# Exploring Data (Step-4)

In Step 4 of our data analysis journey, known as "Exploring Data," we go deep into our collected data for the first time. Think of it as our first date with the data – we want to get to know it better! This step is important because it helps us clean up and prepare the data, getting it ready for the exciting task of finding market segments.

## First Glimpse at the Data:

In Step 4 of the data analysis process, known as "Exploring Data," we take our first look at the collected data. This stage involves conducting exploratory data analysis to clean and, if required, preprocess the data. It serves as a crucial guide for selecting the most appropriate algorithm to extract valuable market segments.

At a more technical level, data exploration accomplishes several key tasks:

1. Identifying the measurement levels of the variables.
2. Examining the individual distributions of each variable.
3. Assessing the relationships and dependencies between variables.

Additionally, data may undergo preprocessing and preparation to make it suitable for various segmentation algorithms. The insights gained during this data exploration phase help us determine the most effective methods for extracting market segments.

## Data Cleaning:

In this step which is the "Data Cleaning" phase, the first crucial step in data analysis is to ensure the data is clean and reliable. This involves several important tasks:

1. **Data Validation**: It starts with checking if all values have been recorded correctly and consistently. For instance, for numeric variables like "age," there are expected ranges of plausible values, such as between 0 and 110. Detecting values outside these ranges can indicate errors during data collection or entry.

2. **Categorical Variable Labels**: For categorical variables like "gender," it's essential to verify that consistent labels have been used. For instance, gender typically has two values: "female" and "male." If other values are present, they should be corrected during data cleaning.
3. **Sorting and Ordering**: Sometimes, data may appear out of order due to factors like data reading methods. Data can be read into factors (categorical variables) in R, which are sorted alphabetically by default. This can lead to unexpected ordering of variables. In such cases, the categories can be re-ordered to match the intended order.
4. **Reproducibility**: All data transformations and cleaning steps are typically documented and implemented using code to ensure reproducibility. This is essential for future analysis, documentation, and collaboration with other data analysts. It also enables the use of the same procedures with new data.
5. **Data Saving**: After cleaning, the data frame is often saved. This allows for easy reloading in future work sessions, maintaining the integrity of the cleaned data.

Data cleaning is a critical part of the data analysis process as it ensures the reliability and quality of the data, which is essential for making accurate and meaningful insights.

## Exploring Data: Descriptive Analysis

Understanding the data is a crucial step in data analysis to prevent misinterpretation of results, especially in complex analyses. Descriptive analysis, including numeric summaries and graphical representations, provides valuable insights into the data.

In this phase, various tools are available for descriptive analysis, with statistical software packages offering a wide array of options. In R, for example, the `summary()` command provides a numeric summary of data, including ranges, quartiles, means for numeric variables, and frequency counts for categorical variables. It also reports the number of missing values for each variable.

Graphical methods play a significant role in data exploration:

- **Histograms**: These visually represent the distribution of numeric variables, revealing characteristics like symmetry or skewness. Histograms group data into bins and display the frequency of observations in each bin.
- **Box-and-Whisker Plots (Boxplots)**: These compress data into minimum, first quartile, median, third quartile, and maximum values, providing a visual summary of a variable's distribution. They help identify skewness and outliers.
- **Bar Plots and Mosaic Plots**: Useful for visualizing categorical variables and associations between multiple categorical variables.

Reproducibility is emphasized through documenting and implementing data cleaning, exploration, and analysis steps using code. This ensures consistency and enables collaboration among data analysts. Additionally, the cleaned data can be saved for future use.

Descriptive analysis is crucial for gaining initial insights into the data's characteristics, guiding subsequent analyses, and making informed decisions about segmentation methods. It helps identify patterns, outliers, and distributional properties, paving the way for more in-depth analyses in later steps.

**Pre-Processing: Handling Categorical Variables**

Categorical variables often require specific pre-processing procedures to prepare them for further analysis. Two common approaches are:

**1. Merging Levels of Categorical Variables:** This step is useful when the original categories have too many subcategories, making the data too granular. For example, consider an income variable with various income ranges. If some of these categories have very few observations, they can be merged into larger categories to create a more balanced distribution. This can enhance the usefulness of the variable for analysis.

**2. Converting Categorical Variables to Numeric:** Sometimes, it makes sense to convert categorical variables into numeric ones, especially when the analysis methods assume numeric data. Ordinal data, such as income categories, can be converted to numeric values if the intervals between categories are reasonably equal. For instance, if income categories represent income ranges of approximately the same length, they can be converted into numeric values.

Similarly, for multi-category scales like Likert scales, where respondents provide answers like "STRONGLY DISAGREE," "DISAGREE," "NEITHER AGREE NOR DISAGREE," "AGREE," and "STRONGLY AGREE," it is assumed that the distances between these answer options are equal. If this assumption holds, the data can be treated as numeric. However, it's essential to be cautious as this assumption may not always be valid due to response styles and cultural factors.

Binary answer options, such as "YES" or "NO," are less susceptible to capturing response styles and are often preferred. They can be easily converted to numeric values (0 and 1) for analysis.

Pre-processing of categorical variables aims to ensure that the data is in a format suitable for the chosen analysis methods. It involves careful consideration of the nature of the variables and the assumptions made during the conversion process.

**Pre-Processing Numeric Variables**

When working with numeric variables in the context of data preprocessing for segmentation analysis, it's essential to consider their range of values. This range can significantly impact the results of distance-based segmentation methods. To ensure that each segmentation variable contributes fairly to the segmentation process, standardization is often employed. Standardization involves transforming the variables onto a common scale.

In statistical terms, standardization is achieved by subtracting the mean ($\bar{x}$) and dividing by the standard deviation (s) of the variable. This transformation creates a new variable (z) for each original variable (x):

$$z_i = (x_i - \bar{x}) / s$$

The result of standardization is that the mean of the new variable (z) becomes 0, and its standard deviation becomes 1. This normalization process ensures that all variables have a consistent scale, making them directly comparable.

It's important to note that when dealing with data that contains outliers (extreme values), alternative standardization methods may be more suitable. These methods use robust estimates for location (e.g., median) and spread (e.g., interquartile range) to prevent outliers from unduly influencing the standardization process. This approach helps maintain the integrity of the segmentation analysis.

**Principal Components Analysis (PCA):**

Principal Components Analysis (PCA) is a powerful technique used to transform a multivariate dataset, typically consisting of metric variables, into a new set of uncorrelated variables known as principal components. These principal components are arranged in descending order of importance, with the first component explaining the most variance in the data, the second component explaining the second most variance, and so on. This transformation allows analysts to explore the data from a different perspective while preserving the relative relationships between observations.

PCA is typically applied to the covariance or correlation matrix of the original variables. When the variables are measured on the same scale and have similar data ranges, the choice between covariance and correlation matrix doesn't significantly impact the results. However, when variables have different scales, it's advisable to use the correlation matrix, which is equivalent to standardizing the data.

The importance of each principal component is determined by its associated standard deviation. The first few principal components often account for the majority of the variance in the data, while subsequent components explain progressively less variance. The proportion of variance explained by each component and the cumulative proportion are critical in assessing the overall contribution of principal components to the dataset's variation.

PCA can be valuable for reducing the dimensionality of high-dimensional data by projecting it into lower-dimensional space for visualization purposes. This is typically done by selecting a subset of the most important principal components.

It's worth noting that while PCA can provide valuable insights into the data's structure and redundancy among variables, it's not typically

recommended for variable reduction in segmentation analysis. Replacing original variables with a subset of principal components can lead to a different basis for segmentation, potentially introducing unintended biases.

In practice, PCA is often used as an exploratory tool to identify highly correlated variables, which can then be removed from the segmentation analysis to reduce dimensionality while retaining the original variables.

In summary, PCA is a valuable technique for exploring data, visualizing high-dimensional data, and identifying correlations among variables.

## Summary of Step 4 - Data Exploration and Pre-processing:

In Step 4, we focused on exploring and preparing the data for the segmentation analysis. Here's a summary of the key actions and considerations:

1. **Data Exploration:**
   - We thoroughly explored the data to identify any inconsistencies or systematic issues.
   - If inconsistencies were found, we took steps to clean the data, ensuring its integrity.
2. **Data Pre-processing:**
   - Pre-processing of the data was conducted to make it suitable for analysis.
   - For numeric variables, we considered standardization to ensure variables were on a common scale.
   - For categorical variables, we merged levels if necessary and converted them to numeric format where applicable.
3. **Segmentation Variable Assessment:**
   - We assessed the number of segmentation variables, ensuring that there were a sufficient number of consumers (a minimum of 100) for each variable.
   - If there were too many segmentation variables, we explored methods to select a manageable subset.
4. **Variable Correlation:**
   - We checked for correlations among segmentation variables.

- o If variables were highly correlated, we selected a subset of uncorrelated variables to reduce redundancy.
5. **Data Preparation:**
   - o Finally, we prepared the cleaned and pre-processed data for use in Step 5, where segments will be extracted based on this data.

By following these steps, we ensured that our data was in optimal condition for the subsequent segmentation analysis.