

Data cleaning is an essential part of data preprocessing, where the focus is on identifying and rectifying errors, inconsistencies, and inaccuracies in the raw data.

Clean data is crucial for obtaining meaningful insights and building accurate machine learning models.

Data cleaning involves several steps and techniques, including:

1. Handling Missing Data:

- **Identify Missing Values:** Detecting and locating missing values in the dataset, represented as NaN (Not a Number) or other placeholders.
- **Removing Rows or Columns:** If the missing data is significant, you may consider removing entire rows or columns that contain too many missing values.
- **Imputation:** Filling in missing values with estimated or imputed values. Common imputation techniques include mean, median, mode imputation, or more advanced methods like k-Nearest Neighbors imputation or regression-based imputation.

2. Dealing with Outliers:

- **Identify Outliers:** Detecting data points that are significantly different from the rest of the data.
- **Handle Outliers:** Depending on the context, outliers can be corrected, removed, or transformed using techniques like truncation or capping.

3. Data Validation:

- **Check Data Integrity:** Ensuring that the data adheres to predefined business rules and

constraints.

- Cross-Field Validation: Verifying that relationships between different fields in the data are consistent and logical.

4. Data Type Conversion:

- Ensure Correct Data Types: Verifying that each feature or attribute is of the correct data type (e.g., numeric, categorical, date, etc.).
- Convert Data Types: Converting data to the appropriate format, such as converting dates from strings to date-time objects.

5. Handling Duplicate Data:

- Identify and Remove Duplicates: Identifying and removing duplicate records to avoid bias and data redundancy.

6. Standardization:

- Scaling Numeric Data: Scaling numerical features to a common scale, typically between 0 and 1 or using z-score normalization.

7. Encoding Categorical Variables:

- Convert Categorical Data: Converting categorical variables into numerical representations that machine learning algorithms can work with. Common techniques include one-hot encoding or label encoding.

8. Text Cleaning (for NLP):

- Tokenization: Breaking text into individual words or tokens.
- Removing Stop Words: Eliminating common words like “the,” “is,” “and” that do not carry significant meaning.
- Lemmatization or Stemming: Reducing words to their base or root form.

Related posts:

1. What is Machine Learning ?
2. Types of Machine Learning ?
3. Applications of Machine Learning
4. Data Preprocessing
5. Handling Missing Data
6. Feature Scaling
7. Labeled data in Machine learning
8. Difference between Supervised vs Unsupervised vs Reinforcement learning
9. Machine learning algorithms for Big data
10. Difference between Supervised vs Unsupervised vs Reinforcement learning
11. What is training data in Machine learning
12. What is Ordinary Least Squares (OLS) estimation
13. Scalar in Machine Learning
14. Scalars in Loss Functions | Machine Learning
15. Linear Algebra for Machine Learning Practitioners
16. Supervised Learning
17. Top Interview Questions and Answers for Supervised Learning
18. Define machine learning and explain its importance in real-world applications.
19. Differences Between Machine Learning and Artificial Intelligence

20. Machine Learning works on which type of data ?
21. What is target variable and independent variable in machine learning
22. Machine Learning Scope and Limitations
23. What is Regression in Machine learning
24. Statistics and linear algebra for machine learning
25. Finding Machine Learning Datasets
26. What is hypothesis function and testing
27. Explain computer vision with an appropriate example
28. Explain Reinforcement learning with an appropriate exaple
29. Reinforcement Learning Framework
30. Data augmentation
31. Normalizing Data Sets in Machine Learning
32. Machine learning models
33. Unsupervised machine learning
34. Neural Network in Machine Learning
35. Recurrent neural network
36. Support Vector Machines
37. Long short-term memory (LSTM) networks
38. Convolutional neural network
39. How to implement Convolutional neural network in Python
40. What is MNIST ?
41. What does it mean to train a model on a dataset ?
42. Can a textual dataset be used with an openCV?
43. Name some popular machine learning libraries.
44. Introduction to Machine Learning
45. Some real time examples of machine learning
46. Like machine learning, what are other approaches in AI ?

47. Statistics and Linear Algebra for Machine Learning ?
48. What is convex optimization in simple terms ?
49. What is data visualization in simple terms ?
50. What is data preprocessing in machine learning ?
51. What are data distributions, and why are they important ?
52. What is data augmentation in machine learning ?
53. What is labelled and unlabelled data set in Machine Learning ?
54. What is neural networks in Machine Learning ?
55. How are convolutional neural networks related to supervised learning ?
56. Fundamentals of Neural Networks
57. Linearity vs non-linearity in Machine Learning ?
58. Machine Learning Short Exam Notes
59. Machine Learning Short Exam Notes – Quick and Easy Revision Guide