

Introduction to Data Cleaning

Data cleaning is the process of fixing or removing incorrect, corrupted, duplicate, or incomplete data.

 by Sanjay Nayak

Importance of Clean Data

1 Improved Decision Making

Clean data leads to more accurate analysis and better decision-making.

2 Increased Efficiency

It saves time and resources by preventing errors and rework.

3 Enhanced Data Quality

High-quality data improves overall operations and customer satisfaction.

Common Data Cleaning Techniques



Deduplication

Identifying and eliminating duplicate data entries.



Normalization

Organizing data to minimize redundancy and dependency.



Validation

Ensuring data accuracy and consistency.

Tools for Data Cleaning

OpenRefine

An open-source tool for working with messy data.

Talend

A comprehensive data integration platform.

Trifacta

Provides a user-friendly approach to data cleaning and wrangling.

Challenges in Data Cleaning

Data Inconsistency	Data Duplicity
Data Accuracy	Data Completeness
Data Integrity	Data Relevancy

CLEANING CHECK

ate data

ld be
e in order to
ximum value
data analysis.

es

IDs indicate
records for one
g. someone
multiple functions
e time.

Os

a labels of all
to see whether
gorical values
ealed.

Missi



Count
analyz
they ar
values
analys
results

Nume



Numer
fairly e
remov
and ma
outlier

Defin



Define
for cat
Define
numeri
Non-m
presun

Best Practices for Data Cleaning

1

Define Data Quality Criteria

Establish standards for data integrity and accuracy.

2

Regular Data Audits

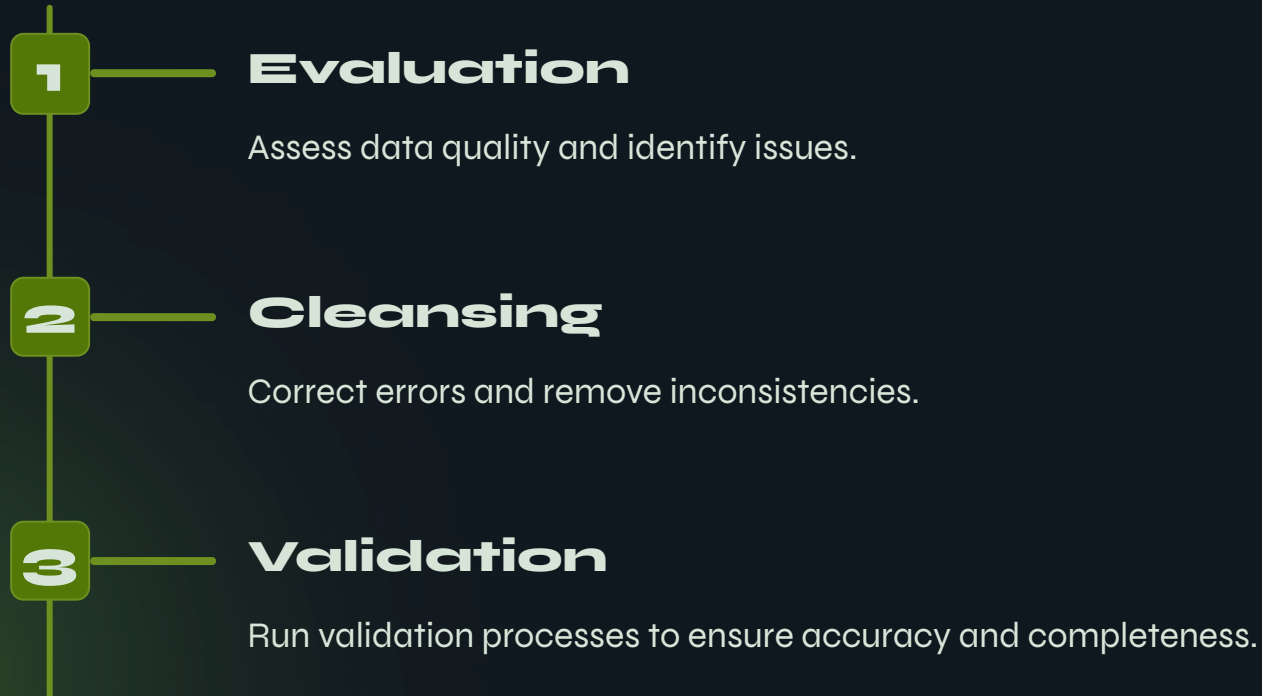
Consistently review and refine data quality processes.

3

Automate Data Cleaning

Utilize automated tools for routine cleaning tasks.

Data Cleaning Process Example



Conclusion and Next Steps

Data cleaning is vital for accurate analytics and informed decision-making. Continual improvement and automation are key for maintaining high-quality, reliable data.