(ordered_attributes)
    ax.set_xlim(0, 1)
    ax.set_title(f"{cluster}: {size} ({percent}%)")

plt.subplots_adjust(hspace=0.3, wspace=0.3)
fig.suptitle('Segment Profile Plot', fontsize=16)
plt.show()
```

**Fig. A.12:** Segment Profile Plot

```
1  from scipy.cluster.hierarchy import linkage, leaves_list
2  np.random.seed(42)
3  attributes = ['disgusting', 'expensive', 'greasy', 'healthy', 'spicy',
   ↪  'fast', 'convenient', 'fattening', 'cheap', 'tasty', 'yummy']
4  data = np.random.rand(100, len(attributes))  # Replace with your actual
   ↪  data
5  segments = np.random.choice([1, 2, 3, 4], size=100)  # Simulated segment
   ↪  membership
6
7  # Convert to DataFrame for easier manipulation
8  df = pd.DataFrame(data, columns=attributes)
9  df['Segment'] = segments
10
11 # Perform PCA
12 pca = PCA(n_components=2)
13 pca_result = pca.fit_transform(df[attributes])
14
15 # Create a DataFrame with PCA results
16 df_pca = pd.DataFrame(pca_result, columns=['principal component 1',
   ↪  'principal component 2'])
17 df_pca['Segment'] = df['Segment']
18
19 # Plotting
20 fig, ax = plt.subplots(figsize=(12, 10))
21 colors = {1: 'tomato', 2: 'mediumseagreen', 3: 'royalblue', 4: 'orchid'}
```

```
22   markers = {1: 'o', 2: '^', 3: 's', 4: 'x'}
23   for segment in sorted(df_pca['Segment'].unique()):
24       subset = df_pca[df_pca['Segment'] == segment]
25       ax.scatter(subset['principal component 1'], subset['principal
         ↪   component 2'],
26                   c=colors[segment], label=f'Segment {segment}',
                     ↪   marker=markers[segment], alpha=0.7, s=60)
27
28   for segment in sorted(df_pca['Segment'].unique()):
29       center = df_pca[df_pca['Segment'] == segment][['principal component
         ↪   1', 'principal component 2']].mean()
30       ax.scatter(center['principal component 1'], center['principal
         ↪   component 2'],
31                   c='black', edgecolor='white', s=100, zorder=5)  # Reduced
                     ↪   size
32       # Label as "1", "2", etc. slightly offset from the center
33       ax.text(center['principal component 1'], center['principal component
         ↪   2'],
34               str(segment), color='white', ha='center', va='center',
                 ↪   fontsize=12, weight='bold', zorder=6)
35
36
37   loadings = pca.components_.T * np.sqrt(pca.explained_variance_)
38   texts = []
39   for i, attr in enumerate(attributes):
40       ax.arrow(0, 0, loadings[i, 0], loadings[i, 1], color='red',
         ↪   alpha=0.7, head_width=0.02, head_length=0.03)
41       # Adjust label positions to avoid overlap
42       texts.append(ax.text(loadings[i, 0] * 1.3, loadings[i, 1] * 1.3,
         ↪   attr, color='red', ha='center', va='center', fontsize=10))
43
44   # Use adjust_text to prevent overlap
45   adjust_text(texts, ax=ax, only_move={'text': 'xy'},
     ↪   arrowprops=dict(arrowstyle="-", color='red', alpha=0.7))
46   ax.set_xlabel('Principal Component 1', fontsize=14)
47   ax.set_ylabel('Principal Component 2', fontsize=14)
48   ax.legend(loc='upper right', fontsize=12, markerscale=1.2)  # Adjusted
     ↪   marker scale
49   ax.grid(True)
50   plt.title('Segment Separation Plot using components 1 and 2',
     ↪   fontsize=16)
51   plt.xticks(fontsize=12)
52   plt.yticks(fontsize=12)
53   plt.tight_layout()
54
55   # Save the plot for inclusion in LaTeX
56   plt.savefig('Segment_Separation_Plot.png')
57   plt.close()
```
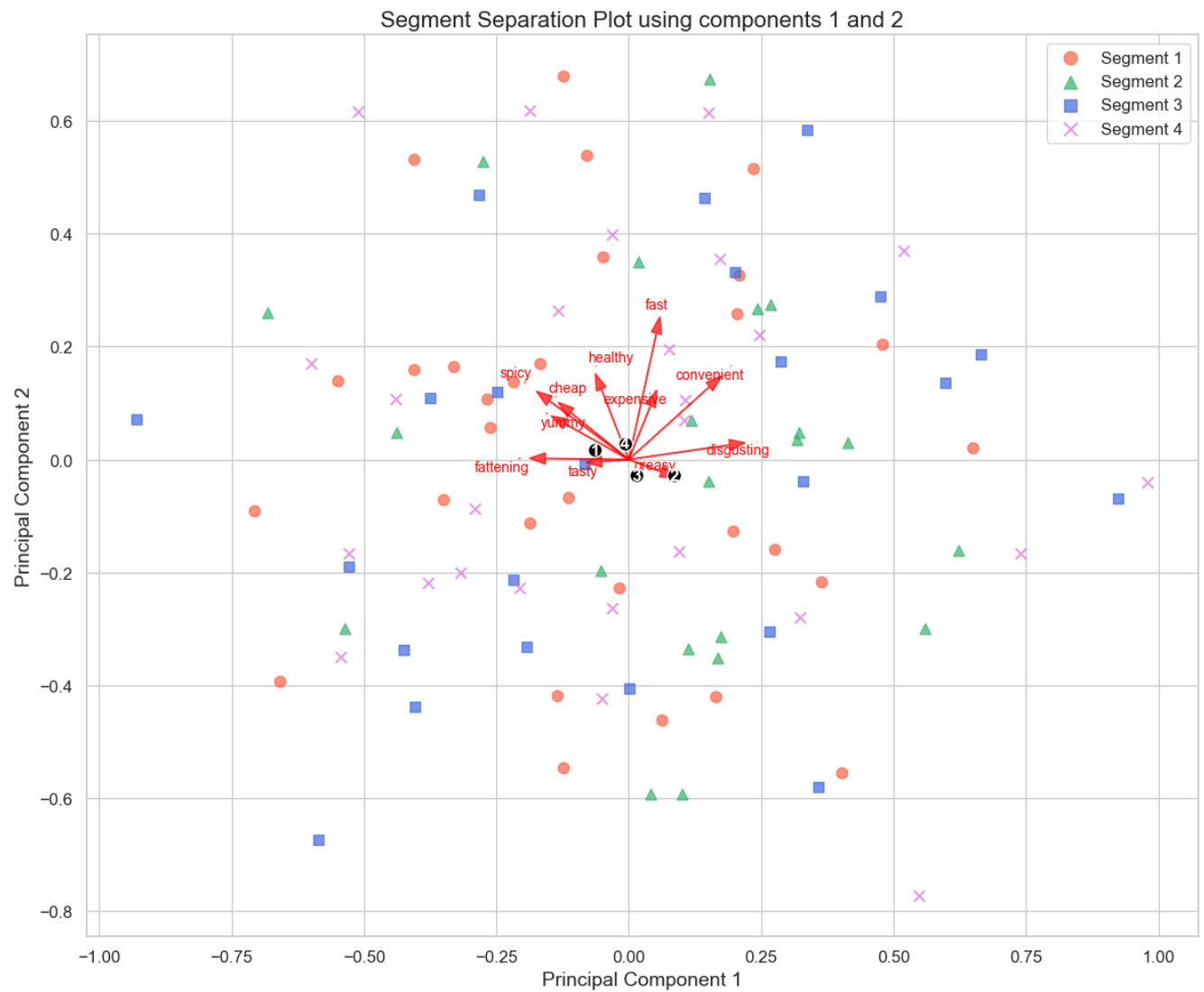
**Fig. A.13:** Segment Separation Plot

## A.4 Describing Segments

```python
from statsmodels.graphics.mosaicplot import mosaic
from scipy.stats import chi2_contingency
import matplotlib.patches as mpatches

# Sample data for demonstration (replace with actual data)
np.random.seed(42)
segments = np.random.choice([1, 2, 3, 4], size=100)  # Simulated segment
    membership
```

```
8    like_hate = np.random.choice(['I hate it!', '-5', '-4', '-3', '-2', '-1',
     ↪    '0', '+1', '+2', '+3', '+4', 'I love it! +5'], size=100)  # Simulated
     ↪    like/hate variable
9
10   # Convert to DataFrame for easier manipulation
11   df = pd.DataFrame({'Segment': segments, 'Like': like_hate})
12
13   # Create a contingency table
14   contingency_table = pd.crosstab(df['Segment'], df['Like'])
15
16   # Perform chi-squared test to get expected frequencies
17   chi2, p, dof, expected = chi2_contingency(contingency_table)
18
19   # Calculate standardized residuals
20   residuals = (contingency_table - expected) / np.sqrt(expected)
21
22   # Define color mapping for standardized residuals
23   def residual_color(value):
24       if value < -4:
25           return '#D73027'
26       elif value < -2:
27           return '#FC8D59'
28       elif value < 0:
29           return '#FEE08B'
30       elif value < 2:
31       return '#D9EF8B'
32       elif value < 4:
33           return '#91CF60'
34       else:
35           return '#1A9850'
36
37   # Function to format the properties of each cell
38   def props(key):
39       segment, like_hate = key
40       value = residuals.loc[int(segment), like_hate]
41       return {'color': residual_color(value)}
42
43   # Generate the mosaic plot with shading for standardized residuals
44   plt.figure(figsize=(24, 18))  # Increase the figure size for better
     ↪    readability
45   mosaic(contingency_table.stack(), gap=0.01, properties=props)
46   plt.xlabel('Segment Number', fontsize=18, labelpad=20)
47   plt.ylabel('Like/Hate McDonald\'s', fontsize=18, labelpad=20)
48   plt.title('Shaded Mosaic Plot for Segment Membership and I LIKE IT for
     ↪    the Fast Food Data Set', fontsize=22, pad=30)
49   plt.xticks(rotation=45, ha='right', fontsize=12)
50   plt.yticks(rotation=0, va='center', fontsize=12)
51
52   # Adding a color legend
```

```
53  legend_labels = ['<-4', '-4:-2', '-2:0', '0:2', '2:4', '>4']
54  colors = ['#D73027', '#FC8D59', '#FEE08B', '#D9EF8B', '#91CF60',
    ↪  '#1A9850']
55  patches = [mpatches.Patch(color=colors[i], label=legend_labels[i]) for i
    ↪  in range(len(legend_labels))]
56  plt.legend(handles=patches, title="Standardized Residuals", loc="center
    ↪  left", bbox_to_anchor=(1, 0.5), fontsize=14, title_fontsize=16)
57  plt.grid(True, which='both', linestyle='--', linewidth=0.5)
58  plt.gca().set_facecolor('#f0f0f0')
59
60  plt.show()
```

## Shaded Mosaic Plot for Segment Membership and I LIKE IT for the Fast Food Data Set
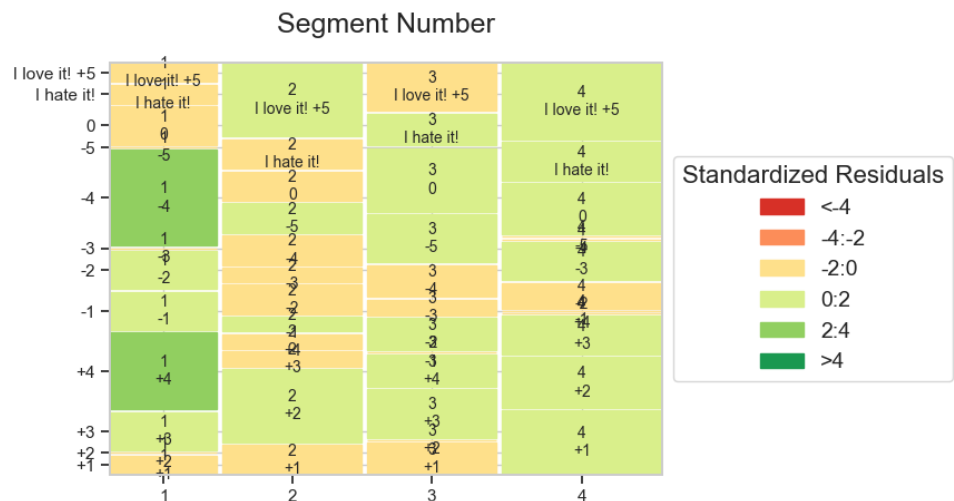


**Fig. A.14:** Shaded Mosaic Plot for Segment Membership and Like/Hate McDonald's

```
1   # Sample data for demonstration (replace with actual data)
2   np.random.seed(42)
3   segments = np.random.choice([1, 2, 3, 4], size=100)  # Simulated segment
    ↪  membership
4   like_hate = np.random.choice(['I hate it!', '-5', '-4', '-3', '-2', '-1',
    ↪  '0', '+1', '+2', '+3', '+4', 'I love it! +5'], size=100)  # Simulated
    ↪  like/hate variable
5   gender = np.random.choice(['Male', 'Female'], size=100)  # Simulated
    ↪  gender variable
6   age = np.random.randint(18, 70, size=100)  # Simulated age variable
7
8   # Convert to DataFrame for easier manipulation
9   df = pd.DataFrame({'Segment': segments, 'Like': like_hate, 'Gender':
    ↪  gender, 'Age': age})
10
```

```
11   # Create a contingency table
12   gender_contingency_table = pd.crosstab(df['Segment'], df['Gender'])
13
14   # Perform chi-squared test to get expected frequencies
15   chi2, p, dof, expected = chi2_contingency(gender_contingency_table)
16
17   # Calculate standardized residuals
18   gender_residuals = (gender_contingency_table - expected) /
     ↪  np.sqrt(expected)
19
20   # Define color mapping for standardized residuals
21   def gender_residual_color(value):
22       if value < -4:
23           return '#D73027'
24       elif value < -2:
25           return '#FC8D59'
26       elif value < 0:
27           return '#FEE08B'
28       elif value < 2:
29           return '#D9EF8B'
30       elif value < 4:
31           return '#91CF60'
32       else:
33           return '#1A9850'
34
35   # Function to format the properties of each cell
36   def gender_props(key):
37       segment, gender = key
38       value = gender_residuals.loc[int(segment), gender]
39       return {'color': gender_residual_color(value)}
40
41   # Generate the mosaic plot with shading for standardized residuals
42   plt.figure(figsize=(16, 12))
43   mosaic(gender_contingency_table.stack(), gap=0.01,
     ↪  properties=gender_props)
44   plt.xlabel('Segment Number', fontsize=18, labelpad=20)
45   plt.ylabel('Gender', fontsize=18, labelpad=20)
46   plt.title('Mosaic Plot for Gender Distribution across Segments',
     ↪  fontsize=22, pad=30)
47   plt.xticks(rotation=45, ha='right', fontsize=14)
48   plt.yticks(rotation=0, va='center', fontsize=14)
49
50   # Adding a color legend
51   legend_labels = ['<-4', '-4:-2', '-2:0', '0:2', '2:4', '>4']
52   colors = ['#D73027', '#FC8D59', '#FEE08B', '#D9EF8B', '#91CF60',
     ↪  '#1A9850']
53   patches = [mpatches.Patch(color=colors[i], label=legend_labels[i]) for i
     ↪  in range(len(legend_labels))]
```

```
54   plt.legend(handles=patches, title="Standardized Residuals", loc="center
     ↪  left", bbox_to_anchor=(1, 0.5), fontsize=14, title_fontsize=16)
55
56   # Adding grid lines and a background color
57   plt.grid(True, which='both', linestyle='--', linewidth=0.5)
58   plt.gca().set_facecolor('#f0f0f0')
59
60   # Display the plot
61   plt.show()
```
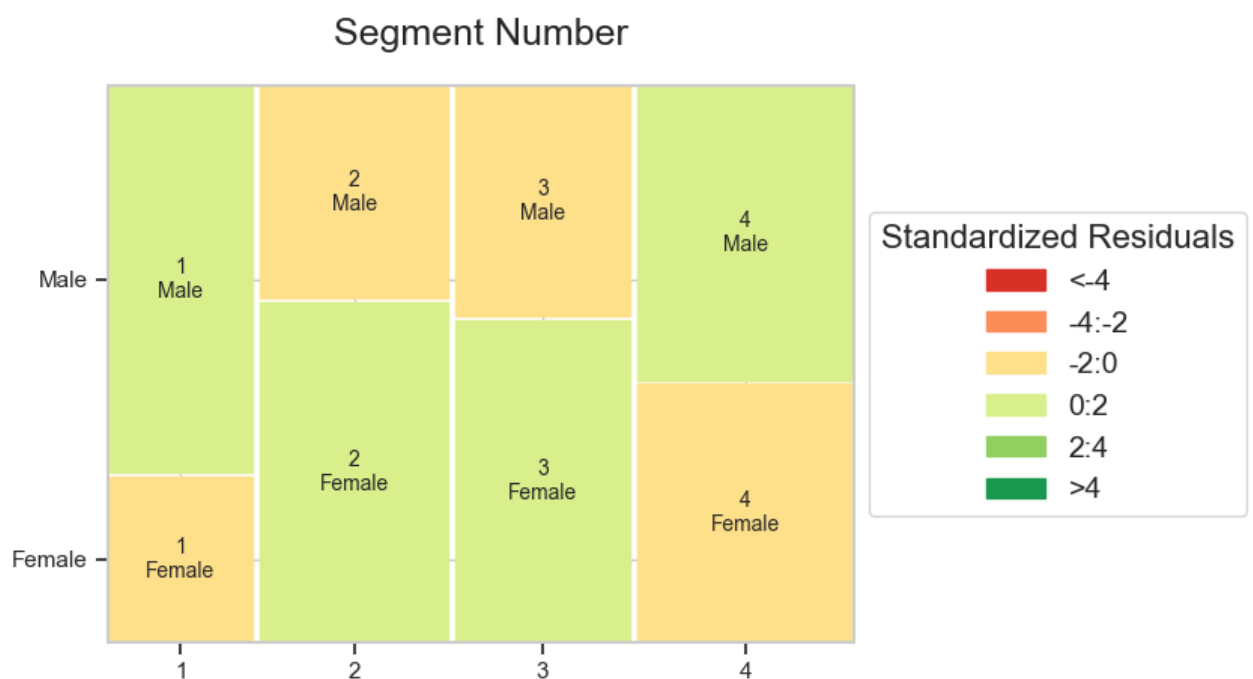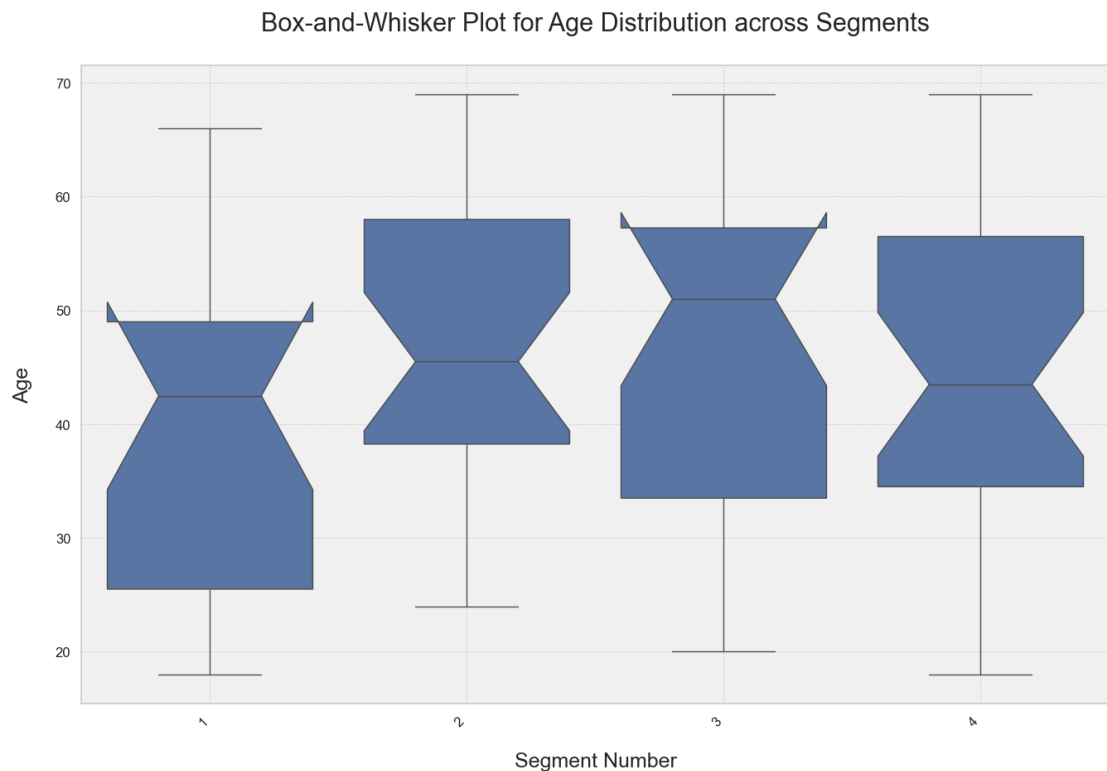


**Fig. A.15:** Mosaic Plot for Gender Distribution across Segments

```
1   plt.figure(figsize=(16, 10))
2   sns.boxplot(x='Segment', y='Age', data=df, notch=True)
3   plt.xlabel('Segment Number', fontsize=18, labelpad=20)
4   plt.ylabel('Age', fontsize=18, labelpad=20)
5   plt.title('Box-and-Whisker Plot for Age Distribution across Segments',
    ↪  fontsize=22, pad=30)
6   plt.xticks(rotation=45, ha='right', fontsize=12)
7   plt.yticks(fontsize=12)
8   plt.grid(True, which='both', linestyle='--', linewidth=0.5)
9   plt.gca().set_facecolor('#f0f0f0')
```
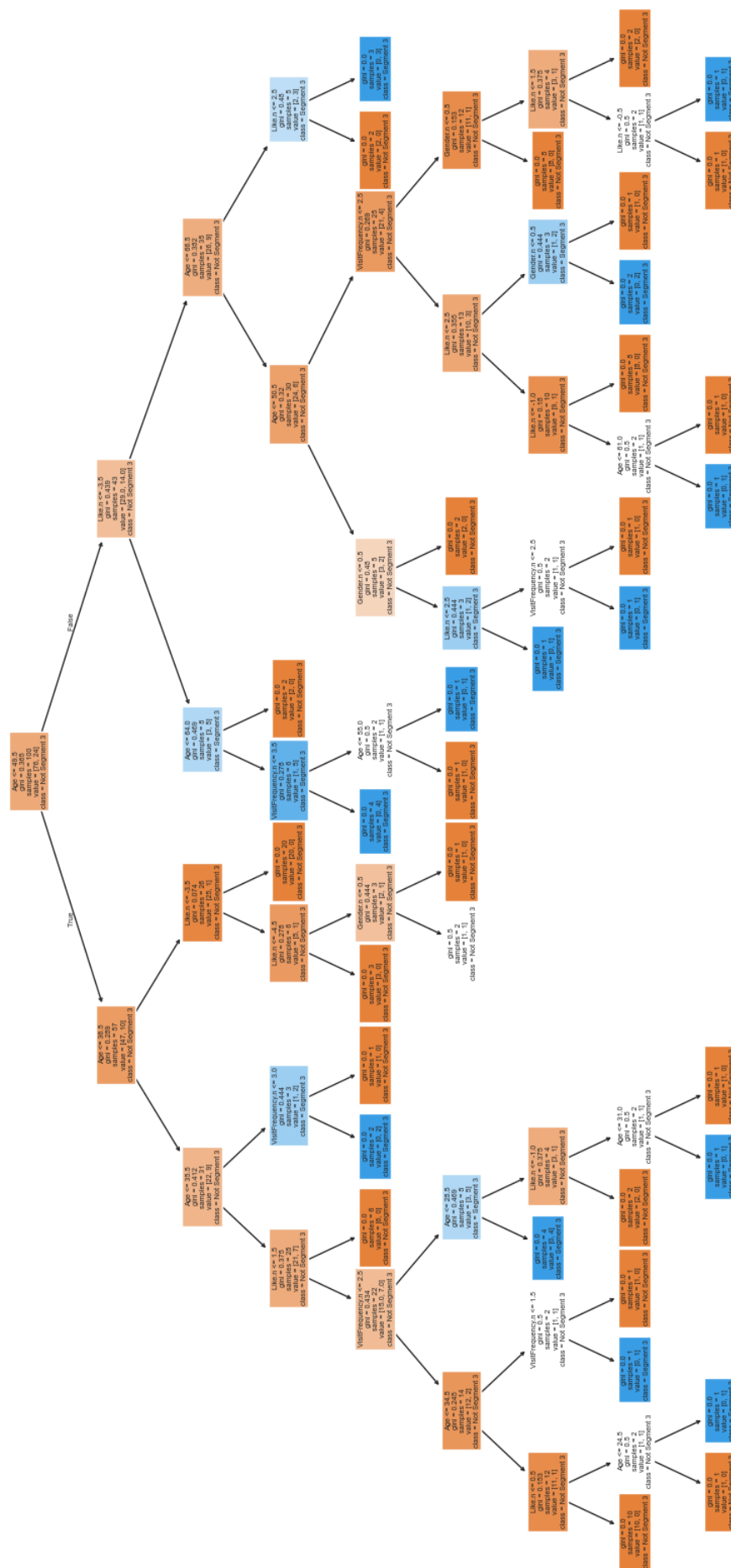
```
10    plt.show()
```

Box-and-Whisker Plot for Age Distribution across Segments



**Fig. A.16:** Box and Whisker Plot for Age

```
1    np.random.seed(42)
2    segments = np.random.choice([1, 2, 3, 4], size=100)  # Simulated segment
     ↪   membership
3    like_hate = np.random.choice(['I hate it!', '-5', '-4', '-3', '-2', '-1',
     ↪   '0', '+1', '+2', '+3', '+4', 'I love it! +5'], size=100)  # Simulated
     ↪   like/hate variable
4    gender = np.random.choice(['Male', 'Female'], size=100)  # Simulated
     ↪   gender variable
5    age = np.random.randint(18, 70, size=100)  # Simulated age variable
6    visit_frequency = np.random.choice(['Once a week', 'Once a month', 'More
     ↪   than once a month', 'Less than once a month'], size=100)
7
8    # Convert categorical variables to numerical codes
9    like_hate_map = {'I hate it!': -5, '-5': -5, '-4': -4, '-3': -3, '-2':
     ↪   -2, '-1': -1, '0': 0, '+1': 1, '+2': 2, '+3': 3, '+4': 4, 'I love it!
     ↪   +5': 5}
10   visit_frequency_map = {'Once a week': 1, 'More than once a month': 2,
     ↪   'Once a month': 3, 'Less than once a month': 4}
11   gender_map = {'Male': 0, 'Female': 1}
```

```
12
13   # Apply mappings
14   df = pd.DataFrame({
15       'Segment': segments,
16       'Like': like_hate,
17       'Gender': gender,
18       'Age': age,
19       'VisitFrequency': visit_frequency
20   })
21   df['Like.n'] = df['Like'].map(like_hate_map)
22   df['VisitFrequency.n'] = df['VisitFrequency'].map(visit_frequency_map)
23   df['Gender.n'] = df['Gender'].map(gender_map)
24   df['Segment3'] = (df['Segment'] == 3).astype(int)
25   X = df[['Like.n', 'Age', 'VisitFrequency.n', 'Gender.n']]
26   y = df['Segment3']
27
28   # Fit the decision tree classifier
29   tree = DecisionTreeClassifier()
30   tree.fit(X, y)
31
32   # Plot the decision tree
33   plt.figure(figsize=(20, 10))
34   plot_tree(tree, feature_names=['Like.n', 'Age', 'VisitFrequency.n',
     ↪  'Gender.n'], class_names=['Not Segment 3', 'Segment 3'], filled=True)
35   plt.title('Decision Tree for Segment 3 Membership Prediction')
36   plt.show()
```

**Fig. A.17:** Decision Tree for segment 3 membership

## A.5 Github Repository

The GitHub link for the jupyter notebook can be found at here.

Thank you for reading this report. Best wishes.