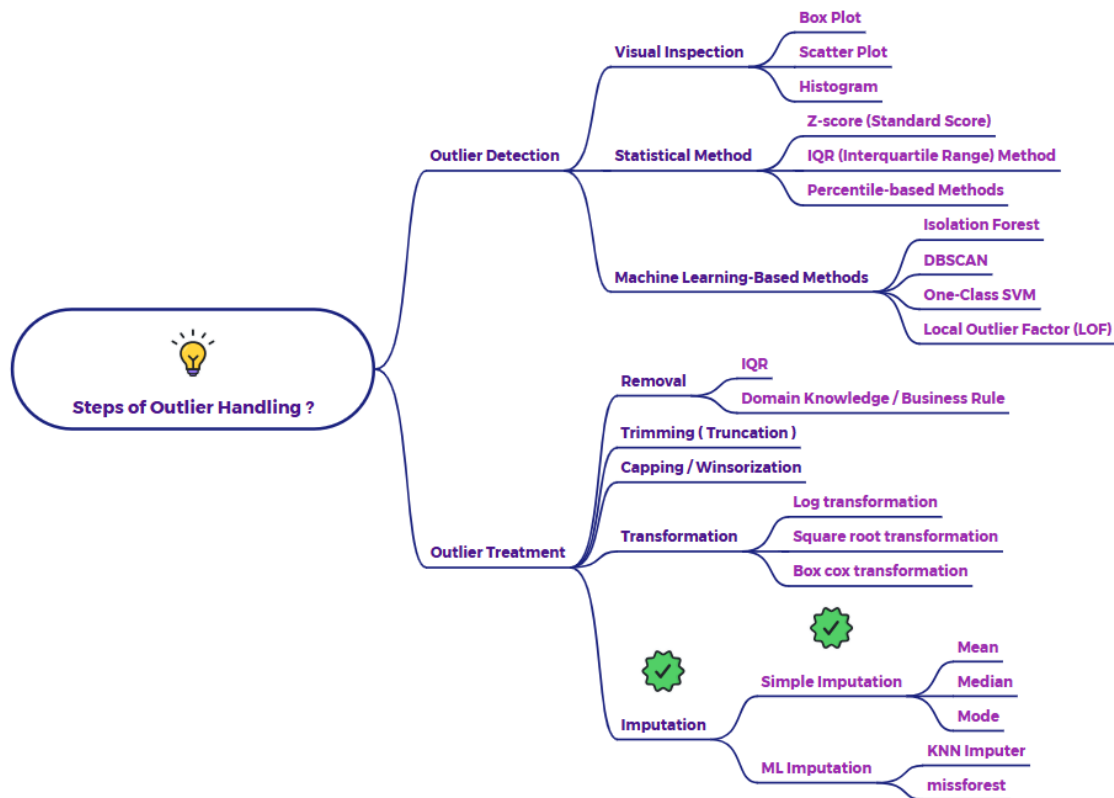


## Explain Outlier treatment through simple imputation



### Outlier Treatment: Simple Imputation

Imputation, in the context of outlier handling, involves replacing outlier values with more plausible estimates. Simple imputation methods use a single, easily calculated value to replace all outliers in a given variable.

#### Methods and Examples:

##### 1. Mean Imputation:

- Replace outliers with the mean of the remaining data (i.e., the data without the outliers).
- **Example:**
  - Dataset: [10, 20, 30, 40, 50, 60, 70, 80, 90, 500]
  - Identify 500 as an outlier (e.g., using the IQR method).
  - Calculate the mean of the remaining data:  $(10 + 20 + 30 + 40 + 50 + 60 + 70 + 80 + 90) / 9 = 50$
  - Replace the outlier: [10, 20, 30, 40, 50, 60, 70, 80, 90, 50]

## 2. Median Imputation:

- Replace outliers with the median of the remaining data.
- **Example:**
  - Dataset: [10, 20, 30, 40, 50, 60, 70, 80, 90, 500]
  - Identify 500 as an outlier.
  - Calculate the median of the remaining data: (10, 20, 30, 40, 50, 60, 70, 80, 90) => 50
  - Replace the outlier: [10, 20, 30, 40, 50, 60, 70, 80, 90, 50]

## 3. Mode Imputation:

- Replace outliers with the mode of the remaining data. This is typically used for categorical data, but can be used for numerical data as well.
- **Example:**
  - Dataset: [10, 20, 30, 40, 50, 60, 70, 80, 90, 10]
  - Identify 90 as an outlier.
  - Calculate the mode of the remaining data: [10, 20, 30, 40, 50, 60, 70, 80, 10] => 10
  - Replace the outlier: [10, 20, 30, 40, 50, 60, 70, 80, 10, 10]

## 4. Constant Value Imputation:

- Replace outliers with a predefined constant value. This value should be chosen based on domain knowledge.
- **Example:**
  - Dataset: [25, 30, 35, 40, 45, 50, 55, 60, 65, 150] (representing age)
  - Identify 150 as an outlier.
  - Based on domain knowledge, we know the maximum reasonable age is 100.
  - Replace the outlier: [25, 30, 35, 40, 45, 50, 55, 60, 65, 100]

## When to Use Simple Imputation for Outliers:

- When you suspect that the outliers are due to data entry errors or other non-representative issues.
- When you want a quick and easy way to handle outliers.
- When the machine learning model you are using requires complete data (no missing values).

### Cautions:

- Simple imputation can distort the distribution of the data, especially if there are many outliers.
- Mean imputation is sensitive to outliers itself.
- It does not account for the relationship between variables.
- Median imputation is more robust to outliers than mean imputation.