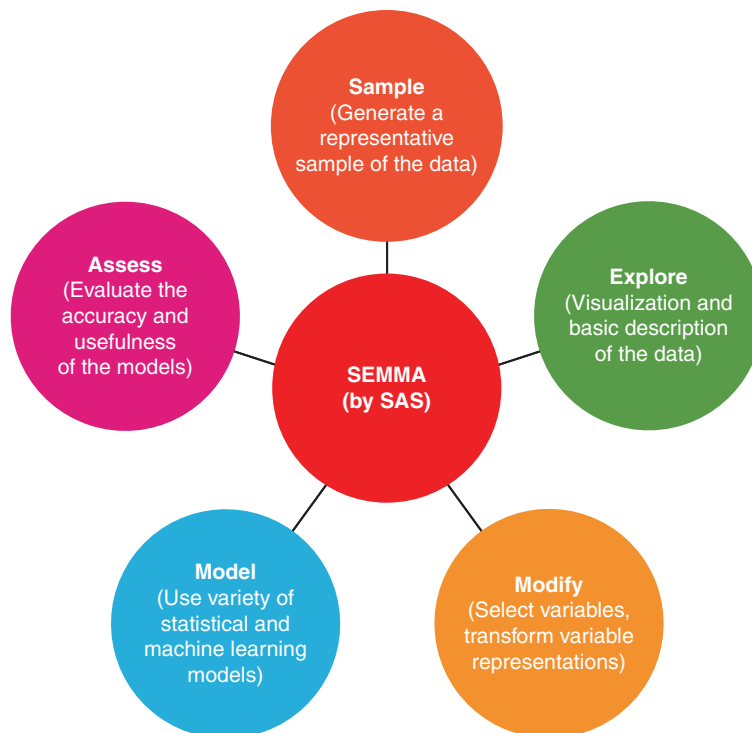


Chapter 4

SEMMA Process



© Kendall Hunt Publishing Company

4.1 Sampling

The initial stage in the SEMMA process is sampling. Upon collecting the data, the decision of whether to sample the data arises. In statistical terms, the complete set of data is referred to as the population, while a portion of this population is termed a sample. While constructing a model using the entire population may yield superior results, the computational demands associated with processing the entire population often present significant challenges. Hence, sampling is employed as a practical solution to manage these computational complexities.

Why Do We Sample?

- Efficiency
 - Sampling allows us to work with a subset of data rather than the entire population, which can be time-consuming and resource-intensive to process, especially for large datasets. By analyzing a representative sample, we can achieve results more efficiently.
- Cost-effectiveness
 - Collecting and analyzing data from the entire population can be expensive. Sampling reduces costs by requiring fewer resources for data collection, storage, and analysis.
- Feasibility
 - In many cases, it may be impractical or impossible to collect data from the entire population due to logistical constraints such as time, budget, or accessibility. Sampling enables us to overcome these limitations by working with a manageable subset of the population.
- Inference
 - A well-designed sample can provide accurate estimates and insights about the population as a whole. By ensuring that the sample is representative of the population, we can make valid inferences and draw conclusions that apply to the broader population.

When Do We Sample?

- Large datasets
 - When dealing with large datasets that are computationally intensive to process, sampling is often employed to work with a manageable subset of the data. Sampling allows for faster analysis without sacrificing the validity of the results.
 - Note: In this textbook, we will not be sampling as the datasets are particularly small, however, we will cover it in theory.
- Resource constraints
 - Sampling is used when there are limitations in resources such as time, budget, or computing power.
 - By analyzing a representative sample instead of the entire population, resources can be conserved while still obtaining valuable insights.
- Quality assurance
 - Sampling may be used in quality assurance processes to assess the accuracy and reliability of data. By randomly selecting samples from the dataset, analysts can verify data integrity, identify errors, and ensure data quality.

Mathematical Approach to Sampling

- Simple random sampling
 - Every member of the population has an equal chance of being selected for the sample. This can be achieved using random number generators or random selection methods.
- Stratified sampling
 - The population is divided into homogeneous groups (strata), and random samples are taken from each stratum. This ensures representation from all subgroups within the population.
- Systematic sampling
 - Members of the population are selected at regular intervals from a randomly chosen starting point. This method is less time-consuming than simple random sampling and ensures systematic coverage of the population.
- Multistage sampling
 - This involves a combination of two or more sampling methods.

Nonmathematical Approach to Sampling

- Convenience sampling
 - Samples are selected based on convenience or accessibility rather than random selection. This method is quick and inexpensive but may result in biased samples.
- Judgmental sampling
 - Samples are selected based on the judgment of the researcher or expert knowledge of the population. This method is subjective and may introduce bias.
- Quota sampling
 - The population is divided into quotas based on certain characteristics (e.g., age, gender), and samples are selected to meet these quotas. Quota sampling is similar to stratified sampling but does not involve random selection.
 - In marketing, this is customer segmentation.

Each sampling method has its advantages and limitations, and the choice of method depends on factors such as the business objectives, the characteristics of the population, and constraints. Once the data is sampled, the next step is to explore the data.

4.2 Exploring

Exploration is a fundamental phase in the data analysis process, providing insights into the structure and characteristics of the dataset. It can be broadly categorized into two aspects: dataset exploration and variable exploration.

Dataset Exploration

In dataset exploration, the focus is on understanding the core aspects of the dataset as a whole. This phase is mostly descriptive analytics. This involves examining various properties and attributes of the dataset to gain an overall understanding of its structure and composition.

- Shape of dataset
 - The shape of the dataset refers to its dimensions, typically represented as the number of rows and columns.
 - Understanding the shape of the dataset helps in assessing its size and complexity, which is essential for determining the appropriate analytical techniques and resources required for analysis.
- Missing values
 - Identifying missing values is crucial for data quality assessment. Missing values can occur due to various reasons, such as data entry errors, incomplete data collection, or system failures.
 - Exploring the presence and extent of missing values helps in deciding how to handle them, whether by imputation, deletion, or other methods.
 - We will cover more of this in the following chapters.
- Data types
 - Examining the data types of variables provides insights into the nature of the data and how it should be treated during analysis.
 - Common data types include numerical (e.g., integers, floats), categorical (e.g., strings, factors), and datetime. Understanding data types is important for data manipulation, visualization, and modeling.
 - We will cover more of this in the following chapters.

Variable Exploration

Variable exploration is a crucial step in the data analysis process, focusing on understanding the characteristics and relationships between individual variables within the dataset. This phase combines elements of descriptive and diagnostic analytics to uncover insights and patterns that can inform further analysis and decision-making. Here are the main methods of exploring variables:

- **Univariate analysis:** In univariate analysis, each variable is examined individually to understand its distribution and characteristics. Here are some common methods used:
 - **Basic statistics:** Descriptive statistics such as mean, median, mode, standard deviation, and range provide summary measures of the variable's central tendency and dispersion.
 - **Histograms:** Histograms are graphical representations of the frequency distribution of numerical variables. They display the frequency of data values within predefined intervals or bins, allowing analysts to visualize the shape and spread of the data distribution.
 - **Boxplots:** Boxplots also known as box-and-whisker plots, provide a visual summary of the distribution of numerical variables. They display the median, quartiles, and potential outliers in the data, making it easy to identify central tendency and variability.
- **Bivariate analysis:** Bivariate analysis involves exploring the relationship between two variables. Here are some common methods used:
 - **Bar plots:** Bar plots are used to visualize the relationship between a categorical variable and a numerical variable. They display the mean or median value of the numerical variable for each category of the categorical variable, allowing analysts to compare group means or medians.
 - **Scatterplots:** Scatterplots are used to visualize the relationship between two numerical variables. Each data point represents a combination of values from both variables, allowing analysts to identify patterns, trends, or clusters in the data.
 - **Line charts:** Line charts are useful for visualizing the relationship between two numerical variables over time or across different categories. They display data points connected by line segments, making it easy to observe trends or changes in the data over time.
- **Multivariate analysis:** Multivariate analysis involves exploring the relationship between three or more variables simultaneously. Here are some common methods used:
 - **Stacked bar charts:** Stacked bar charts are used to visualize the relationship between two categorical variables while showing the proportion of each category within each group. They provide insights into the distribution of categorical variables across different groups or categories.
 - **Scatterplots with color:** Scatterplots with color coding are used to visualize the relationship between three numerical variables. Each data point is represented by a dot, with color indicating the value of a third categorical variable. This allows analysts to explore the relationship between multiple variables while accounting for additional factors.

After exploring the data, the next step is to modify it. Usually exploration and modification go hand-in-hand, as we will see when going through the example in the following chapters.

4.3 Modifying

In the modifying phase of data preparation, adjustments are made to the dataset to prepare it for modeling purposes. This phase involves several key steps:

- **Data types modification**
 - The first step involves assessing and correcting any incorrect data types within the dataset. This may include converting categorical variables to a numerical format or vice versa, ensuring consistency and compatibility for subsequent analysis.

- Missing value treatment
 - The handling of missing values is critical for maintaining data integrity and model performance.
 - There are several approaches to treating missing values, including
 - Keeping them as is: In some cases, missing values may carry meaningful information and should be retained.
 - Deleting rows: Rows with a high proportion of missing values (e.g., exceeding 5%) may be removed from the dataset to preserve data quality.
 - Imputation: Missing values can be replaced with estimated values based on statistical methods such as mean, median, or mode imputation.
- Outlier treatment
 - Outliers, or observations that deviate significantly from the rest of the data, can skew statistical analysis and model results.
 - Options for handling outliers include
 - Keeping them as is: Outliers may represent genuine data points and should be retained if they are valid and not the result of data entry errors or measurement issues.
 - Deleting outliers: In some cases, outliers may be removed from the dataset to prevent them from unduly influencing model outcomes.
- Data transformation
 - Data transformation involves standardizing variables to ensure they are on the same scale, facilitating comparison and analysis across different features.
 - Common techniques for data transformation include
 - Standardization: Scaling numerical variables to have a mean of zero and a standard deviation of one, ensuring that variables are comparable and contribute equally to the analysis.
 - Normalization: Scaling variables to a specific range (e.g., between 0 and 1), preserving the relative relationships between data points while constraining them within a defined range.

By addressing these modifications in the dataset, analysts can ensure that the data is appropriately structured and prepared for subsequent modeling and analysis tasks. Now, the data is ready to model.

4.4 Modeling

In the modeling phase, the dataset, having undergone thorough exploratory data analysis and preprocessing, is now prepared for building predictive models. This phase marks the transition from data exploration to model development, where the focus shifts toward leveraging the insights gained earlier to construct predictive algorithms.

Linear Regression

Linear regression is a fundamental statistical method used for modeling the relationship between a dependent variable (often denoted as “y”) and one or more independent variables (denoted as “x”). In this textbook, linear regression is employed when the target variable (y) is numerical or continuous in nature. The model aims to fit a linear equation to the data, allowing for the prediction of the dependent variable based on the values of the independent variables.

Logistic Regression

Logistic regression, on the other hand, is utilized when the target variable (y) is categorical or binary in nature, typically representing two outcomes (e.g., yes/no, pass/fail). In logistic regression, the model estimates the probability that an observation belongs to a particular category, given the values of the independent variables. This technique is commonly used for binary classification problems.

In the modeling phase, the following steps are typically undertaken:

- Model selection
 - Based on the nature of the target variable and the problem at hand, the appropriate modeling technique (linear regression or logistic regression) is selected.
- Model building
 - The selected model is built using the cleaned and modified input variables (x) and the target variable (y). The model parameters are estimated using techniques such as ordinary least squares (for linear regression) or maximum likelihood estimation (for logistic regression).

4.5 Assessing

In the final phase of model assessment, each constructed model undergoes rigorous evaluation to gauge its performance and predictive accuracy. Several model evaluation metrics are employed to provide a comprehensive understanding of the model's strengths and weaknesses. Recommendations are then derived based on the assessment outcomes to guide decision-making processes.

Model Evaluation

Upon completion of model construction, an array of metrics is utilized to evaluate the model's performance and predictive accuracy.

For linear regression models, metrics such as R -squared are commonly employed to assess the goodness of fit, measuring the proportion of variance explained by the model.

For logistic regression models, evaluation metrics include the confusion matrix, accuracy, precision, recall, and F1-score, which collectively offer insights into the model's classification performance and ability to discriminate between classes.

Model Validation

To ascertain the reliability and generalizability of the model, rigorous validation techniques such as cross-validation or holdout validation are employed. These validation techniques help assess the model's performance on unseen data, mitigate the risk of overfitting, and ensure its robustness across different datasets.

Model Interpretation

In the final stage of assessment, the model results are meticulously interpreted to extract meaningful insights and elucidate the relationships between the independent and dependent variables.

This involves scrutinizing the model coefficients, significance levels, and confidence intervals of the estimated parameters. By examining the magnitude and direction of the coefficients, analysts can discern the impact of each independent variable on the dependent variable, thereby gaining deeper insights into the underlying mechanisms driving the observed outcomes.

Incorporating Business Context

In addition to model evaluation, validation, and interpretation, it's imperative to incorporate business context into the assessment process. This contextual understanding enriches the analysis by aligning model outcomes with organizational goals and requirements, thereby enhancing the relevance and applicability of the results.

By integrating business context into the analysis, analysts can tailor recommendations to address specific challenges or capitalize on emerging opportunities, ensuring that the models deliver actionable insights that resonate with the strategic objectives and operational realities of the business.