

The **analysis of data follows a structured, iterative framework** that guides the transformation of raw data into actionable insights. This process encompasses stages such as data acquisition, preparation, analysis, modeling, and communication, ensuring methodological rigor and reproducibility in data-driven projects. Frameworks like CRISP-DM (“Cross-Industry Standard Process for Data Mining”), OSEMN (“Obtain, Scrub, Explore, Model, iNterpret”), and the model of the data science process, described in “R for Data Science” (<https://r4ds.hadley.nz/intro.html>). Below is a concise, scientifically grounded outline of the data science process:

1. Load Required Packages

Initiate by importing essential libraries (e.g., `tidyverse`, `data.table`, `haven`, `psych`, ...) to facilitate data manipulation, visualization, and analysis.

2. Load the Dataset

Import data from various sources such as CSV files, databases, or APIs using functions like `read_csv()` or `dbConnect()`.

3. Load Custom Functions

Incorporate user-defined functions or external scripts to streamline repetitive tasks and enhance code modularity.

4. Data Preparation

a. Initial Data Exploration

Conduct exploratory data analysis to understand data structure, distributions, and potential anomalies. Techniques include summary statistics, histograms, and scatter plots.

b. Handling Missing Values

Identify missing data patterns and decide on appropriate strategies:

- **Deletion:** Remove incomplete cases using `na.omit()`.
- **Imputation:** Estimate missing values through methods like mean substitution or multiple imputation (read: <https://stefvanbuuren.name/fimd/sec-MCAR.html>).

Note: When subsetting, ensure to exclude NA values explicitly (e.g., `age > 95 & !is.na(age)`).

c. Subsetting the Dataset

Apply inclusion/exclusion criteria to focus on relevant subsets, enhancing the specificity of subsequent analyses.

d. Data Transformation

- **Tidy Data:** Reshape data into a long format where each variable is a column, each observation is a row, and each type of observational unit forms a table.
- **Feature Engineering:** Create new variables or composite indicators that encapsulate underlying patterns or criteria.
- **Outlier Detection:** Identify and address anomalies using domain knowledge (e.g., reaction times < 300ms) or statistical methods like Mahalanobis distance.

e. Further Data Exploration

Re-examine the transformed data to validate the effects of preprocessing steps and to uncover additional insights.

5. Descriptive Statistics

Generate a comprehensive overview of the dataset:

- **Univariate Analysis:** Assess individual variables using measures like mean, median, and standard deviation.
- **Bivariate/Multivariate Analysis:** Explore relationships between variables through cross-tabulations and correlation matrices.
- **Normality Assessment:** Evaluate distributional assumptions using tests (e.g., Shapiro-Wilk, Kolmogorov-Smirnov) and visualizations (e.g., Q-Q plots, density plots).
- **Group Comparisons:** Examine differences across groups to identify potential effects or trends.

6. Inferential Statistics

Apply statistical models to test hypotheses and infer population parameters:([wired.com](https://www.wired.com))

- **Hypothesis Testing:** Utilize t-tests, ANOVA, or non-parametric equivalents to assess group differences.
- **Regression Analysis:** Implement linear or logistic regression models to examine relationships between variables.
- **Hierarchical Models:** Employ mixed-effects models to account for nested data structures.
- **Factor Analysis:** Explore underlying latent constructs within the data.
- **Cluster Analysis:** Identify natural groupings within the data using methods like k-means or hierarchical clustering.
- ...

7. Communication of Results

Synthesize findings into a coherent narrative, emphasizing clarity and adherence to reporting standards such as APA 7. This includes:

- **Data Visualization:** Create informative plots and charts to illustrate key results.
- **Statistical Reporting:** Present estimates, confidence intervals, and p-values with appropriate interpretation.
- **Reproducibility:** Document the analysis workflow to facilitate replication and verification by others. At best write all your analyses within a dynamic script (e.g., using the <https://quarto.org/> framework).