Data Preparation in Data Science

Data preparation is a critical stage in the data science process. It involves a series of steps to transform raw data into a format suitable for analysis and modeling.



Data Collection

1

Sources

Data can come from multiple sources like databases, APIs, web scraping, sensors, social media, and more.



Considerations

The collection process should be efficient, reliable, and compliant with data privacy and security regulations.



Types

Data types can be structured, semi-structured, or unstructured, depending on the source and format.



Tools

Various tools and technologies are available for data collection, ranging from scripting languages to specialized platforms.

Data Discovery and

Profiling Understanding

Data
This step involves
exploring the collected
data to gain insights into
its characteristics and
patterns.

Data Profiling

This involves generating descriptive statistics, identifying data types, detecting outliers, and analyzing data distribution.

Data Quality Assessment

This helps identify
potential issues like
missing values,
inconsistent data, and
data duplication.

Data Cleansing

Missing Value

Handling massing values can involve deletion, mean/median imputation, or using machine learning techniques.

Data Standardization and Normalization

Standardization ensures all variables have the same scale, while normalization rescales data to a specific range.







Outlier Detection and Handling

Outliers can distort analysis results. They can be removed, replaced, or handled by transforming data using techniques like winsorization.

Data Consistency

Data consistency checks ensure data follows defined rules and formats, improving data accuracy and reliability.

Missing Value Handling Methods

- row deletion
- o column deletion
- mean/median/mode imputation
- o machine learning imputation
- provided constant value imputation

Example

Name	Age	Income
Alice	30	50,000
Bob	25	N/A
Charlie	35	60,000

Types of Outliers

Outliers can be classified into different types, depending on their cause and impact on the dataset. These categories include univariate outliers, multivariate outliers, and contextual outliers.

Univariate Outliers

These outliers are identified in a single variable, such as height or weight. They may be due to measurement errors or natural variation.

Multivariate

Outliers
These outliers are
identified based on their
position in multiple
dimensions, considering
the relationships between
variables. They can be hard
to detect visually.

Contextual Outliers

These outliers are unusual based on their context. For example, a high temperature reading in the Arctic would be considered an outlier.

Identifying Outliers

Several techniques can be used to detect outliers, including statistical methods and visualization tools.

Method	Description
Box Plots	Visualizes data distribution and highlights points outside the interquartile range.
Z-Scores	Measures the number of standard deviations a data point is from the mean.
IQR Rule	Identifies outliers based on the interquartile range (IQR) and its multiples.

Handling Outliers

Once outliers are identified, it's important to decide how to handle them. The approach depends on the type of outlier and its potential impact.









Deletion

Remove outliers
from the
dataset if they
are clearly
erroneous or
have a
significant
impact on the
analysis.

Transformation

Transform the data using methods like log transformation or square root transformation to reduce the impact of outliers.

Winsorizatio

Replace
extreme
outliers with the
nearest
non-outlier
value, reducing
their influence
on the data.

Imputation

Replace outliers
with imputed
values based on
other data
points,
minimizing data
loss and
preserving
information.

Standardization: Definition and

Federal Particular Standard S

Formula

z = (x - mean) / standard deviation

Where

z: standardized value, x: original value,
mean: mean of the data, standard deviation:
standard deviation of the data.

Standardization:

Example With heights of individuals measured in centimeters. After standardization, the data will be scaled to have a mean of 0 and a standard deviation of 1. This ensures that all heights are represented on a common scale, regardless of the original unit of measurement.

Original Height (cm)	Standardized Height (z-score)
170	0.5
180	1.0
165	-0.5

Normalization: Definition and

Formula is a specific range, typically between 0 and 1. This technique is useful for algorithms that are sensitive to the magnitude of features, such as distance-based algorithms like k-nearest neighbors.

Formula

x' = (x - min) / (max - min)

Where

x': normalized value, x: original value, min: minimum value in the dataset, max: maximum value in the dataset.

Normalization:

Example Entaining the prices of different products, ranging from \$10 to \$100.

Normalization rescales these prices to a range between 0 and 1, making all prices comparable on a common scale.

Original Price (\$)	Normalized Price (0-1)
10	0.0
50	0.44
100	1.0



When to Use

Stadparate Mhen data has a large variance or different scales. It helps algorithms focus on the relationship between features rather than being influenced by their magnitudes.



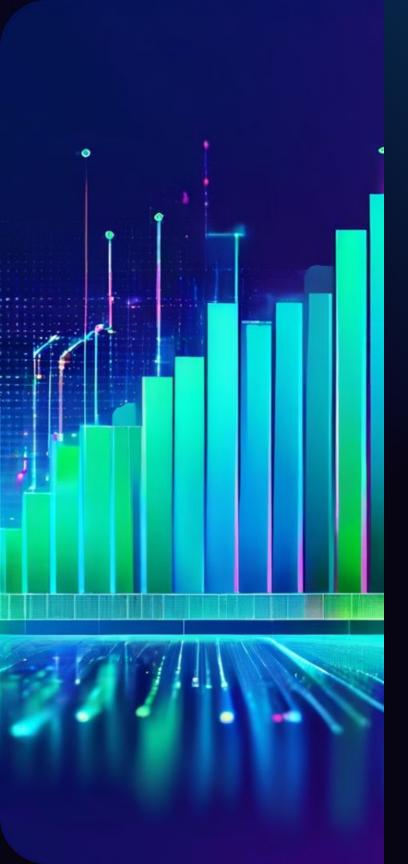
Algorithms Sensitive to

Component Analysis (PCA), Support Vector Machines (SVM).



Data with Large

Starion Sec ation helps to reduce the impact of outliers and improve the performance of algorithms.



When to Use

algorithms that rely on distance calculations or when comparing values across different scales.



Distance-Based Algerithmsighbors, K-Means

Clustering.



Neural Networks

Normalization helps prevent the vanishing gradient problem, which can hinder the training process.

Data Transformation and

Equiplement

Data Aggregation

Combining data from multiple sources or summarizing data at different levels of granularity to create new insights.

Feature

from existing data, transforming variables, or combining variables to improve model performance.

Data Encoding

Converting categorical variables into numerical formats suitable for machine learning algorithms, using techniques like one-hot encoding or label encoding.

Categorical Encoding

Original Feature

ID	Color	
1	Blue	
2	Green	
3	Red	
4	Yellow	

Label Encoding

ID	Color
1	0
2	1
3	2
4	3

One-hot Encoding

ID	Red	Blue	Green	Yellow
1	0	1	0	0
2	0	0	1	0
3	1	0	0	0
4	0	0	0	1



Data Validation

Data Integrity

Ensuring data is accurate, complete, consistent, and free from errors.

Business Rule

Verlightigata against predefined business rules to ensure it meets specific criteria and constraints.

Data Consistency

identify potential discrepancies.

Data Governance and

Matadata

Data Governance Policies and procedures for managing

data, ensuring data quality, security, and

compliance.

Metadata Management Documenting data characteristics,

relationships, and lineage to track data

provenance and facilitate data discovery

and understanding.

Data Security and Privacy Implementing measures to protect data

from unauthorized access, use, or

disclosure, ensuring compliance with

relevant regulations.

Tools and Technologies

Python Libraries

Pandas, NumPy,
Scikit-learn, and
others provide
powerful data
manipulation and
analysis
capabilities.

R Packages

Tidyverse, dplyr, and data.table offer comprehensive data preparation and analysis functions.



SQL

Structured Query
Language (SQL) is
essential for
working with
relational
databases and
extracting data for
analysis.



Cloud Platforms

AWS, Azure, and
Google Cloud offer
cloud-based data
preparation and
analysis services
with scalability and
cost-effectiveness.