

# Thinking: The Core of Machine Learning

Before diving into data collection and analysis, critical thinking is essential in machine learning. Understanding the problem, defining objectives, and formulating the right questions are key steps in building a successful model.

1. **Problem Understanding** – Clearly define the business or research problem before looking for data.
2. **Hypothesis Formation** – Think about potential factors that may influence the target variable.
3. **Feature Relevance** – Consider which types of data are meaningful and how they might impact predictions.
4. **Data Limitations** – Evaluate potential biases, missing values, and inconsistencies in available datasets.
5. **Model Feasibility** – Think about which models would be best suited for the data and task at hand.

The ability to think critically about data and its real-world implications is what separates an average machine learning project from a truly effective one.

# Data Analysis

First, we need to determine our **main objective or target** for the project. Once we have defined the goal, the next step is to search for a relevant dataset.

[Kaggle](#) is one of the best sources for high-quality datasets. We can search for datasets there based on our target. Once we find a potential dataset, it's crucial to **validate its quality and relevance** to our goal. This involves checking for factors like:

- **Completeness** of data.
- **Correctness** of data (whether it represents the problem correctly).
- **Size** of the dataset (whether it is large enough to train a model).
- **Data type compatibility** (whether the features align with the type of analysis you plan to do).

If the dataset meets these requirements, we can download it. If it's not suitable, we either need to find another one or, in some cases, **collect additional data** from other sources. After ensuring the dataset is valid, we can move on to the next step, which is **data preprocessing**.

# Data Visualization

First, we will import the necessary libraries, such as **pandas** for data manipulation and **matplotlib/seaborn** for visualization. And load the dataset.

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
import seaborn as sns
```

```
# Load the dataset
```

```
df = pd.read_csv("data.csv")
```

Next, we will visualize each feature in the dataset to understand any **patterns or relationships** between them. I understand that for beginners, it might be difficult to grasp why **data visualization** is necessary. Let me explain:

## Why Visualize the Data?

- **Identify trends and patterns:** Visualizations help reveal hidden patterns and trends in the data that are difficult to detect just by looking at the raw numbers.
- **Spot outliers:** It's easier to spot extreme values (outliers) visually than numerically.
- **Understand distributions:** Plots like histograms or box plots can help us understand the distribution of data (normal, skewed, etc.).
- **Relationship between variables:** Visualizing data can show how different features correlate with each other. For instance, scatter plots can reveal if two features have a linear or non-linear relationship.
- **Improve decision-making:** Data visualization supports better and more informed decision-making because it makes complex data more accessible and easier to interpret.

## Example Code to Visualize Data:

Let's look at a **scatter plot** to understand the relationship between two numerical variables (e.g., "age" and "salary").

```
# Scatter plot example  
  
sns.scatterplot(x="age", y="salary", data=df)  
  
plt.title('Age vs. Salary')  
  
plt.show()
```

Another useful plot is the **histogram**, which shows the distribution of a single variable (e.g., "age").

```
# Histogram for age distribution  
  
sns.histplot(df['age'], kde=True)  
  
plt.title('Age Distribution')  
  
plt.show()
```

## Types of Plots:

- **Scatter Plot:** Useful for visualizing relationships between two continuous variables.
- **Histogram:** Helps visualize the distribution of a single variable.
- **Box Plot:** Useful for identifying outliers and understanding the distribution.
- **Bar Plot:** Great for comparing categorical data.
- **Heatmap:** Displays correlations between different features.

Here is an example of all the parameters in **SeaBorn**

```
sns.barplot(x=None, y=None, data=None, hue=None, order=None, hue_order=None,  
            estimator=<function: mean>, ci="sd", n_boot=1000, units=None,  
            orient=None, color=None, palette=None, errcolor=None, errwidth=None,  
            capsize=None, dodge=True, width=0.8, dodge=None)
```

## Matplotlib:

```
plt.bar(x, height, width=0.8, bottom=None, align='center', data=None, color='blue',  
        edgecolor='black', linewidth=1, tick_label=None, log=False, orientation='vertical')
```

# Data Preprocessing & Feature Engineering

Before applying machine learning models, we must clean and prepare our dataset. Raw data often contains **missing values, duplicates, outliers, categorical variables, and inconsistent formats** that must be handled properly.

## 1. Understanding the Dataset

Before preprocessing, we must understand the data:

```
import pandas as pd
```

```
# Load dataset
```

```
df = pd.read_csv("data.csv")
```

```
# Check first few rows
```

```
print(df.head())
```

```
# Check column information
```

```
print(df.info())
```

```
# Check for missing values
```

```
print(df.isnull().sum())
```

# Check for duplicate values

```
print(df.duplicated().sum())
```

# Check statistical summary of numerical features

```
print(df.describe())
```

## 2. Handling Missing Values

### A. Removing Missing Values

If the missing data is too much and cannot be inferred, we **drop the rows or columns**.

# Drop rows with missing values

```
df.dropna(inplace=True)
```

# Drop columns with too many missing values

```
df.dropna(axis=1, inplace=True)
```

### B. Filling Missing Values (Imputation)

Instead of removing missing values, we can **fill** them with appropriate values.

**For Numerical Data:**

- **Mean Imputation** (for normally distributed data)
- **Median Imputation** (for skewed data)
- **Mode Imputation** (for categorical-like numerical data)

# Fill missing values with mean

```
df['column_name'].fillna(df['column_name'].mean(), inplace=True)
```

# Fill missing values with median

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

# Fill missing values with mode (most frequent value)

```
df['column_name'].fillna(df['column_name'].mode()[0], inplace=True)
```

**For Categorical Data:**

- Replace missing values with **mode** (most frequent category).
- If the category has **too many missing values**, we create a new category like **"Unknown"**.

```
df['category_column'].fillna('Unknown', inplace=True)
```

### 3. Handling Duplicate Values

Duplicate values can introduce bias in the model. We should remove them before training.

# Check for duplicates

```
print(df.duplicated().sum())
```

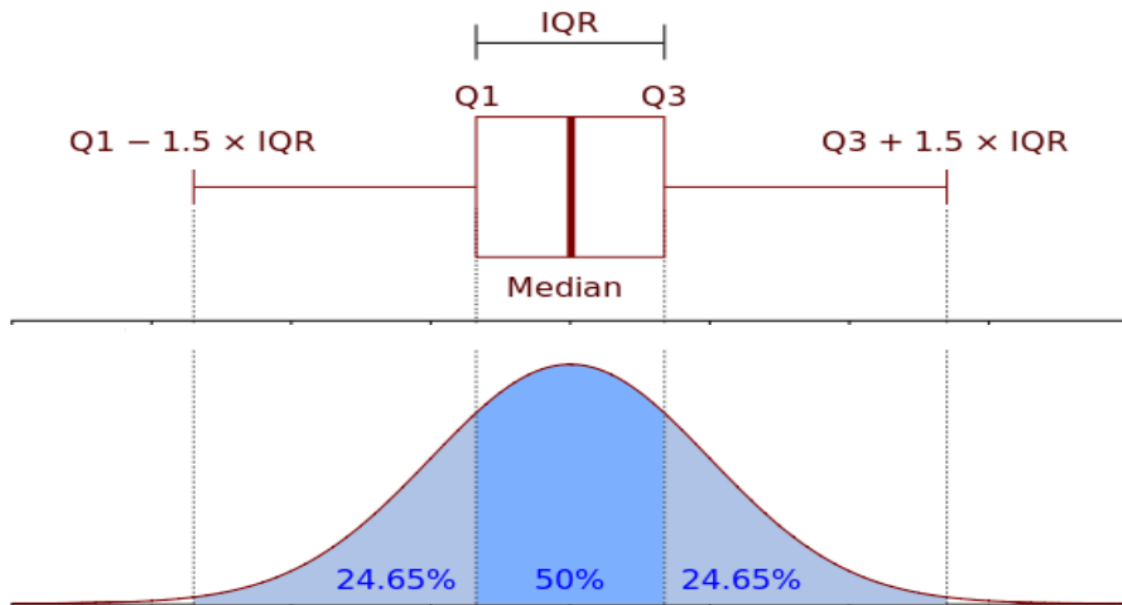
# Remove duplicates

```
df.drop_duplicates(inplace=True)
```



## 4. Handling Outliers

Outliers can **skew model predictions** and must be handled carefully. Here is visual of IQR.



### A. Detecting Outliers (Using IQR Method)

```
import numpy as np
```

```
# Compute IQR (Interquartile Range)
```

```
Q1 = df['feature'].quantile(0.25)
```

```
Q3 = df['feature'].quantile(0.75)
```

```
IQR = Q3 - Q1
```

```
# Define outlier boundaries
```

```
lower_bound = Q1 - 1.5 * IQR
```

```
upper_bound = Q3 + 1.5 * IQR
```

```
# Identify outliers
```

```
outliers = df[(df['feature'] < lower_bound) | (df['feature'] > upper_bound)]
```

```
print(outliers)
```

## B. Removing or Transforming Outliers

- **Remove outliers** if they are due to incorrect data entry.
- **Cap outliers** (set extreme values to lower/upper bound).
- **Use log transformation** if the feature is highly skewed.

```
# Capping outliers
```

```
df['feature'] = np.where(df['feature'] > upper_bound, upper_bound,  
                        np.where(df['feature'] < lower_bound, lower_bound, df['feature']))
```

## 5. Handling Categorical Data

Machine learning models work with numerical data, so categorical variables must be converted into numbers.

### A. Label Encoding (for ordinal data)

Used when categories have a meaningful order (Low < Medium < High).

```
from sklearn.preprocessing import LabelEncoder
```

```
# Apply Label Encoding
```

```
encoder = LabelEncoder()
```

```
df['category_column'] = encoder.fit_transform(df['category_column'])
```

## B. One-Hot Encoding (for nominal data)

Used when categories have **no specific order** (e.g., Red, Blue, Green).

```
from sklearn.preprocessing import OneHotEncoder
```

```
# One-hot encode categorical column
```

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

## 6. Feature Scaling (Standardization & Normalization)

Some models (like logistic regression, KNN, and SVM) perform better when features are **scaled**.

### A. Standardization (StandardScaler)

- Converts data to **mean = 0** and **standard deviation = 1**.
- Useful for models that assume normally distributed data (e.g., logistic regression, SVM).

```
from sklearn.preprocessing import StandardScaler
```

```
scaler = StandardScaler()
```

```
df[['feature1', 'feature2']] = scaler.fit_transform(df[['feature1', 'feature2']])
```

### B. Normalization (MinMaxScaler)

- Scales values between **0 and 1**.
- Useful for distance-based models (e.g., KNN, neural networks).

```
from sklearn.preprocessing import MinMaxScaler

scaler = MinMaxScaler()

df[['feature1', 'feature2']] = scaler.fit_transform(df[['feature1', 'feature2']])
```

## 7. Feature Engineering

**Definition:** Feature engineering is the process of transforming raw data into meaningful features that improve machine learning model performance. It involves creating, modifying, selecting, and encoding features.

### 1. Why is Feature Engineering Important?

- ✓ A well-engineered dataset boosts accuracy and helps models find patterns easily.
- ✓ Poor feature selection adds noise, leading to overfitting or underfitting.
- ✓ Good features can reduce the need for complex models.

### 2. How to Decide What Features to Add or Remove?

Before adding or modifying features, consider:

#### A. Understanding the Dataset

- Check the domain knowledge (e.g., if predicting house prices, features like "number of rooms" or "location" matter).
- Analyze the data distribution (Are some features irrelevant or highly correlated?).
- Identify missing values and decide whether to drop or impute them.

#### B. Using Feature Importance

- Run correlation analysis to check relationships between features.

- Use machine learning models like Decision Trees, Random Forest, or SHAP to see which features are most important.
- Remove features with low variance (features that remain mostly the same across all data points).

### C. When to Create a New Feature?

- When an interaction between features improves the model (e.g., *"age squared"* in a salary prediction model).
- When domain knowledge suggests a derived metric (e.g., *"BMI"* from height and weight).
- When encoding categorical variables (e.g., *One-Hot Encoding* or *Label Encoding*).

## 3. Types of Feature Engineering

### A. Handling Missing Values (examples in Data Preprocessing No. 2)

**Drop columns** if they contain too many missing values.

**Impute values** using:

- Mean/median/mode for numerical data.
- Forward-fill/backward-fill for time series.
- KNN imputation for better accuracy.

### B. Encoding Categorical Variables (examples in Data Preprocessing No. 5)

- **One-Hot Encoding (OHE)** – Used for nominal categories (e.g., "Color": Red, Blue, Green).
- **Label Encoding** – Used for ordinal categories (e.g., "Size": Small, Medium, Large).
- **Target Encoding** – Maps categorical values to their target mean in classification problems.

### C. Feature Scaling (Normalization & Standardization) (examples in DP No.6)

- **Normalization (Min-Max Scaling)** – Scales values between 0 and 1. Best for deep learning models.

- **Standardization (Z-score scaling)** – Centers the data with mean = 0 and std = 1.  
Best for SVM, Logistic Regression, and PCA.

## D. Creating New Features (Feature Generation)

### 1. Polynomial Features

- Creating **squared, cubic, or interaction terms** to capture non-linear relationships.
- Often useful for models like **linear regression** when relationships aren't strictly linear.

- ♦ **Example:** Adding  $x_2^2$  and  $x_3^2$  features

```
from sklearn.preprocessing import PolynomialFeatures

import numpy as np

X = np.array([[2], [3], [4]])

poly = PolynomialFeatures(degree=2, include_bias=False) # Generates  $x$  and  $x^2$ 

X_poly = poly.fit_transform(X)

print(X_poly)
```

### 2. Binning (Discretization)

- Converts **continuous variables** into categorical bins.
- Helps models like **decision trees** or **logistic regression** when data distribution is skewed.
- Example: Instead of using **exact age**, we group them into bins like "young," "middle-aged," and "old."

- ♦ **Example:**

```
import pandas as pd

df = pd.DataFrame({'Age': [15, 22, 35, 50, 65]})

df['AgeGroup'] = pd.cut(df['Age'], bins=[0, 18, 35, 50, 100],
labels=['Teen', 'Young Adult', 'Middle-Aged', 'Senior'])

print(df)
```

### 3. Date-Based Features

- Extracts information from **timestamps** (day, month, year, hour, weekday, weekend, holidays).
- Essential for **time-series models** and business forecasting.

#### ♦ Example: Extracting features from a timestamp

```
df['Date'] = pd.to_datetime(df['Date']) # Convert to datetime format

df['Year'] = df['Date'].dt.year

df['Month'] = df['Date'].dt.month

df['Weekday'] = df['Date'].dt.weekday # Monday=0, Sunday=6

df['IsWeekend'] = df['Weekday'].apply(lambda x: 1 if x >= 5 else 0)

print(df.head())
```

### 4. Aggregations (Grouping Data)

- Used in **time-series** and **grouped analysis** (e.g., average customer purchase per month).
- Helps in **reducing dimensionality** while keeping important patterns.

♦ **Example: Average purchase amount per customer**

```
df_grouped = df.groupby('CustomerID').agg({'PurchaseAmount': ['sum',  
    'mean', 'count']})  
  
df_grouped.columns = ['TotalSpent', 'AverageSpent', 'NumberOfPurchases']  
  
print(df_grouped.head())
```

## 5. Log Transformation (Handling Skewed Data)

- Converts skewed distributions to **more normal distributions**.
- Useful for **linear models** that assume normality.
- Often applied to features like **income, house prices, and population size**.

♦ **Example: Applying Log Transformation**

```
import numpy as np  
  
df['Log_Salary'] = np.log1p(df['Salary']) # log(1 + x) to avoid log(0)  
error
```

## 6. Text-Based Feature Extraction

- Converts **text data** into **numeric representations** (word counts, TF-IDF, embeddings).

♦ **Example: Converting text to word count features**

```
from sklearn.feature_extraction.text import CountVectorizer  
  
text_data = ["Machine learning is great", "I love machine learning"]  
  
vectorizer = CountVectorizer()
```



```
X = vectorizer.fit_transform(text_data)

print(pd.DataFrame(X.toarray(),
columns=vectorizer.get_feature_names_out()))
```

## 7. Interaction Features

- Captures **relationships** between two or more variables.
- Often used when feature importance suggests two variables interact.

### ♦ Example: Multiplying two features to create a new one

```
df['Income_to_Debt_Ratio'] = df['Income'] / df['Debt']

df['Experience_Level'] = df['YearsExperience'] * df['EducationYears']
```

## 8. Frequency Encoding

- Replaces **categorical values** with the frequency of their occurrences.
- Works well when **certain categories appear frequently** and have a direct impact on predictions.

### ♦ Example: Frequency encoding of categorical variable

```
df['Category_Freq'] = df['Category'].map(df['Category'].value_counts() /
len(df))
```

## 9. Target Encoding (Mean Encoding)

- Replaces categories with the **mean of the target variable**.
- Useful for **ordinal categorical features** and helps with **predictive power**.

- ◆ **Example: Encoding 'City' based on average house price**

```
df['City_encoded'] = df.groupby('City')['HousePrice'].transform('mean')
```

## 10. Lag Features (For Time-Series Data)

- Uses **past values** to predict future values in time-series data.
- Common in **stock market prediction, demand forecasting, and anomaly detection**.

- ◆ **Example: Creating a lag feature of 1 day**

```
df['Sales_Lag1'] = df['Sales'].shift(1)
```

```
df['Sales_Lag7'] = df['Sales'].shift(7) # Previous week's sales
```

## 11. Rolling Window Features (Moving Averages)

- Uses a **rolling window** to calculate moving averages, standard deviations, etc.
- Useful in **time-series forecasting** (e.g., stock price moving averages).

- ◆ **Example: Creating a rolling mean of last 7 days**

```
df['Sales_Rolling7'] = df['Sales'].rolling(window=7).mean()
```

## 12. Ratio-Based Features

- Instead of using raw values, creating **ratios** often makes more sense (e.g., **income-to-loan ratio** is more predictive than just income).

- ◆ **Example: Creating a ratio between two variables**

```
df['Debt_to_Income'] = df['Debt'] / df['Income']

df['Price_per_SquareFoot'] = df['HousePrice'] / df['HouseSize']
```

### 13. Clustering-Based Features

- Assigns a new **cluster label** using K-Means clustering, allowing models to capture **group behavior**.
- Works well when **categorical variables are not well-defined**.

#### ♦ Example: Assigning cluster labels using KMeans

```
from sklearn.cluster import KMeans

import numpy as np

X = np.array([[1, 2], [2, 3], [3, 4], [8, 9], [9, 10]]) # Example data

kmeans = KMeans(n_clusters=2)

df['Cluster'] = kmeans.fit_predict(X)
```

### 14. Text Sentiment Analysis as a Feature

- Converts **text reviews or social media posts** into a **sentiment score**.
- Used in **customer review analysis, fraud detection, and recommendation systems**.

#### ♦ Example: Extracting sentiment score from text

```
from textblob import TextBlob

df['Sentiment'] = df['Review'].apply(lambda x:
TextBlob(x).sentiment.polarity)
```

## 15. Image-Based Feature Engineering

- Extracts features like **color histograms, texture, and edges** from images.
- Used in **computer vision models**.

### ♦ Example: Extracting average pixel intensity from images

```
import cv2

img = cv2.imread("image.jpg", cv2.IMREAD_GRAYSCALE)

avg_pixel = img.mean()

print("Average Pixel Intensity:", avg_pixel)
```

## 4. Checking Feature Importance (Feature Selection) (we will see this in future)

Once features are created, we need to **eliminate irrelevant ones** to improve efficiency.

### Methods to Check Feature Importance

- **Correlation Matrix** – Identify highly correlated features and remove redundancy.
- **Mutual Information** – Measures how much information one variable provides about another.
- **Tree-Based Models** – Feature importance scores from Decision Trees or Random Forest.
- **Lasso Regression (L1 Regularization)** – Shrinks less important features to zero.

## 8. Data Splitting

Before training a model, we need to **split the dataset** into training and testing sets.

```
from sklearn.model_selection import train_test_split

X = df.drop('target_column', axis=1) # Features
```

```
y = df['target_column'] # Target variable

# Split data (80% training, 20% testing)

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)
```

## 9. Data Imbalance Handling (For Classification Problems)

If the dataset has **imbalanced classes**, some techniques can help:

### A. Undersampling (Reducing Majority Class)

```
from imblearn.under_sampling import RandomUnderSampler

undersample = RandomUnderSampler(sampling_strategy='auto')

X_resampled, y_resampled = undersample.fit_resample(X, y)
```

### B. Oversampling (Increasing Minority Class)

```
from imblearn.over_sampling import SMOTE

smote = SMOTE(sampling_strategy='auto')

X_resampled, y_resampled = smote.fit_resample(X, y)
```

## Advantages & Disadvantages of Resampling Methods

## ✓ Advantages

- ✓ **Better Model Performance:** More balanced predictions.
- ✓ **Higher Recall for Minority Class:** Helps detect rare cases.
- ✓ **Reduces Model Bias:** Ensures both classes are learned properly.

## ✗ Disadvantages

- ✗ **Oversampling Can Cause Overfitting:** The model memorizes the synthetic data instead of learning patterns.
- ✗ **Undersampling Can Lose Important Data:** Reducing the majority class may remove useful information.
- ✗ **Longer Training Time:** With oversampling, more data means slower training.

## Uses of Resampling

Some models handle imbalance better than others.

### ✓ Tree-Based Models (Best for Imbalance)

- Decision Trees: Can naturally handle imbalance but may overfit.
- Random Forest: Can work well but needs class weighting.
- XGBoost / LightGBM: Work well with imbalance adjustments like `scale_pos_weight`.

### ✓ Logistic Regression (With Class Weights)

- Assigns a higher weight to the minority class.
- Works well when data is linearly separable.

### ✓ Anomaly Detection Models (Best for Extreme Imbalance)

- Isolation Forest: Works when the minority class is very rare and different.
- One-Class SVM: Best for extreme cases (e.g., rare fraud detection).