

## Machine Learning life cycle

---

- ◆ **1. Problem Definition**
- ◆ **2. Data Collection**
- ◆ **3. Data Cleaning**
  - Handle:
    - Missing values
    - Duplicates
    - Inconsistent formats
    - Outliers
- ◆ **4. Data Preprocessing**
  - Convert raw data into a usable format:
    - Normalize/standardize numerical data
    - Encode categorical variables (One-hot, Label Encoding)
    - Convert dates/times
    - Feature engineering (create new features)
    - Feature selection
- ◆ **5. Exploratory Data Analysis (EDA)**
- ◆ **6. Model Selection**
- ◆ **7. Model Training**
- ◆ **8. Model Evaluation**
- ◆ **9. Hyperparameter Tuning**
- ◆ **10. Model Deployment**
- ◆ **11. Monitoring & Maintenance**

## Data cleaning and data preprocessing

### ❖ Data Cleaning: Overview, Operations, Example, and Limitations

---

#### What is Data Cleaning?

**Data cleaning** is the process of detecting and correcting (or removing) errors, inconsistencies, and inaccuracies in a dataset to improve its quality and make it suitable for analysis or modeling.

---

#### Operations in Data Cleaning

Here are common data cleaning tasks:

##### 1. Handling Missing Values

- Methods:
  - Remove rows/columns
  - Impute using:
    - Mean/Median/Mode
    - Interpolation
    - Predictive models (e.g., KNN imputation)

##### 2. Removing Duplicates

- Identifying and dropping duplicate rows using `.drop_duplicates()` in pandas.

##### 3. Fixing Data Types

- Convert incorrect data types (e.g., string to datetime, float to int).

##### 4. Handling Outliers

- Use:
  - Z-score
  - IQR (Interquartile Range)
  - Visual tools like boxplots

##### 5. Standardizing Data

- Make formats consistent (e.g., date formats, country codes, units).

##### 6. Correcting Typos or Inconsistencies

- Normalize values (e.g., unify "male", "Male", "MALE" to "Male").

## 7. Filtering Irrelevant Data

- Remove unneeded columns or rows that do not contribute to the goal.
- 

### Data Cleaning Operations: Summary Table

Cleaning Task	Methods Used	Example (Before → After)
<b>Handling Missing Values</b>	- Remove rows/columns with NaN - Impute using: <ul style="list-style-type: none"> <li>• Mean / Median / Mode</li> <li>• Interpolation</li> <li>• Predictive models (e.g., KNN)</li> </ul>	Age = [25, NaN, 35] → Age = [25, 30, 35] (filled with median = 30)
<b>Removing Duplicates</b>	- <code>.drop_duplicates()</code> in pandas	['Alice', 'Bob', 'Alice'] → ['Alice', 'Bob']
<b>Fixing Data Types</b>	- Convert types with <code>pd.to_datetime()</code> , <code>astype()</code>	'2020-01-01' (string) → Timestamp('2020-01-01')
<b>Handling Outliers</b>	- Z-score ( $\pm 3$ std dev) - IQR (Q1–1.5×IQR to Q3+1.5×IQR) - Boxplots	Salary = [50000, 1000000] → Remove 1000000 as outlier
<b>Standardizing Data</b>	- Normalize formats (dates, units, codes)	'FEMALE', 'female', 'Female' → 'Female', 'Female', 'Female'
<b>Correcting Typos/Inconsistencies</b>	- Use <code>.str.lower()</code> , <code>.str.capitalize()</code> , mapping	'Mle', 'male', 'MALE' → 'Male', 'Male', 'Male'
<b>Filtering Irrelevant Data</b>	- Drop columns with <code>.drop()</code> - Filter rows with conditions	Drop column: 'Name' Remove row where 'Salary' < 0

### Limitations of Data Cleaning

Limitation	Description
<b>Time-consuming</b>	Can take more time than model building, especially with large or messy datasets.
<b>Requires domain knowledge</b>	Some fixes need understanding of the context (e.g., imputation choices).
<b>Risk of data loss</b>	Dropping missing or duplicate values may remove valuable data.
<b>Bias introduction</b>	Improper imputation or removal of outliers can distort the dataset.
<b>Not fully automatable</b>	Some inconsistencies need human judgment (e.g., typos in text).

## Data Preprocessing: What It Is, Operations, Examples, and Limitations

---

### What is Data Preprocessing?

**Data preprocessing** is the step that comes after data cleaning. It involves **transforming the cleaned data** into a suitable format for analysis or machine learning models.

It prepares the data by scaling, encoding, and structuring it in a way that algorithms can understand and perform efficiently.

---

### Key Operations in Data Preprocessing

#### 1. Feature Scaling

- Ensures that numerical features are on a similar scale.

**Techniques:**

- Normalization (Min-Max Scaling):** Scales data to [0, 1]  
$$x' = (x - \text{min}) / (\text{max} - \text{min})$$
- Standardization (Z-score Scaling):** Centers data around mean = 0, std = 1  
$$x' = (x - \text{mean}) / \text{std}$$

#### 2. Encoding Categorical Variables

- Converts non-numeric labels into numeric form.

**Types:**

- Label Encoding** (e.g., Male = 0, Female = 1)
- One-Hot Encoding** (creates binary columns for each category)

#### 3. Feature Engineering

- Create new features or transform existing ones (e.g., extract year from date, combine features).

#### 4. Feature Selection

- Remove irrelevant or redundant features.
- Use techniques like:
  - Correlation analysis
  - Recursive Feature Elimination (RFE)
  - Feature importance from models

## 5. Data Splitting

- Divide data into:
  - **Training set**
  - **Validation set**
  - **Test set**

Typically: 70% train / 15% validation / 15% test

---

**Data Preprocessing Operations: Summary Table**

Operation	Purpose	Techniques	Example (Before → After)
Feature Scaling	Normalize numerical features to a common scale	♦ <b>Min-Max Scaling</b> ♦ <b>Z-score Standardization</b>	Salary = [30k, 50k, 100k] → Min-Max: [0.0, 0.33, 1.0] → Z-score: [-1.05, -0.18, 1.23]
Encoding Categorical Vars	Convert text labels into numbers	♦ <b>Label Encoding</b> ♦ <b>One-Hot Encoding</b>	Gender = ['Male', 'Female'] → Label: [1, 0] → One-Hot: ['Male'=1, 'Female'=0]
Feature Engineering	Derive new features from existing ones	♦ Extract year/month ♦ Combine or transform variables	'Join_Date' = '2021-01-01' → Join_Year = 2021
Feature Selection	Reduce dimensionality by removing irrelevant/redundant features	♦ <b>Correlation Matrix</b> ♦ <b>RFE (Recursive Feature Elimination)</b> ♦ <b>Model-Based Importance</b>	Selected top 3 features from 10 using ANOVA or RandomForest
Data Splitting	Divide dataset for training, validating, and testing model performance	♦ <code>train_test_split()</code> in sklearn ♦ Manual 70/15/15 or 80/20 split	Split 100 rows into: <ul style="list-style-type: none"> <li>• 70 train</li> <li>• 15 validation</li> <li>• 15 test</li> </ul>

## Limitations of Data Preprocessing

Limitation	Description
Risk of Overfitting	Poor preprocessing (e.g., encoding) may lead to models that fit training data too well.
Computational Cost	Scaling and encoding large datasets can be expensive.
Data Leakage	Must avoid using future data in transformation (e.g., fitting scaler on full data before split).
Loss of Interpretability	Transformed data may be less understandable (e.g., after PCA or encoding).
Dependence on context	Preprocessing methods vary greatly by domain (e.g., finance vs. text).

## Interview Questions:

### 1. What are the different types of missing data?

In data science, **missing data** can be categorized into **three main types**:

- ◆ **1. MCAR (Missing Completely At Random)**

- The missing values are entirely independent of any other data (observed or unobserved).
- **No pattern** in the missing data.

**Example:**

In the Titanic dataset, if some passengers randomly didn't fill in their age without any reason related to other features.

```
df['Age'].isnull().sum() # 177 missing
```

If the missing Age values are evenly distributed across Sex, Pclass, etc., it's likely **MCAR**.

---

- ◆ **2. MAR (Missing At Random)**

- The missingness is **related to other observed data**, but **not the missing value itself**.

**Example:**

Suppose younger passengers (Pclass=3) are less likely to have their Age recorded. Here, Age missingness depends on Pclass.

```
sns.boxplot(x='Pclass', y='Age', data=df)
```

This shows if Age is more frequently missing in certain classes, suggesting **MAR**.

---

- ◆ **3. MNAR (Missing Not At Random)**

- The missingness **depends on the missing value itself**.
- This is the hardest to handle and often needs domain knowledge.

**Example (Hypothetical):**

If wealthier passengers tend to withhold their Fare values because of privacy — the missingness depends on how high the fare is → **MNAR**.

---

### 🧠 Why It Matters:

Understanding the type of missing data helps you choose the right **imputation strategy**:

- **MCAR** → deletion or mean/median imputation is okay.

- **MAR** → better to use statistical or ML-based imputation.
  - **MNAR** → may need custom modeling or deeper investigation.
- 

## 2. How do you handle categorical variables?

Categorical variables are **non-numeric columns** representing categories or labels (e.g., "male/female", "yes/no", "Embarked").

Since most ML models can't work with strings directly, we **convert** categorical variables into **numeric form** using **encoding**.

### Main Techniques to Handle Categorical Variables

---

#### ◆ 1. Label Encoding

- Converts each category into a unique integer.
- Useful for **binary** or **ordinal** categories (where order matters).

**Example — Titanic Sex column:**

```
df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})  
print(df['Sex'].head())
```

**Output:**

```
0 0  
1 1  
2 1  
3 1  
4 0
```

---

#### ◆ 2. One-Hot Encoding

- Creates **new columns** for each category with 0 or 1.
- Useful for **nominal** categories (no order).

**Example — Titanic Embarked column:**

```
df = pd.get_dummies(df, columns=['Embarked'], drop_first=True)
```

## Output Columns:

Embarked\_Q Embarked\_S

0	1
0	0
1	0

Here:

- We dropped one column using `drop_first=True` to avoid **dummy variable trap** (redundancy).
  - Embarked had values: C, Q, S → we kept Q and S as new columns.
- 

## 🤔 When to Use Which?

Technique	Use When
Label Encoding	Binary or ordered categories (e.g., "low" < "medium" < "high")
One-Hot Encoding	Unordered (nominal) categories with few unique values
Target Encoding (advanced)	When category values have a meaningful relationship to the target (e.g., mean of target per category)

---

## 🧠 Why Encoding Is Important?

- ML models (like logistic regression, SVM, etc.) require numeric input.
  - Encoding prevents misinterpretation of categories as text.
  - Ensures better model accuracy and interpretability.
- 

## 3. What is the difference between normalization and standardization?

Both **normalization** and **standardization** are feature scaling techniques used to bring numeric columns to a similar scale — which is essential for machine learning models (especially distance-based ones like KNN, SVM, and Gradient Descent optimizers).

---

### ◆ Normalization (Min-Max Scaling)

- **Definition:** Scales data between **0 and 1**.
- **Formula:**

$$x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

- **Use Case:** When you want **bounded** values (e.g., neural networks).

**Example (Titanic Fare column):**

python

```
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()
df['Fare_normalized'] = scaler.fit_transform(df[['Fare']])
print(df[['Fare', 'Fare_normalized']].head())
```

 **Output:**

markdown

	Fare	Fare_normalized
0	7.2500	0.01415
1	71.2833	0.13914
2	7.9250	0.01566
..		

◆ **Standardization (Z-score Scaling)**

- **Definition:** Transforms data to have **mean = 0** and **standard deviation = 1**.
- **Formula:**

$$x' = \frac{x - \mu}{\sigma}$$

- **Use Case:** When data follows **normal distribution** or outliers are present.

**Example (Titanic Age and Fare):**

python

```
from sklearn.preprocessing import StandardScaler

scaler = StandardScaler()
df[['Age', 'Fare']] = scaler.fit_transform(df[['Age', 'Fare']])
print(df[['Age', 'Fare']].head())
```

 **Output:**

markdown

	Age	Fare
0	-0.565736	-0.502445
1	0.663861	0.786845
2	-0.258337	-0.488854
..		

## Summary Comparison Table

Feature	Normalization	Standardization
Scale Range	0 to 1	Mean = 0, Std = 1
Formula	$(x - \min) / (\max - \min)$	$(x - \text{mean}) / \text{std}$
Sensitive to Outliers	Yes	Less sensitive
Use Case	Neural nets, bounded input	SVM, KNN, linear models
Example Tool	MinMaxScaler()	StandardScaler()

### Key Insight:

Always **fit the scaler only on training data** and **transform both training and test data** to avoid data leakage.

---

## 4. How do you detect outliers?

Outliers are data points that **deviate significantly** from the rest of the dataset. They can **skew** the results of machine learning models, especially those that rely on averages or distances (e.g., linear regression, KNN).

---

### Common Outlier Detection Methods

---

#### ◆ 1. Boxplot (IQR Method — Most Used)

**How it works (behind the scenes):**

- It uses the **Interquartile Range (IQR)**, which is the range between the **25th (Q1)** and **75th percentile (Q3)**.
- Any data point outside this range is considered an outlier.

#### Formulas:

$$IQR = Q3 - Q1$$

$$\text{Lower bound} = Q1 - 1.5 \times IQR$$

$$\text{Upper bound} = Q3 + 1.5 \times IQR$$

### Example (Titanic Fare column):

```
import seaborn as sns
import matplotlib.pyplot as plt

# Boxplot for visual
sns.boxplot(x=df['Fare'])
plt.title("Boxplot for Fare (Outliers shown as dots)")
plt.show()

# Behind the scenes IQR calculation
Q1 = df['Fare'].quantile(0.25)
Q3 = df['Fare'].quantile(0.75)
IQR = Q3 - Q1

lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

outliers = df[(df['Fare'] < lower_bound) | (df['Fare'] > upper_bound)]
print("Number of outliers:", len(outliers))
```

### Output:

Number of outliers: 116

These 116 records have Fare values that are **either too low or too high** compared to the majority.

---

### ◆ 2. Z-score Method (Standard Deviation)

#### How it works (behind the scenes):

- Measures how far a point is from the mean in terms of standard deviations.
- Data points with a **Z-score > 3 or < -3** are considered outliers.

### Example (Titanic Age column):

```
python

from scipy.stats import zscore

df['Age_z'] = zscore(df['Age'])
outliers = df[(df['Age_z'] > 3) | (df['Age_z'] < -3)]
print("Age outliers using Z-score:", len(outliers))
```

### Output:

```
cpp

Age outliers using Z-score: 4
```

---

### ◆ 3. Visualization Techniques

- **Boxplot:** Quick visual indicator of spread and outliers.
  - **Scatterplots:** Good for multivariate outliers.
  - **Histogram/Distplot:** Helps spot skewed distributions.
- 

### ⌚ Which One Should You Use?

Method	Use Case	Pros	Cons
IQR	Skewed data or with outliers	Easy, no assumption of normality	Doesn't scale well to high-dim data
Z-score	Normally distributed data	Statistically sound	Sensitive to mean and std
Boxplot	Exploratory data analysis (EDA)	Quick visual check	Not for multivariate outliers

---

### 🧠 Real-World Use:

Outliers are often **not just bad data** — they can represent rare but important events (e.g., VIP passengers in Titanic). Always analyze **context** before removing them.

---

## 5. Why is preprocessing important in Machine Learning (ML)?

**Preprocessing** is the foundation of a successful machine learning pipeline. Raw data is often messy, inconsistent, and contains irrelevant patterns. If you skip preprocessing, even the most powerful model will **underperform or give misleading results**.

---

### ✿ Key Reasons Why Preprocessing Is Crucial

---

#### ◆ 1. Handles Missing Values

Most ML algorithms **can't handle nulls**. Missing data needs to be imputed or removed to ensure training doesn't fail.

#### Example (Titanic Dataset):

```
df['Age'].isnull().sum() # Returns 177
```

```
df['Age'].fillna(df['Age'].median(), inplace=True)
```

If you skip this, models like LogisticRegression() will throw an error.

---

## ◆ 2. Encodes Categorical Variables

ML models need **numeric input**. Categorical features (like "Sex", "Embarked") must be converted.

**Example:**

```
df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})  
df = pd.get_dummies(df, columns=['Embarked'], drop_first=True)
```

Skipping this would cause the model to crash or misinterpret the data.

---

## ◆ 3. Scales Numerical Data

Features like Fare and Age might have different ranges. Models based on **distance or gradient descent** (e.g., KNN, SVM, Neural Nets) perform poorly without scaling.

**Example:**

```
from sklearn.preprocessing import StandardScaler  
scaler = StandardScaler()  
df[['Age', 'Fare']] = scaler.fit_transform(df[['Age', 'Fare']])
```

This avoids one feature dominating the learning process due to a larger scale.

---

## ◆ 4. Detects and Handles Outliers

Outliers can **distort model weights** or increase error. Preprocessing helps find and fix them.

**Example:**

```
# Remove extreme Fare values using IQR  
Q1 = df['Fare'].quantile(0.25)  
Q3 = df['Fare'].quantile(0.75)  
IQR = Q3 - Q1  
df = df[(df['Fare'] >= Q1 - 1.5*IQR) & (df['Fare'] <= Q3 + 1.5*IQR)]
```

---

## ◆ 5. Improves Model Accuracy & Stability

Well-preprocessed data leads to:

- Faster convergence
  - Less overfitting
  - Better generalization
- 

#### ◆ 6. Reduces Noise and Irrelevance

Unnecessary features or inconsistencies can **confuse the model**. Preprocessing allows:

- Dropping irrelevant columns
  - Cleaning corrupted text/data
  - Consistent formatting
- 

#### 📌 Summary Table

Problem	Preprocessing Solution
Missing values	Imputation or removal
Categorical variables	Label/One-hot encoding
Unequal feature scales	Standardization/Normalization
Outliers	IQR or Z-score detection
Noisy or messy data	Cleaning/formatting

Model failure or poor accuracy Handled via full pipeline

---

## 6. What is One-Hot Encoding vs Label Encoding?

Both **One-Hot Encoding** and **Label Encoding** are techniques used to convert **categorical (non-numeric)** data into **numerical format**, which is necessary for machine learning models.

---

#### ◆ Label Encoding — Logical View

#### 📌 What it does:

- Assigns a **unique integer to each category**.
- Transforms categories into **ordinal numbers** (e.g., 0, 1, 2, ...).

## Logical Interpretation:

Label encoding introduces **implicit ordering**, even if no natural order exists.

### Example — Sex column in Titanic dataset:

```
df['Sex'] = df['Sex'].map({'male': 0, 'female': 1})
```

#### Result:

male → 0

female → 1

 **Issue:** This is okay for binary columns, but for City = [Delhi, Mumbai, Chennai], encoding them as 0, 1, 2 falsely implies *Delhi < Mumbai < Chennai*, which has **no logical meaning**.

---

### ◆ One-Hot Encoding — Logical View

#### What it does:

- Creates a **separate binary column for each category**.
- Only one column is '1' per row; rest are '0'.

## Logical Interpretation:

One-hot encoding treats categories as **independent and unrelated** — which is often correct.

### Example — Embarked column in Titanic:

Values: C, Q, S

```
df = pd.get_dummies(df, columns=['Embarked'], drop_first=True)
```

#### Result Columns:

Embarked_Q	Embarked_S
0	1
1	0
0	0 (means C)

 This avoids the **false ordering** and lets the model learn **true patterns** in the data.

## Key Differences (Logically)

Feature	Label Encoding	One-Hot Encoding
Format	Integer values (0, 1, 2...)	Binary columns (0 or 1)
Category assumption	Assumes <b>order</b> (even if none)	Assumes <b>no order</b> (correct)
Columns added	1	n - 1 or n (if drop_first=False)
Memory efficient?	Yes (fewer columns)	No (more columns for high cardinality)
When to use	Binary or <b>ordinal</b> categorical	<b>Nominal</b> (non-ordinal) categorical

## Which one to use when?

Case	Use
Binary category (male/female)	Label or One-hot
Ordered category (low < med < high)	Label Encoding
Unordered categories (red, blue, green)	One-Hot Encoding
Many categories (>100)	Consider <b>Target Encoding or Embedding</b>

## Warning:

Using **Label Encoding** on **unordered categories** can trick ML models (especially tree-based or linear models) into thinking some values are more important or ranked.

## 7. How do you handle data imbalance?

**Data imbalance** occurs when the number of samples in one class greatly outnumbers the other(s) — for example, in the Titanic dataset, more people died (0) than survived (1). This can cause models to become biased toward the majority class, ignoring the minority class.

## Techniques to Handle Data Imbalance

## ◆ 1. Resampling Techniques

### a. Oversampling

- Increases the number of minority class samples by duplicating or synthetically generating them (e.g., using SMOTE).
- Helps balance the dataset and improves recall for minority class.

### b. Undersampling

- Reduces the number of majority class samples by removing examples.
  - Simple but can discard useful data.
- 

## ◆ 2. Use Class Weights

- Many ML algorithms allow you to assign **more weight to minority class samples** during training.
  - Makes the model pay more attention to the underrepresented class.
- 

## ◆ 3. Use Appropriate Evaluation Metrics

- Instead of relying on accuracy, use:
    - **Precision**
    - **Recall**
    - **F1-Score**
    - **ROC-AUC**
  - These give better insight into model performance on the minority class.
- 

## ◆ 4. Generate Synthetic Samples (SMOTE)

- SMOTE (Synthetic Minority Over-sampling Technique) creates **synthetic data points** for the minority class by interpolating between existing examples.
  - More powerful than simple duplication.
- 

## ◆ 5. Ensemble Methods

- Use **bagging or boosting** (e.g., Random Forest, XGBoost) which are more robust to class imbalance.
- Some algorithms like **BalancedRandomForestClassifier** are designed for imbalance.

---

## 8. Can preprocessing affect model accuracy?

**Yes — preprocessing can significantly affect model accuracy.**

It plays a crucial role in shaping the quality of the data, which directly impacts how well a model can learn and make predictions.

---

### Logic Behind Why Preprocessing Affects Accuracy

---

#### 1. Garbage In, Garbage Out (GIGO)

If your input data is noisy, inconsistent, or poorly formatted, your model learns from flawed information — leading to poor accuracy.

---

#### 2. Handling Missing Values

- If missing values are left untreated, many algorithms will crash or behave unpredictably.
  - **Imputation (mean/median/fill)** provides complete data for the model to train on effectively.
- 

#### 3. Encoding Categorical Variables

- Algorithms can't understand text — so categories must be **converted to numbers**.
  - Wrong encoding (e.g., using label encoding on non-ordinal data) can confuse the model by implying an order that doesn't exist.
- 

#### 4. Feature Scaling (Normalization/Standardization)

- Many ML algorithms (like KNN, SVM, Logistic Regression) are sensitive to **feature magnitudes**.
  - If features are not scaled, models might favor larger-valued features, leading to biased predictions and reduced accuracy.
- 

#### 5. Outlier Removal

- Outliers can skew the model's understanding of normal patterns.
  - Removing or capping outliers makes the model **more stable and generalizable**.
-

## 6. Data Imbalance Handling

- If the model sees mostly one class, it may just predict that class every time.
  - Proper balancing (e.g., SMOTE, class weights) ensures the model learns to recognize **both classes**, improving **recall, precision, and accuracy**.
- 

## 7. Noise Reduction

- Cleaning irrelevant or misleading data points allows the model to focus on actual patterns.
  - This improves **training efficiency** and final performance.
- 

### Summary

Preprocessing is not just a formality — it **directly affects the quality of features** that a model uses. Well-preprocessed data gives the model **clearer patterns, better generalization, and ultimately higher accuracy**.