# Introduction to Machine Learning

## The Fundamentals of Data in ML

Lecture 3

Spring 2024
Hemin Ibrahim, PhD
hemin.ibrahim@tiu.edu.iq

# Outline

- What is Data?
- Why data in ML?
- Types of data
- Data representation
- Using Data in Machine Learning
- Feature Extraction

# Objectives

- Understand the fundamental concept of data
- Distinguish between the core categories of structured and unstructured data.
- Recognize the importance of choosing suitable data representations for effective machine learning algorithms.
- Understand the idea that machine learning algorithms fundamentally rely on data patterns to learn
- Understand the crucial step of splitting data into training, validation, and testing sets to ensure accurate model evaluation
- Understand that Feature extraction extracts essential data points for better machine learning.

The concept of machine learning as algorithms "learning" from "data" rather than being explicitly programmed for everything.

# What is Data?

## "data"

Data is a collection of raw facts, figures, measurements, or observations.
In machine learning, it's the information used to train the algorithms.

- Numerical values (e.g., house prices, test scores)
- Text (e.g., product reviews, tweets)
- Images (e.g., medical scans, photographs)
- Audio (e.g., speech recordings, music)
- Video (eg., movie, a tutorial, a recorded live event)

# Data

# "data"
is **power**, **control**, **money**, **dominance**, ..



PCWorld

NEWS ⌄   BEST PICKS ⌄   REVIEWS ⌄   HOW-TO ⌄   DEALS ⌄   LAPTOPS   WINDOWS   SECURITY   MORE ⌄   NEWSLETTERS

## Time to update your password: 26 billion personal records just leaked

The data set clocks in at a massive 12TB.

By Alaina Yee
Senior Editor, PCWorld  |  JAN 24, 2024 11:16 AM PST

# Why Data in ML?

- Data is the foundation of machine learning algorithms.

- **Learning from patterns:** Machine learning methods can find complex patterns and correlations in data that may be difficult for humans.

- High-quality, relevant data leads to accurate models.

- **Making predictions:** Based on learned patterns, models predict outcomes **unseen data** points (e.g., weather, future sales).

- **Driving decisions:**Data guides decisions, making industries more efficient and effective.

# Types of data

- **Structured Data:** highly organized and follows a specific format.
  - Common in business applications, financial records, health data, and relational databases.
- **Unstructured Data:** lacks a predefined data model or structure.
  - Text documents (e.g., emails, articles, social media posts).
  - Multimedia files (e.g., images, audio, video).
- **Semi-Structured Data:** it has some organizational properties but doesn't conform to a strict structure like structured data.
  - XML, API, and JSON

# Data Collection and Preparation

Data collection and preparation are crucial steps in the data analysis and machine learning pipeline.

- **Collecting data from various sources:** sensors, databases, APIs, web scraping, surveys, logs.. Each source may have its own format and structure.
- **Data preprocessing**: cleaning, normalization, handling missing values, and feature engineering (generate new features)
- **Why it is a challenge?**
  - availability
  - quality
  - relevance
  - privacy

# Data Representation

- Tabular data: Rows and columns

- Text data: Bag-of-words, TF-IDF.

- Image data: Pixel values.

- Sound data: Waveforms, spectrograms.

# Data Representation

Tabular data: Rows and columns

features — labels

| Family History | Resting Blood Pressure | Cholesterol | Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |

data points

Text data: Bag-of-words, TF-IDF.



I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!

| | |
|---|---|
| it | 6 |
| I | 5 |
| the | 4 |
| to | 3 |
| and | 3 |
| seen | 2 |
| yet | 1 |
| would | 1 |
| whimsical | 1 |
| times | 1 |
| sweet | 1 |
| satirical | 1 |
| adventure | 1 |
| genre | 1 |
| fairy | 1 |
| humor | 1 |
| have | 1 |
| great | 1 |
| … | … |

Image data: Pixel values.

# Data Representation

Sound data: Waveforms, spectrograms.

# Using Data in Machine Learning

- Splitting data into training, validation, and testing sets.
- Cross-validation for model evaluation.
- Handling imbalanced datasets.

# Training, Validation and Testing Data

- **Training data**: Used to train machine learning models.

- **Validation data**: Used to fine-tune model parameters.

- **Testing data**: Used to evaluate model performance.

# Training, Validation and Testing Data

https://teachablemachine.withgoogle.com/

# Cross-validation

Cross-validation is a method to check how well a model works by repeatedly testing it on different parts of the data.

**K-Fold Cross-Validation:** The dataset is divided into k subsets/folds of approximately equal size.
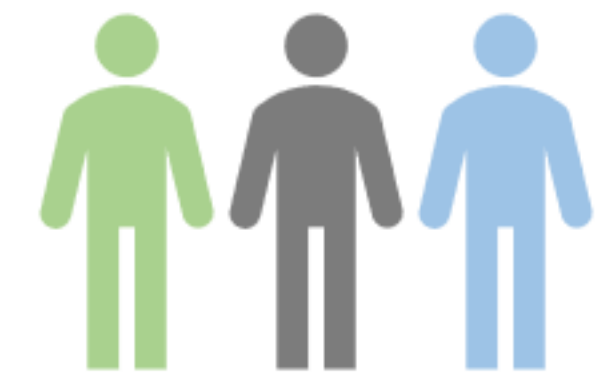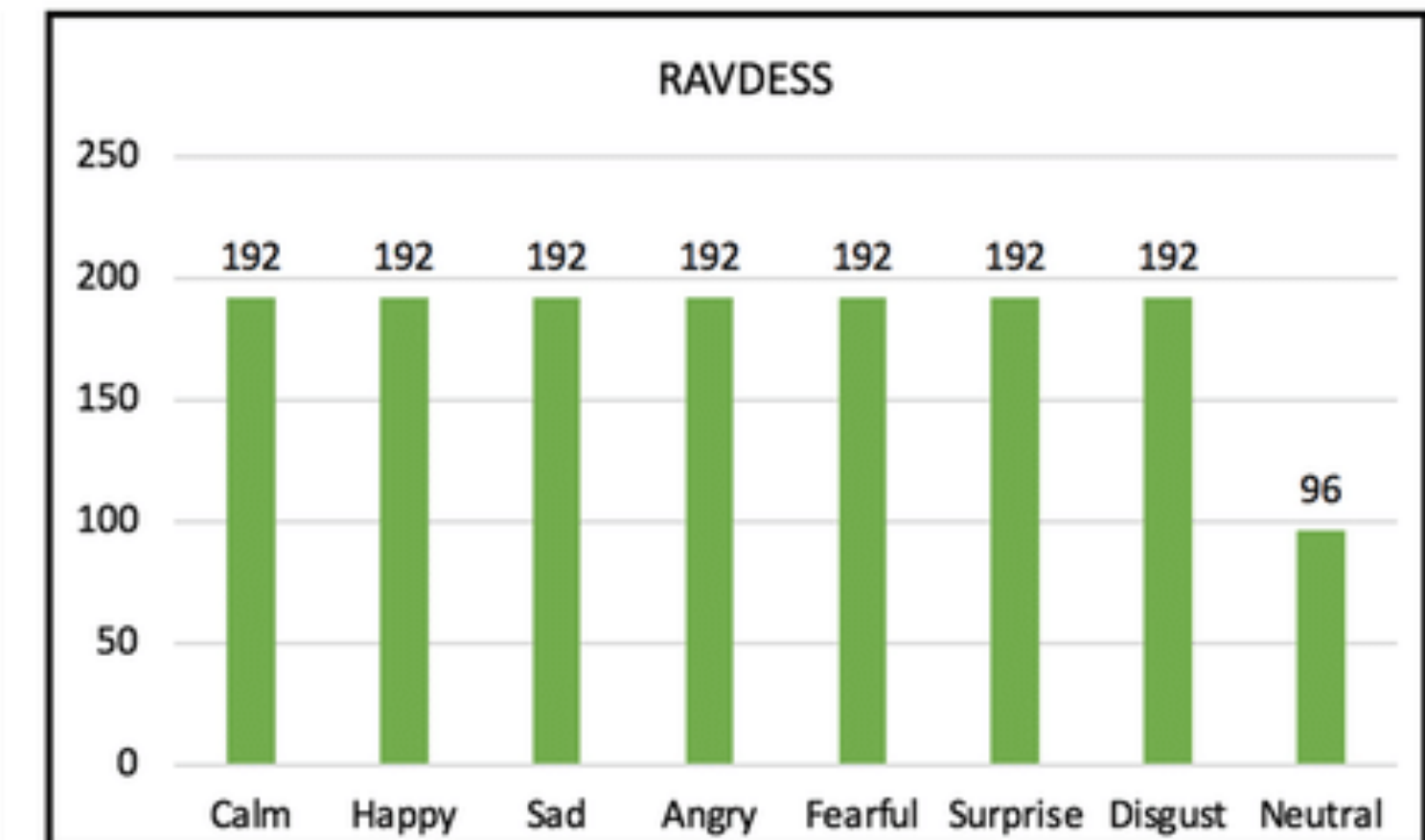


4-fold validation (k=4)
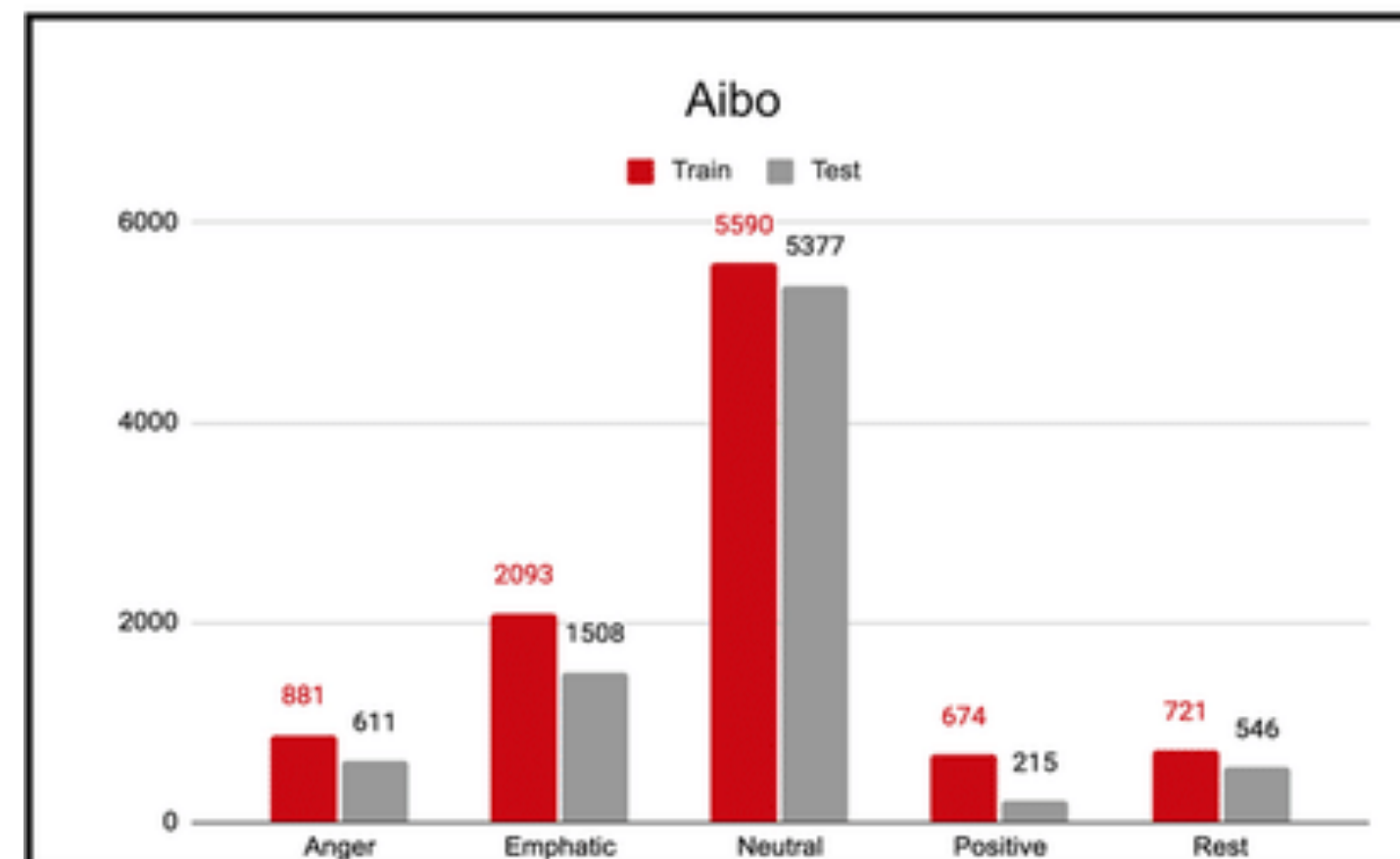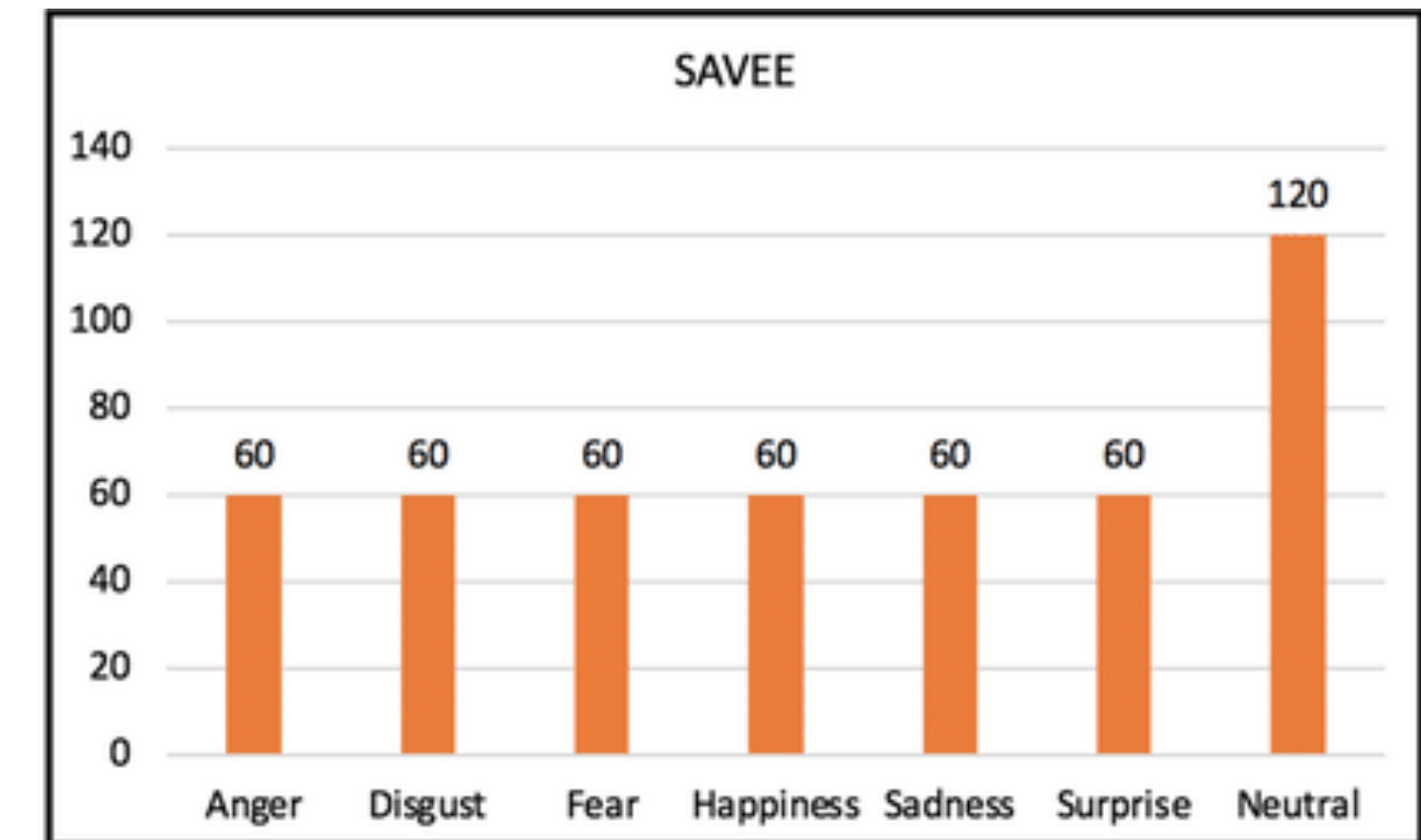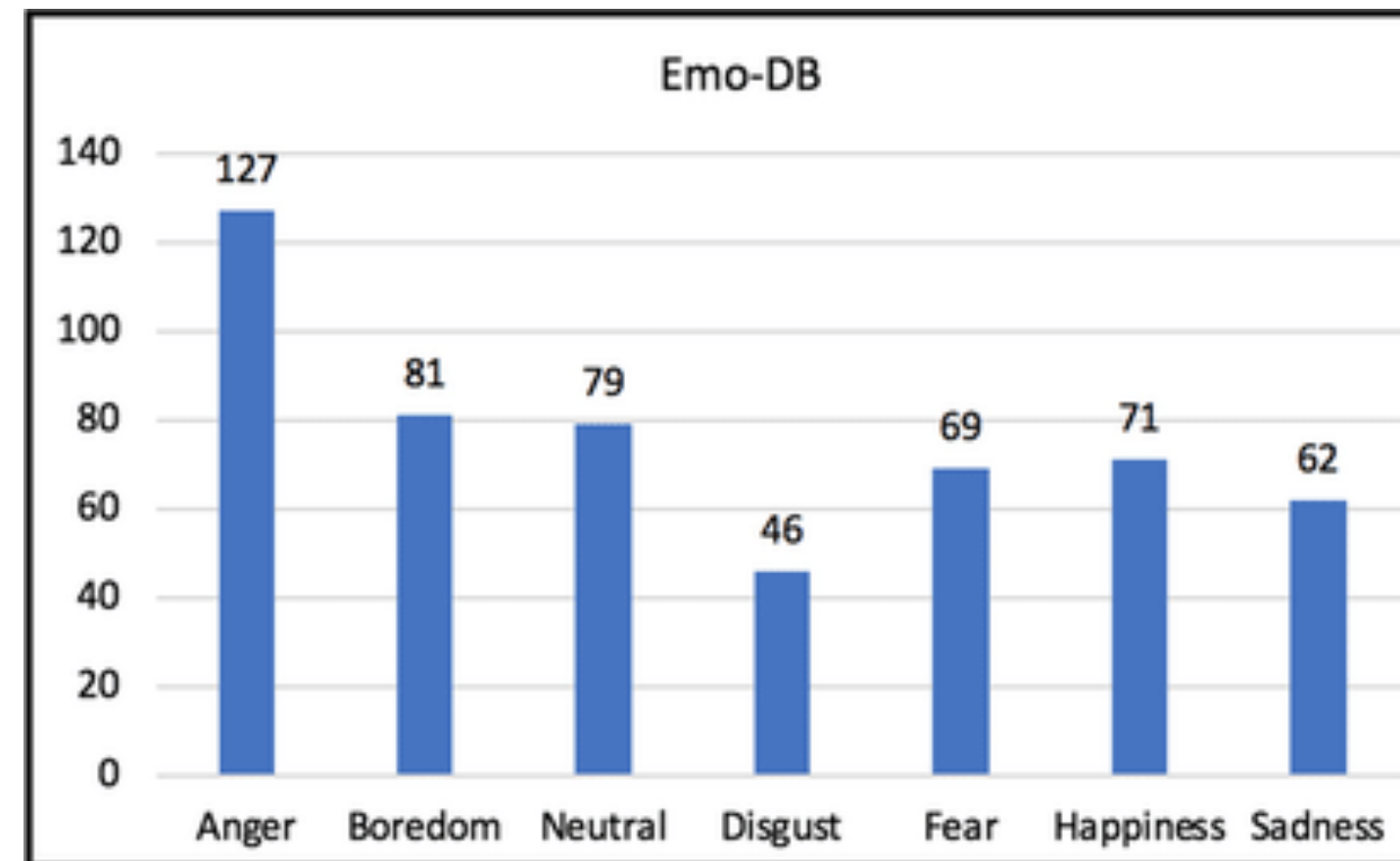
# Cross-validation

**Leave-One-Out:** The model is trained on all but one data point and tested on the one left out. This process is repeated for each data point in the dataset.

# Imbalanced Datasets

Imbalanced datasets refer to datasets where one class has significantly fewer samples compared to the other(s)
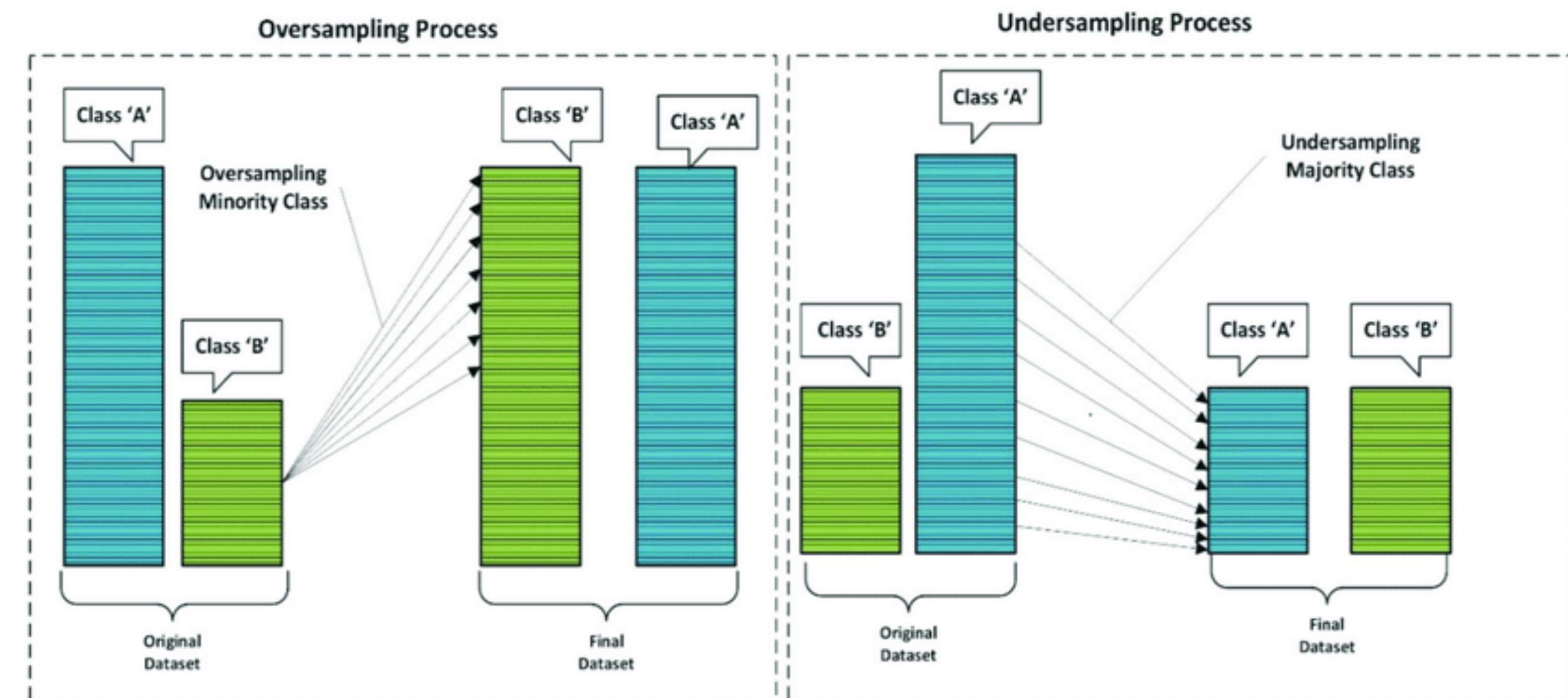
# Handling Imbalanced Datasets

Resampling Methods:

- **Undersampling**: Randomly removing samples from the majority class to balance the class distribution.

- **Oversampling**: Duplicating samples from the minority class or generating synthetic samples to increase its representation.
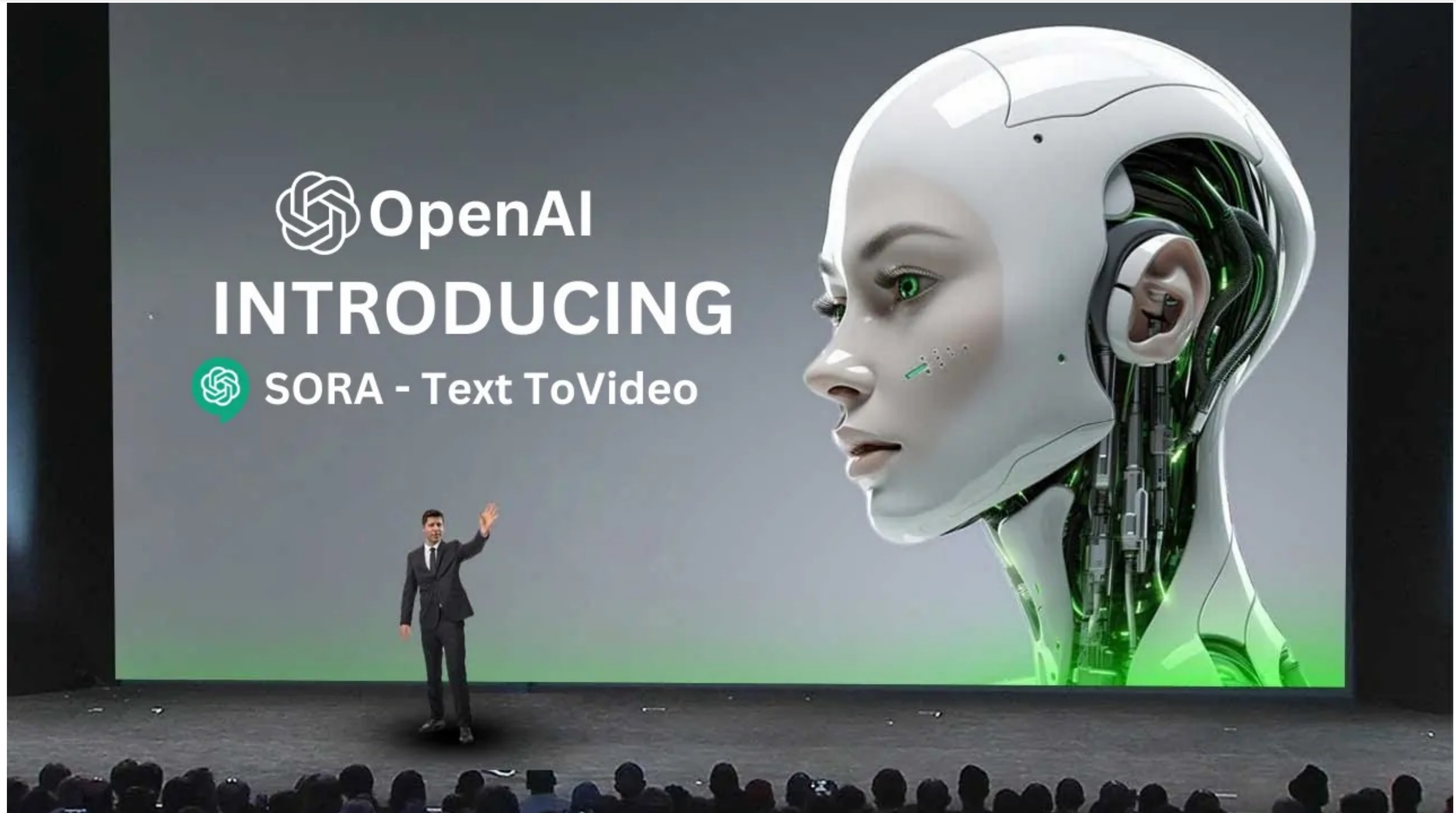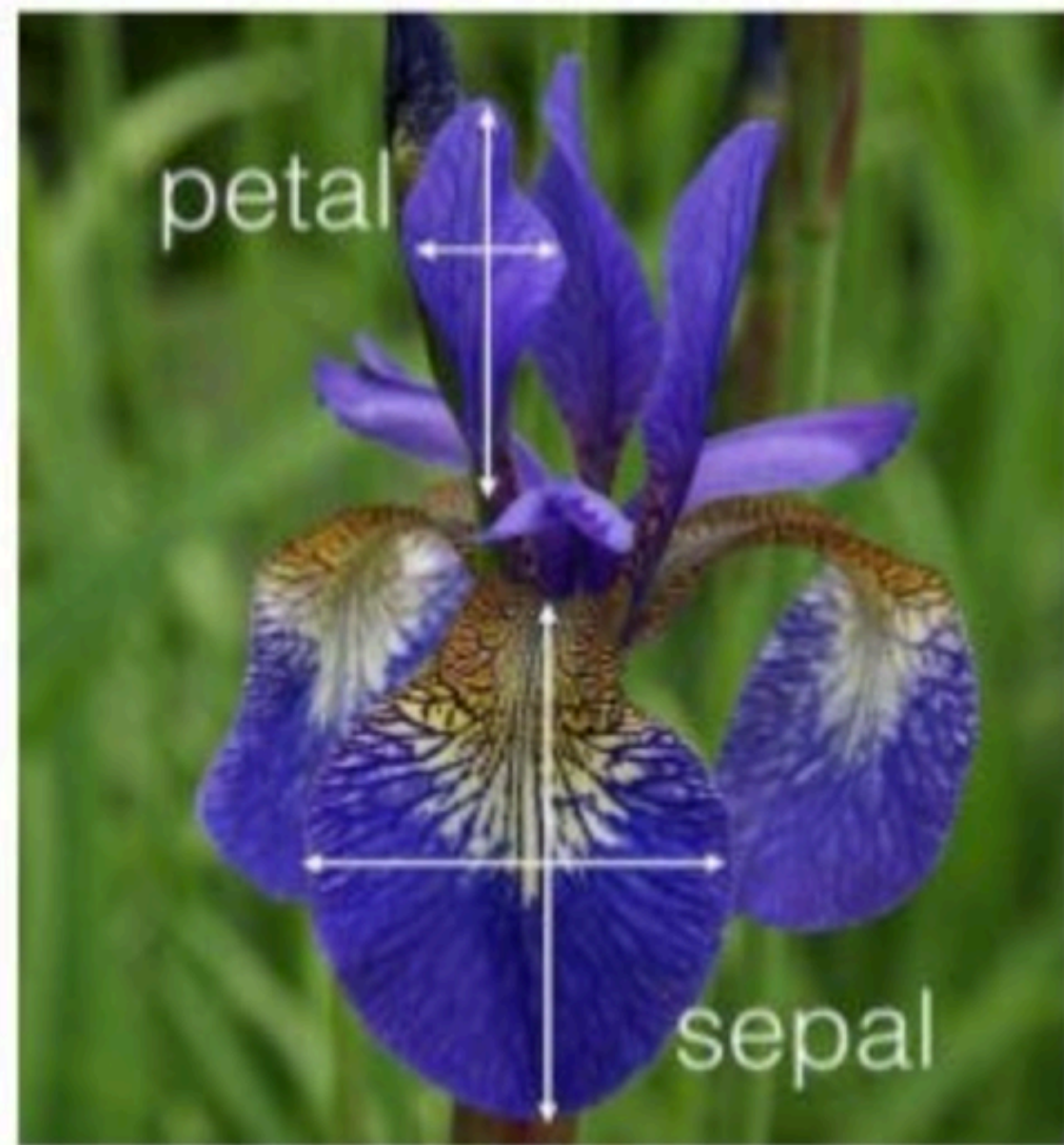
# AI State of the art Applications

# Features

In the context of machine learning, features refer to individual measurable properties or characteristics of the data that are used as inputs for a predictive model.



### Training / test data

| Features | | | | Labels |
|---|---|---|---|---|
| Sepal length | Sepal width | Petal length | Petal width | Species |
| 5.1 | 3.5 | 1.4 | 0.2 | Iris setosa |
| 4.9 | 3.0 | 1.4 | 0.2 | Iris setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | Iris versicolor |
| 6.4 | 3.2 | 4.5 | 1.5 | Iris versicolor |
| 6.3 | 3.3 | 6.0 | 2.5 | Iris virginica |
| 5.8 | 3.3 | 6.0 | 2.5 | Iris virginica |

features — labels

| Family History | Resting Blood Pressure | Cholesterol | Heart Disease? |
|---|---|---|---|
| Yes | Low | Normal | No |
| No | Medium | Normal | No |
| No | Low | Abnormal | Yes |
| Yes | Medium | Normal | Yes |
| Yes | High | Abnormal | Yes |

data points

# Feature Extraction

Feature extraction is a fundamental aspect of machine learning, particularly in tasks involving structured or unstructured data.

It is the process of selecting, transforming, and combining raw data into a set of features that capture relevant information for a given task or problem



Raw Data     Feature Learning     Algorithm     Assignment

(a)

Raw Data     Feature Learning & Algorithm     Assignment

(b)

# Practical - Extract features from image

- Edge Detection: Detects edges in the image using the Canny edge detector algorithm. (Calculates the number of detected keypoints)

- Keypoint Detection: Identifies interest points (corners) in the image.

- HOG Features: Calculates the Histogram of Oriented Gradients (HOG) features. (Calculates the mean and standard deviation of the HOG features)