# Supervised Learning: Classification, Regression and Decision Trees



Lida Kanari

Mathematical Institute University of Oxford

Introduction to Machine Learning, October 2025









## Table of Contents



In this lecture, we will cover the following topics:

- ► Key concepts of Supervised Learning
- Classification
- ► Bayes Classifier
- Decision Trees
- Regression

## Supervised Learning



Supervised learning consists of algorithms that are trained on labeled data. It is suitable for solving problems where the goal is to predict the output value based on the input data.

- ▶ There is a clear relationship between the input data and the output labels.
- ► Sufficient labeled data is available to train the algorithm.
- ► The problem can be defined as a classification or regression task.
- ► The goal is to make accurate predictions on new, unseen data.



**Definition.** Let  $\mathcal{X}$  be the input space and  $\mathcal{Y}$  the output space. Given a training set

$$S = \{(x_i, y_i)\}_{i=1}^n \subset \mathcal{X} \times \mathcal{Y},$$

the goal of supervised learning is to find a function

$$f: \mathcal{X} \to \mathcal{Y}$$

such that f generalizes well to unseen data:

$$S' = \{(x_i', y_i')\}_{i=1}^m \subset \mathcal{X} \times \mathcal{Y},$$



**Setup.** Given the training data

$$(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}, \quad i = 1, \ldots, n,$$

the goal is to predict a new  $y \in \mathcal{Y}$  from a previously unseen  $x \in \mathcal{X}$ . **Prediction rule.** A supervised learning algorithm outputs a function

$$f: \mathcal{X} \to \mathcal{Y}$$
.

**Challenge:** Good predictions on unseen (test) data.



**Question:** What would be the best possible predictor for  $p(\mathcal{Y} \mid \mathcal{X})$  if the joint distribution  $p(\mathcal{X}, \mathcal{Y})$  is known?

- ▶ Random variables  $(\mathcal{X}, \mathcal{Y}) \sim p$  on  $\mathcal{X} \times \mathcal{Y}$ .
- ▶ Prediction function  $f: \mathcal{X} \to \mathcal{Y}$ .
- ▶ Performance measured by a *loss function*  $\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}$ .



**Definition.** A loss function measures the penalty for predicting  $z \in \mathcal{Y}$  when the true label is  $y \in \mathcal{Y}$ :

$$\ell: \mathcal{Y} \times \mathcal{Y} \to \mathbb{R}_{\geq 0}, \qquad (y, z) \mapsto \ell(y, z).$$

## Interpretation:

- ▶  $\ell(y, z) = 0$  when prediction z is correct.
- ▶ Larger  $\ell(y, z)$  means worse prediction.

**Role in learning:** The choice of  $\ell$  determines what we consider a "good" predictor and drives both theoretical analysis and practical algorithms.



Example	Task	<b>Loss function</b> $\ell(y, z)$
Binary classification	$y \in \{-1, 1\}$	$1_{y\neq z}$ (0–1 loss)
Multiclass classification	$y \in \{1, \ldots, K\}$	$1_{y\neq z}$ (0–1 loss)
Regression	$y \in \mathbb{R}$	$(y-z)^2$ (squared loss)
Regression	$y \in \mathbb{R}$	y-z  (absolute loss)

Oxford Supervised Learning October 2025 8/67
Mathematics



## **Examples of supervised learning problems:**

- ▶ Binary classification: Is this image a Cat y = 1 or a Dog y = -1?
- ▶ Multiclass classification: Given a set of medical exam measurements classify the patient into the following categories: Healthy y=1, Disease 1 y=2, Disease 2 y=3,...
- ▶ **Regression:** Predicting house prices from a list of features (house size, location, dates).

Oxford Supervised Learning October 2025 9 / 67
Mathematics



## **Expected risk (generalization error):**

$$R(f) = \mathbb{E}[\ell(\mathcal{Y}, f(\mathcal{X}))] = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, f(x)) dp(x, y).$$

## **Empirical risk (training error):**

$$\hat{R}(f) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f(x_i)).$$

## Classification of a Medical Condition



We want to predict the condition of a patient, based on some medical observations. If the patient is healthy (H), the doctor does not need to prescribe anything. If the patient has a virus (V), the doctor will prescribe painkillers. If the patient has a bacterial infection (B) the doctor needs to prescribe antibiotics. The doctor will assess the condition based on the following features:

- ▶ Temperature
- Blood pressure
- Coughing frequency



We want to predict the condition of a patient, based on some medical observations.

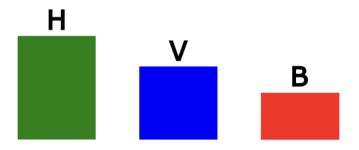


Figure: Patients distributed in three medical conditions: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 12 / 67
Mathematics





We want to predict the condition of a patient, based on some medical observations:

	PatientID	Temperature	BloodPressure	CoughFreq	Condition
0	81	39.73	118.5	14.0	В
1	85	38.61	133.7	13.6	В
2	34	35.75	115.9	2.9	Н
3	82	40.70	139.3	13.8	В
4	94	38.80	147.1	15.0	В
5	18	36.85	128.0	0.6	Н
6	37	36.77	127.3	0.4	Н
7	83	39.83	130.1	16.7	В
8	70	38.70	114.6	6.4	V
9	66	38.21	114.1	3.9	V

Figure: Medical conditions: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 13/67



#### Distribution of Temperature

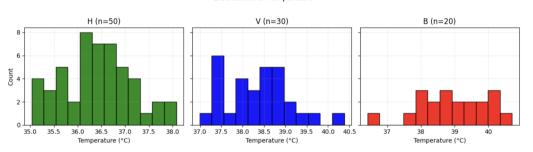


Figure: Temperature by condition: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 14/67



#### Density estimates for Temperature

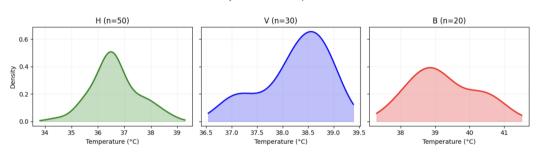


Figure: Temperature density by condition: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 15 / 67
Mathematics

## Probabilistic view on the classification task



We will make some assumptions so that we can use probabilities on our problem.

- Independence between the features.
- ► Each feature within each condition is statistically distributed according to a specific known distribution.



Let's assume that the features can be characterized by a normal distribution. For example, the temperature of each patient is normally distributed according to the following:

$$p(T \mid c, \mathbf{w}) \approx N(\mu_c, \sigma_c^2)$$

where  $\boldsymbol{w}$  is a set of parameters:

$$\mathbf{w} = (\mu_H, \mu_V, \mu_B, \sigma_H^2, \sigma_V^2, \sigma_B^2)$$



Let's assume that the features can be characterized by a normal distribution for each condition:

$$N(\mu_H, \sigma_H^2)$$

$$N(\mu_V, \sigma_V^2)$$

$$N(\mu_B, \sigma_B^2)$$



Let's assume that the features can be characterized by a normal distribution for each condition:

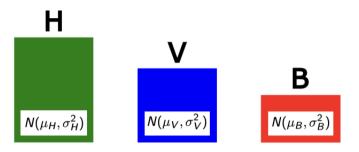


Figure: Patients distributed in three medical conditions: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 19 / 67



The probability of an event represents the frequency over a large number of samples. For example, if I keep taking measurements of the temperature T of patients, the proportion of temperatures of healthy patients in the interval  $[37,38]^{\circ}C$  degrees will converge to:

$$P(T = [37, 38] \mid H, \mathbf{w}) = \int_{37}^{38} p(T \mid H, \mathbf{w}) dT$$



The probability of an event, represents the degree of belief about the event. For example, if I take the temperature of a patient, I cannot predict the exact value, but I can estimate the probability that it will be between [37, 38]°C degrees as:

$$P(T = [37, 38] \mid H, \mathbf{w})$$

Furthermore, if  $P(T = [37, 38] \mid H, \mathbf{w}) > P(T = [38, 39] \mid H, \mathbf{w})$  it is more likely that the measurement will be between  $[37, 38]^{\circ}C$  degrees.

# Bayes classification task



How do we infer the posterior  $p(\mathbf{w} \mid \mathcal{X})$  from the input data  $\mathcal{X}$ ?



23 / 67

How do we infer the posterior  $p(\mathbf{w} \mid \mathcal{X})$  from the input data  $\mathcal{X}$ ?

Let's recall from Bayes probability:

$$p(\mathcal{X}, \mathbf{w}) = p(\mathbf{w} \mid \mathcal{X})p(\mathcal{X})$$



How do we infer the posterior  $p(\mathbf{w} \mid \mathcal{X})$ ?

Let's recall from Bayes probability:

$$p(\mathcal{X}, \mathbf{w}) = p(\mathbf{w} \mid \mathcal{X})p(\mathcal{X}) = p(\mathcal{X} \mid \mathbf{w})p(\mathbf{w})$$



How do we infer the posterior  $p(\mathbf{w} \mid \mathcal{X})$ ?

Let's recall from Bayes probability:

$$p(\mathcal{X}, \mathbf{w}) = p(\mathbf{w} \mid \mathcal{X})p(\mathcal{X}) = p(\mathcal{X} \mid \mathbf{w})p(\mathbf{w})$$

Bayes theorem:

$$p(\mathbf{w} \mid \mathcal{X}) = \frac{p(\mathcal{X} \mid \mathbf{w})p(\mathbf{w})}{p(\mathcal{X})}$$



Bayes theorem:

$$p(\mathbf{w} \mid \mathcal{X}) = \frac{p(\mathcal{X} \mid \mathbf{w})p(\mathbf{w})}{p(\mathcal{X})}$$

where,

 $p(\mathbf{w})$  is the prior representing the initial belief about  $\mathbf{w}$ .

 $p(\mathbf{w} \mid \mathcal{X})$  is the posterior representing the updated belief about  $\mathbf{w}$  after observations of data D.

 $p(\mathcal{X} \mid \mathbf{w})$  is the likelihood that the data  $\mathcal{X}$  are observed for the distributions with parameters  $\mathbf{w}$ .



## Two types of uncertainty:

1. Epistemic: which model is correct, given the data we have observed? The more patient data I collect, the more certain I am about  $\mu_H$ ,  $\sigma_H^2$ .



## Two types of uncertainty:

- 1. Epistemic: which model is correct, given the data we have observed? The more patient data I collect, the more certain I am about  $\mu_H$ ,  $\sigma_H^2$ .
- 2. Aleatoric: related to the stochastic nature of the variables. I will never be able to predict a precise temperature measurement even when I know exactly the condition of the patient.

# Bayes classification task



### Based on our assumptions:

- ► Independence of features
- ► Normal distribution of features within a class



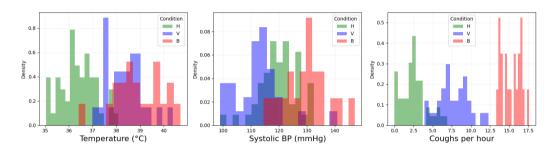
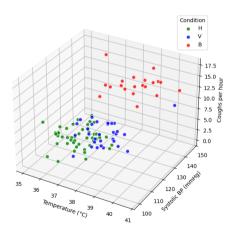


Figure: Temperature, blood pressure and coughing frequency for: healthy (H), viral (V) or bacterial (B) infection.

Oxford Supervised Learning October 2025 30 / 67
Mathematics







**Definition.** Suppose  $\mathcal{Y} = \{1, \dots, K\}$  is a finite label set. The *Bayes classifier* assigns to each input  $x \in \mathcal{X}$  the class with the highest posterior probability:

$$f^*(x) = \arg \max_{y \in \mathcal{Y}} \ \mathbb{P}(\mathcal{Y} = y \mid \mathcal{X} = x).$$

**Intuition:** The Bayes classifier uses complete knowledge of the data distribution to minimize the classification error.

**Key property:** Assuming perfect knowledge of p(x, y) then  $f^*$  is the optimal classifier, achieving the lowest possible error (the *Bayes error*).

Oxford Authermatics Supervised Learning October 2025 32 / 67



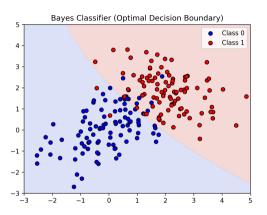


Figure: Example of a Bayes classifier in a binary classification problem



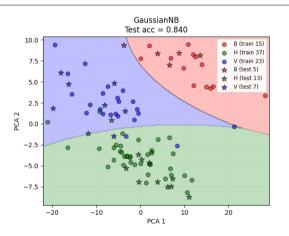


Figure: Bayes classifier to distinguish medical conditions.

Oxford Supervised Learning October 2025 34/67 Mathematics



**Definition.** The lowest achievable risk given full knowledge of *p* is the *Bayes risk*:

$$R^* = \mathbb{E}_X \left[ \inf_{z \in \mathcal{Y}} \mathbb{E}[\ell(y, z) \mid \mathcal{X}] \right].$$

Excess risk is defined as:

$$R(f)-R^*\geq 0.$$



**Definition.** The Bayes optimal predictor minimizes conditional risk:

$$R^*(x) = \arg\min_{z \in \mathcal{Y}} \ \mathbb{E}[\ell(Y, z) \mid X = x].$$

## **Special cases:**

▶ Classification: 0-1 loss

► **Regression:** squared loss



Let 
$$\eta(x) = p(\mathcal{Y} = 1 \mid \mathcal{X} = x)$$
.

### Bayes classifier:

$$f^*(x) = \begin{cases} 1 & \text{if } \eta(x) > 1/2, \\ -1 & \text{if } \eta(x) < 1/2. \end{cases}$$

### Bayes risk:

$$R^* = \mathbb{E}[\min\{\eta(x), 1 - \eta(x)\}].$$



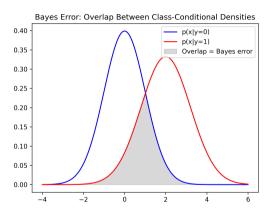


Figure: Example of a Bayes error



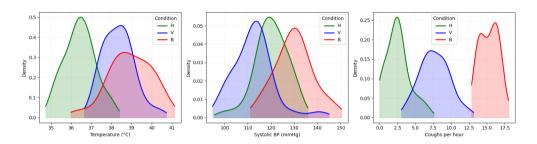


Figure: Bayes error in the classification of medical conditions.

Oxford Supervised Learning October 2025 39 / 67
Mathematics



**Key idea:** The Bayes predictor is optimal but requires full knowledge of p(x, y). **In practice:** The distribution p(x, y) is unknown.

**Objective:** Approximate the Bayes predictor from finite training data:

- ▶ **Generative methods:** approximate p(x | y) and p(y), then apply Bayes rule.
- ▶ **Discriminative methods:** approximate  $f^*$  which corresponds to  $p(y \mid x)$  directly from the data.

Oxford Supervised Learning October 2025 40 / 67



#### Discriminative models:

- ▶ Model the conditional distribution  $p(y \mid x)$  directly.
- Learn a decision boundary between classes.
- Examples: Logistic regression, SVMs, neural networks.



#### Generative models:

- ▶ Model the joint distribution  $p(x, y) = p(y) p(x \mid y)$ .
- ► Capture how the data x is generated for each class y.
- ▶ Use Bayes' rule to compute  $p(y \mid x)$ :

$$p(y \mid x) = \frac{p(y) p(x \mid y)}{\sum_{y'} p(y') p(x \mid y')}.$$

► Examples: Naive Bayes, LDA, QDA.



We model:

$$p(x, y \mid \theta, \pi) = p(y \mid \pi) p(x \mid y, \theta)$$

where

$$p(y = c \mid \pi) = \pi_c, \qquad \sum_{c} \pi_c = 1.$$

#### Procedure:

- 1. Choose a model for  $p(x \mid y = c, \theta_c)$  (e.g. Gaussian, multinomial).
- 2. Estimate parameters  $\pi_c, \theta_c$  from training data.
- 3. Predict using:

$$\hat{y} = \arg \max_{c} \pi_{c} p(x \mid y = c, \theta_{c}).$$



**Assumption:** Features are conditionally independent given the class.

$$p(x \mid y = c, \theta_c) = \prod_{j=1}^{D} p(x_j \mid y = c, \theta_{jc})$$

### **Advantages:**

- Simple and fast to train.
- Surprisingly effective for text and categorical data.

### **Limitations:**

- Independence assumption rarely true.
- Model is often misspecified but works well for classification.

# Linear Discriminant Analysis (LDA)



**Goal:** classify by modelling how each class *generates* the data, then use Bayes' rule. Key idea:

- ► For each class c assume the features  $x \in \mathbb{R}^D$  are Gaussian with mean  $\mu_c$  and a shared covariance  $\Sigma$ .
- ▶ Because covariances are shared, class-conditional contours are ellipses with the same shape, only their centers differ.
- Decision boundaries between classes become linear in x.



#### Assume for each class c:

$$p(x \mid y = c) = \mathcal{N}(x \mid \mu_c, \Sigma).$$

Prior class probabilities:  $\pi_c = p(y = c)$ ,  $\sum_c \pi_c = 1$ .

Posterior (unnormalised):

$$p(y = c \mid x) \propto \pi_c \exp\left(-\frac{1}{2}(x - \mu_c)^{\top} \Sigma^{-1}(x - \mu_c)\right).$$

Algebra (rearrange quadratic term):

$$p(y = c \mid x) \propto \exp(\beta_c^{\top} x + \gamma_c),$$

with

$$\beta_c = \Sigma^{-1} \mu_c, \qquad \gamma_c = -\frac{1}{2} \mu_c^{\mathsf{T}} \Sigma^{-1} \mu_c + \log \pi_c.$$



- $\blacktriangleright$   $\mu_c$ : empirical mean of class c (centre of class c).
- $\triangleright$   $\Sigma$ : shared covariance across classes: controls feature scaling and correlation.
- $\beta_c = \Sigma^{-1}\mu_c$ : a weight vector: direction in which class c pulls the log-odds.
- $ightharpoonup \gamma_c$ : offset combining class prior  $\pi_c$  and a quadratic correction  $-\frac{1}{2}\mu_c^{\top}\Sigma^{-1}\mu_c$ .
- ▶ Decision function for class c:  $g_c(x) = \beta_c^\top x + \gamma_c$ . Predict class with largest  $g_c(x)$ .



### Training (ML estimates):

- 1. Compute class counts  $N_c$  and priors  $\hat{\pi}_c = N_c/N$ .
- 2. Compute class means:  $\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} x_i$ .
- 3. Compute pooled covariance:

$$\hat{\Sigma} = \frac{1}{N-C} \sum_{c} \sum_{i:y_i=c} (x_i - \hat{\mu}_c) (x_i - \hat{\mu}_c)^{\top}$$

4. Set  $\hat{\beta}_c = \hat{\Sigma}^{-1} \hat{\mu}_c$ ,  $\hat{\gamma}_c = -\frac{1}{2} \hat{\mu}_c^{\top} \hat{\Sigma}^{-1} \hat{\mu}_c + \log \hat{\pi}_c$ .



Goal: Similar to LDA but allows more flexibility in the shapes.

- ▶ Each class c has its own covariance  $\Sigma_c$  and mean  $\mu_c$ .
- Class-conditional contours can have different shapes and orientations.
- Decision boundaries between classes are in general quadratic surfaces (ellipses, parabolas, etc).

Oxford Supervised Learning October 2025 49 / 67



Assume for each class c:

$$p(x \mid y = c) = \mathcal{N}(x \mid \mu_c, \Sigma_c).$$

Using Bayes' rule:

$$p(y=c\mid x) \propto \pi_c |\Sigma_c|^{-1/2} \exp\left(-\frac{1}{2}(x-\mu_c)^{\top} \Sigma_c^{-1}(x-\mu_c)\right).$$

Decision rule (take log and compare):

$$\hat{y} = \arg \max_{c} \left[ -\frac{1}{2} (x - \mu_c)^{\top} \Sigma_c^{-1} (x - \mu_c) - \frac{1}{2} \log |\Sigma_c| + \log \pi_c \right].$$

Because of the quadratic term  $x^{\top}\Sigma_c^{-1}x$  that depends on c, the boundary is quadratic.

Oxford Supervised Learning October 2025 50 / 67 Mathematics



- $\blacktriangleright \mu_c$ : mean (center) of class c.
- $\triangleright$   $\Sigma_c$ : class-specific covariance (controls shape/orientation of class c's density).
- $ightharpoonup -\frac{1}{2}(x-\mu_c)^{\top}\Sigma_c^{-1}(x-\mu_c)$ : log-likelihood for class c.
- $ightharpoonup -\frac{1}{2}\log|\Sigma_c|$ : volume term, penalises classes with large covariance.
- ▶  $\log \pi_c$ : prior log-probability (accounts for class imbalance).
- ▶ The quadratic term in x depends on  $\Sigma_c^{-1}$ , hence different  $\Sigma_c$  leads to quadratic decision boundaries.



### **Training (ML estimates):**

- 1. Compute class counts  $N_c$  and priors  $\hat{\pi}_c = N_c/N$ .
- 2. Compute class means:  $\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} x_i$ .
- 3. Compute class covariances:

$$\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:v:=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^{\top}.$$

**Prediction:** for a test x compute scores

$$s_c(x) = -\frac{1}{2}(x - \hat{\mu}_c)^{\top} \hat{\Sigma}_c^{-1}(x - \hat{\mu}_c) - \frac{1}{2} \log |\hat{\Sigma}_c| + \log \hat{\pi}_c,$$

and predict  $\arg \max_c s_c(x)$ .



### **Summary table:**

Model	Covariance	Decision boundary
LDA	$\Sigma$ (shared)	Linear
QDA	$\Sigma_c$ (class-specific)	Quadratic
Naïve Bayes	Diagonal / independent	Linear/naive

Oxford Mathematics Supervised Learning October 2025 53/67



#### Trade-offs:

- ▶ QDA more flexible (can model complex class shapes) but needs more data to estimate  $\Sigma_c$ .
- ightharpoonup LDA is a compromise: fewer parameters ightharpoonup lower variance but possibly higher bias if covariances truly differ.
- ► Naïve Bayes imposes stronger independence assumptions and is very parameter-efficient.

Practical rule: use QDA if you have lots of data and evidence that the covariances differ.



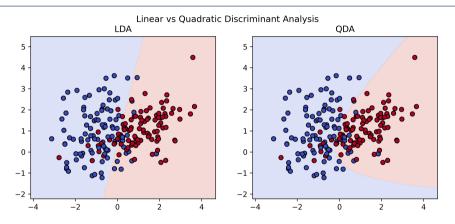
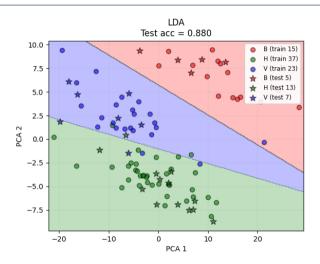
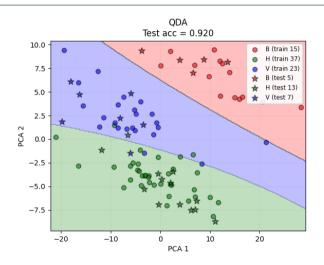


Figure: Example of LDA vs QDA











**Definition.** In regression, the output space is continuous,  $\mathcal{Y} \subseteq \mathbb{R}$ .

The goal is to find a function  $f: \mathcal{X} \to \mathbb{R}$  that predicts  $\mathcal{Y}$  from  $\mathcal{X}$  by minimizing the expected squared error:

$$f^*(x) = \arg\min_{f} \mathbb{E}[(\mathcal{Y} - f(\mathcal{X}))^2 \mid \mathcal{X} = x].$$

**Result:** The optimal predictor is the conditional expectation:

$$f^*(x) = \mathbb{E}[\mathcal{Y} \mid \mathcal{X} = x].$$

**Example:** Linear regression assumes  $f(x) = w^{T}x + b$ .



Given  $x_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$ , we model

$$f(x) = w^{\top}x + b.$$

The least-squares estimator minimizes

$$\hat{R}_n(w,b) = \frac{1}{n} \sum_{i=1}^n (y_i - w^\top x_i - b)^2.$$

Solution:

$$\hat{w} = (X^{\top}X)^{-1}X^{\top}y, \qquad \hat{b} = \bar{y} - \hat{w}^{\top}\bar{x},$$

provided  $X^{\top}X$  is invertible.



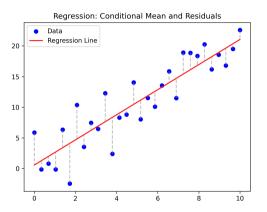


Figure: Example of Linear Regression



**Setup:** Binary classification with  $\mathcal{Y} = \{-1, 1\}$ .

**Model:** Posterior probability

$$p(y = 1 \mid x) = \sigma(w^{\top}x + b), \quad \sigma(t) = \frac{1}{1 + e^{-t}}.$$

Decision rule:

$$f(x) = \operatorname{sign}(w^{\top}x + b).$$

**Learning:** Parameters (w, b) are estimated by minimizing the *logistic loss*, a convex surrogate of the 0-1 loss.

Supervised Learning October 2025 61 / 67



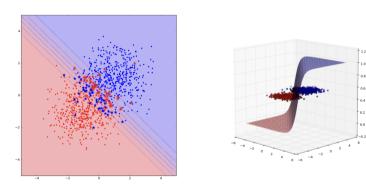


Figure: Example of Logistic Regression (image from Prof. Seth Flaxman)



**Idea:** Partition the input space  ${\mathcal X}$  into regions and assign a prediction in each region.

#### Construction:

- ightharpoonup Recursively split  $\mathcal X$  along feature dimensions.
- At each node, choose the split that maximizes class separation (e.g., Gini index, entropy).

## **Decision Trees**



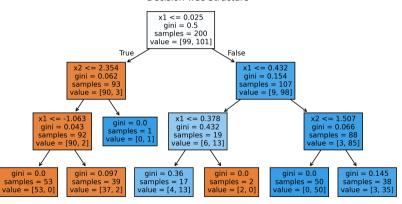
**Prediction:** A new point x is assigned to the region it falls into; prediction is the majority label (classification) or mean response (regression).

#### Remarks:

- Easy to interpret.
- ▶ High variance (often reduced with ensembles: random forests, boosting).



#### **Decision Tree Structure**





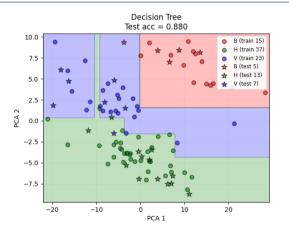


Figure: Decision Trees of medical example

Thank you! Questions?