



# Data Wrangling: Transforming Raw Data into Insights

Data wrangling is the crucial process of cleaning, structuring, and transforming raw data into a format that can be effectively analyzed. This essential step lays the foundation for uncovering valuable insights and informed decision-making.

# Understanding Data Formats and Structures

## Structured Data

Organized data with a clear schema, such as spreadsheets and databases.

## Unstructured Data

Free-form data without a defined structure, like text files, images, and social media posts.

## Semi-Structured Data

Data with some structure, like XML and JSON, which can be transformed into a tabular format.



# Cleaning and Preprocessing Data

## Handling Errors

Identifying and correcting invalid, incomplete, or duplicate data entries.

## Standardizing Formats

Ensuring consistent data types, units, and conventions across the dataset.

## Dealing with Outliers

Detecting and addressing extreme or anomalous values that may skew analysis.

## Enriching Data

Supplementing the dataset with additional relevant information from external sources.



# Handling Missing and Erroneous Values

## 1 Imputation Techniques

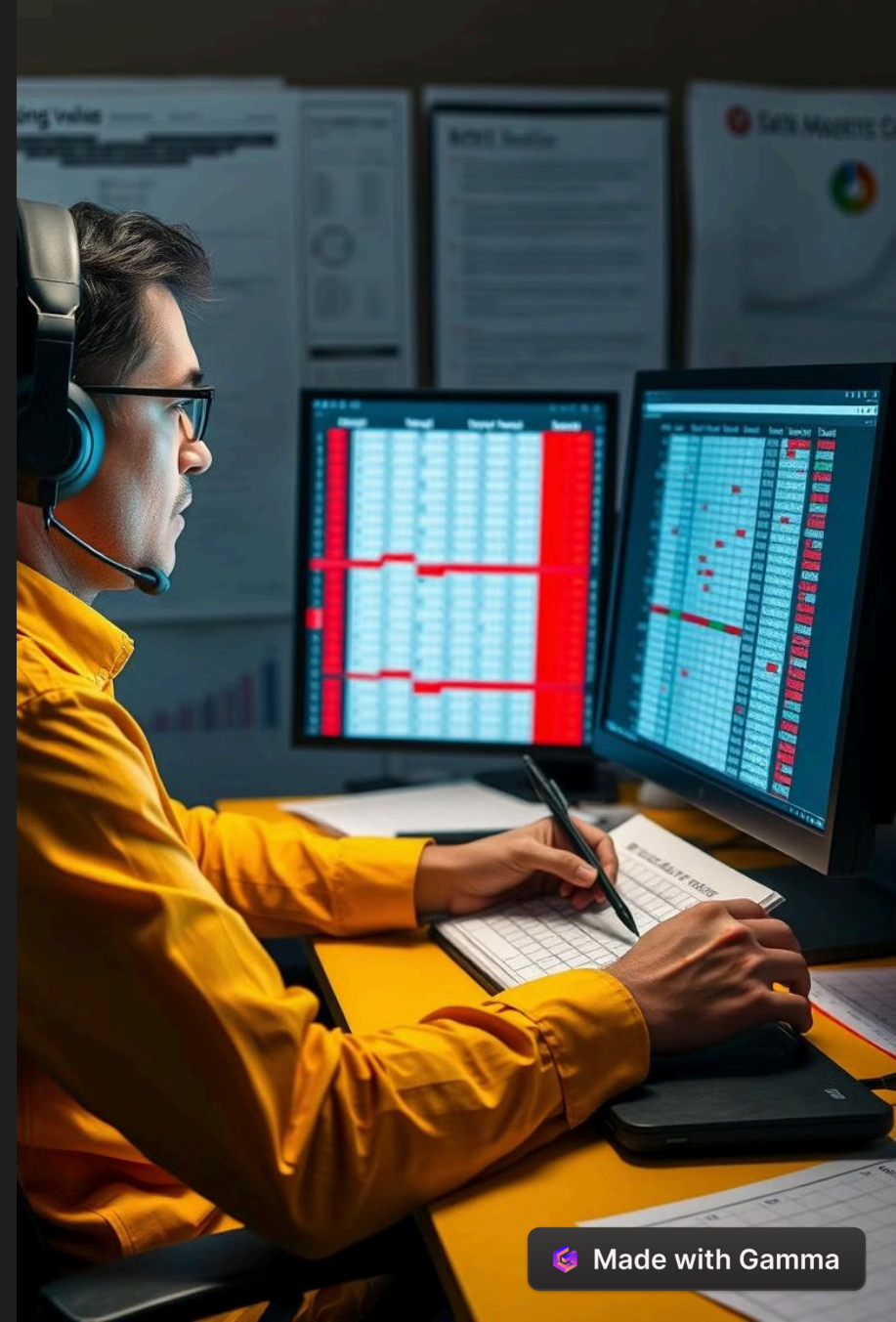
Replacing missing values with estimated data based on statistical methods or patterns in the dataset.

## 2 Flagging Errors

Identifying and marking problematic data points for further review or exclusion from analysis.

## 3 Balancing Trade-offs

Considering the impact of data cleaning decisions on the overall quality and representativeness of the dataset.



# Merging and Joining Datasets

1

## Vertical Merge

Combining datasets with the same structure but different observations, like sales data over time.

2

## Horizontal Merge

Combining datasets with the same observations but different features, like customer data and transaction history.

3

## Database Joins

Linking datasets based on common identifiers, such as customer IDs or product codes.





# Transforming and Reshaping Data



## Filtering

Selecting relevant subsets of data based on specific criteria.



## Pivoting

Restructuring data from long to wide format, or vice versa.



## Aggregation

Summarizing data by calculating metrics like sums, averages, or counts.



## Derivation

Creating new features or variables from existing data, such as percentages or ratios.





# Automating Data Wrangling Workflows

1

## Data Ingestion

Automatically collecting and importing data from multiple sources.

2

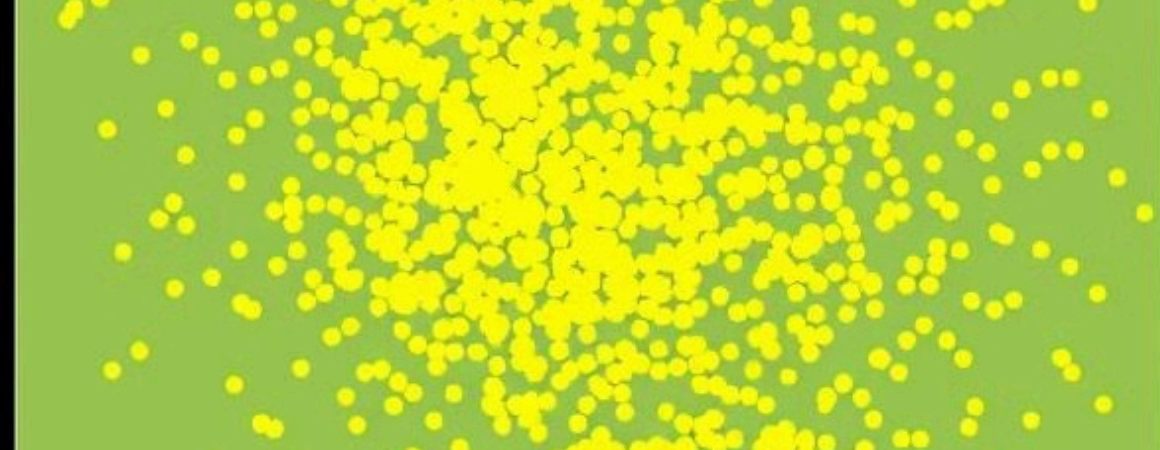
## Data Transformation

Applying a series of cleaning, merging, and reshaping steps.

3

## Data Storage

Storing the transformed data in a centralized repository for further analysis.



# Visualizing the Impact of Data Wrangling

Before Data Wrangling	After Data Wrangling
Inconsistent data formats	Standardized and cleaned data
Missing and erroneous values	Imputed and flagged issues
Disjointed datasets	Merged and integrated data