

Step 03

# DATA CLEANING

A Guide to the Data Cleaning Process



Data cleaning is a critical step in building a machine learning model as it ensures that the data used for training the model is accurate, complete, consistent, and relevant.

Data cleaning has several advantages while building a machine learning model. Some of them are,

**Improves accuracy :**

Improve the accuracy of the machine learning model by removing the inconsistencies and errors in the dataset.

**Reduces bias :**

It helps to reduce the bias in the machine learning model by removing the irrelevant and duplicate data.

**Improves performance :**

Improve the performance of the machine learning model by removing the noise and outliers that can affect the model's ability to make accurate predictions.

### **Saves time and resources :**

By cleaning the data before building the machine learning model, you can save time and resources by avoiding errors and inconsistencies that can arise during the modeling process.

### **Provides better insights :**

Data cleaning can help to uncover patterns and relationships in the data that may have been hidden by errors and inconsistencies. This can provide better insights into the underlying processes and variables that affect the outcome being predicted.

### **Enhances decision-making :**

Clean data provides more accurate insights, which ultimately leads to better decision-making and improved business outcomes.

so now, Are you wondering how to perform data cleaning in your data?

Let's break it down into the key tasks involved in the process.

### **Handling missing data:**

Missing data is a common problem in datasets and can be handled using techniques like imputation.

### **Handling duplicates:**

Duplicates can distort analysis and it is important to remove them to get an accurate picture of the data.

### **Handling outliers:**

Outliers are extreme values that can skew the results of the analysis. They can be handled by removing or transforming them.

### **Handling inconsistent data:**

Inconsistent data can arise due to data entry errors, incorrect data conversions, etc. It can be handled by standardizing the data or using regular expressions.

### **Handling invalid data:**

Invalid data can arise due to data entry errors or data that falls outside of the expected range. It can be handled by removing or correcting the data.

### **Handling noisy data:**

Noisy data is data that contains errors or outliers that can make analysis difficult. It can be handled by smoothing the data or using interpolation techniques.

### **Handling data formatting:**

Data formatting involves making sure that the data is in a consistent format. It can be handled by converting data types, changing date formats, etc.

### **Handling date and time data :**

Date and time data can be complex and difficult to work with. It can be handled by formatting the data correctly or converting it to a more manageable format.

### **Handling text data :**

Text data can be unstructured and messy. It can be handled by cleaning and preprocessing the text or using natural language processing techniques.

### **Handling categorical data :**

Categorical data is non-numeric data that can be difficult to work with. It can be handled by converting it to numerical data or using techniques like one-hot encoding.

By performing the mentioned tasks and other best practices, we can improve the quality and reliability of our data and achieve more accurate and robust machine learning models.

Data cleaning can be a time-consuming and complex process, but it is critical to the success of our machine learning project. By investing time and effort into data cleaning, we can ensure that our model is built on a solid foundation of accurate and reliable data.

Happy learning! 😊

# THANK YOU

@ Gangadhar Neelam