

# Lecture 8

Decision tree

Neural Network

L1 - edited Log. Log.

LIV - posted

Q1 - Thu 2/12 4:02

don't be late

HW3 due tomm.

Sol 2 - Now ✓

Sol 1 - on thu

HW4 tomorrow  
ungraded

## Predictors

- K-Nearest-Neighbor
- Linear - for regression  
- for classification
- Logistic regression ↗
- Perceptron, LDA
- Decision Trees (CART)

## Algorithms

- LS Regression
- Logistic Regression by  
Gradient Ascent/Descent
- CART greedy algorithm

## Concepts

- Decision Region, Dec. Boundary
- Training error, Test error  
Expected error ↗
- Variance, Bias
- Loss functions - training/  
test/expected loss
- Max likelihood

# Lecture III: Classification and Decision Trees (CART)

Marina Meilă  
mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath  
Cheriton School of Computer Science  
University of Waterloo

January 27, 2026

Classification and regression tree(s) (CART) ✓

Learning a CART ←

Predicting with a CART ✓

Some issues with CART ←

**Reading** HTF Ch.: 9.2 CART, Murphy Ch.: 16.2.1–4 CART, Bach Ch.:

## Classification and regression trees (CART)

- ▶ A **classification tree** or (**decision tree**) is built recursively by splitting the data with hyperplanes parallel to the coordinate axes.
  - ▶ At each split, try to separate + examples from - examples as well as possible.
  - ▶ Eventually, all the regions will be "pure", i.e. will contain examples from one class only.
- ▶ Classification trees can be used in multiway classification as well (there we try to create a pure region on at least one side of the split)
- ▶ A **regression tree** uses the same principle for regression  
here we try to obtain regions where the outputs are nearly the same

## Some advantages of CART

- ▶ Natural and easy to interpret (if small)  $\nexists$   $X_{1,2,\dots,d}$  interpretable
- ▶ Can approximate any function (with enough leaves). FLEXIBLE  $\heartsuit$

1) any decision region (or function)

2) Predict any (finite)  $\mathcal{D}$  with  $\mathcal{O}$  error

POWERFUL

LOW BIAS

Variance  $\uparrow$

Reduce variance  $\Rightarrow$  increase Bias  
 limit-depth of tree  
 - # leaves

## Some advantages of CART

- ▶ Natural and easy to interpret (if small)
- ▶ Can approximate any function (with enough leaves)

Ex 1: predict recidivism:  $x$  incarcerated person -  
release?

$x =$  age  
gender  
crime  
time incarcerated  
...



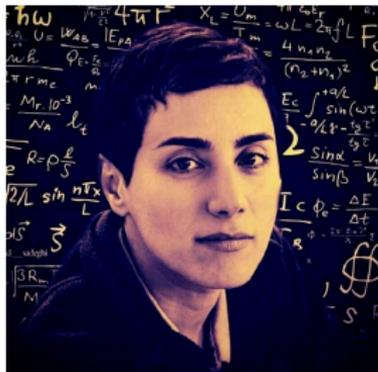
## Some advantages of CART

- ▶ Natural and easy to interpret (if small)
- ▶ Can approximate any function (with enough leaves)

Ex 2 image classification

$X =$  image of face

$Y = \text{😊, 😞, 😐}$



$X \in \mathbb{R}^{500 \times 500}$

→ D Tree → 😐

NOT interpretable

Mayam Mirzahani

# "Learning" a CART

A standard algorithm for building a decision tree works recursively in top-down fashion.

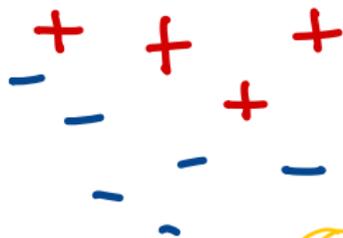
**Input** Training set  $D$  of size  $n$

**Initialize** with a tree with only one region, that contains all the data

**Repeat** until all leaves are pure (or until desired purity is attained in all leaves)

- Find the "optimal" split over all leaves  $R_q$  and all possible splits of  $R_q$ .  
"Optimal" is defined in terms on purity (e.g split the least pure leaf, find the split that makes one of the new leaves pure)
- Perform the "optimal" split and add the two new leaves to the tree

This is a greedy algorithm. Sometimes, trees obtained this way are **pruned** back to smaller sizes.



Init  $R_0, D, n, \hat{y}(R_0) = y_0 = -, p_0 = \frac{6}{10}$   
Impurity:  $I_0$

Recursively:

find split with  $\max I_q - (I_{q^1} + I_{q^2})$   
 $R_q \rightarrow R_{q^1} + R_{q^2}$

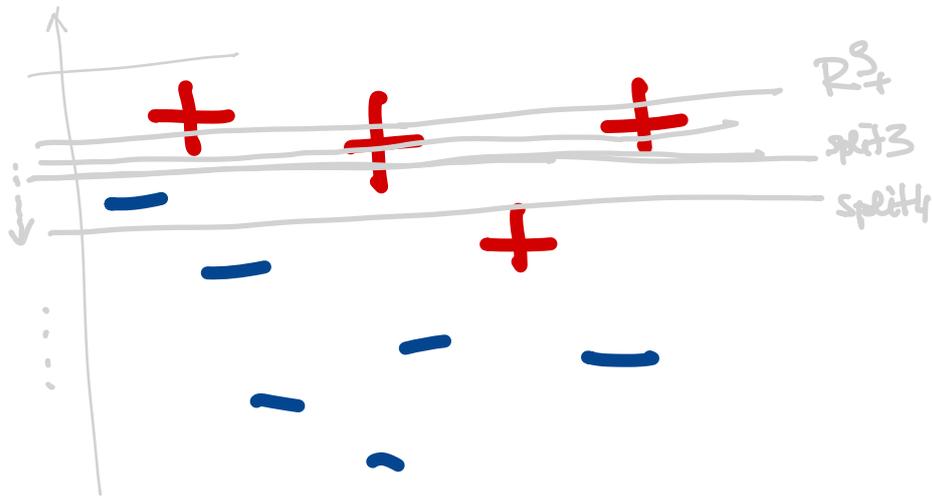
See next page

Until  $I_q \leq \text{threshold}$  for all leaves  $q$

$$n = 10$$

$$n_- = 6 \Rightarrow p_0 = \frac{6}{10}$$

$$n_+ = 4$$



split 3:  $P(R_+^3) = 1$      $P(R_-^3) = \frac{6}{7} \Rightarrow I(R_+^3) + I(R_-^3)$   
 split 4:  $P(R_+^4) = \frac{3}{4}$      $P(R_-^4) = \frac{5}{6} \Rightarrow \dots \leq ?$

Compute time  $\propto d \cdot n \cdot L$   
 $\uparrow$   
 levels

Find best  
Split

$x \in \mathcal{R}^d$

for leaf  $g = 1:n$  leaves

for  $j = 1:d$  dimension of  $x$

for  $i = 1:n-1$

Nodes of leaf  $g$

compute  $\Delta I(g, j, i)$

Split leaf  $g^* = \underset{g, j, i}{\operatorname{argmax}} \Delta I(g, j, i)$

$$\Delta I(g, j, i) = I_g - (I_{g'(j, i)} + I_{g''(j, i)})$$

decrease in impurity

# Purity

- ▶ Natural ways to set  $y_q$  based on the data, once the partition  $\mathcal{T}$  has been fixed:
  - ▶ denote  $Y_q = \{y^i \mid x^i \in R_q, i = 1 : N\}$  the set of labels at a leaf  $R_q$
  - ▶ Regression  $y_q =$  average of  $Y_q$
  - ▶ Classification  $y_q =$  majority label of  $Y_q$
- ▶ a leaf  $R_q$  is **pure** if all labels are the same, i.e. if  $|Y_q| = 1$
- ▶ criteria for the **(im)purity** of a leaf  $R_q$ 
  - ▶ Regression impurity = sample variance of  $Y_q$
  - ▶ Classification let  $p_q$  be the frequency of  $y_q$  in  $Y_q$

$$\text{impurity} = \begin{cases} \text{Misclassification error} & 1 - p_q \\ \text{Gini} & p_q(1 - p_q) \\ \text{Entropy} & - (p_q \ln p_q + (1 - p_q) \ln(1 - p_q)) \end{cases} \quad (2)$$

$\frac{4}{10} = I$   
 $\rightarrow I = \frac{4.6}{10^2}$

These measures generalize naturally to the multiclass setting.

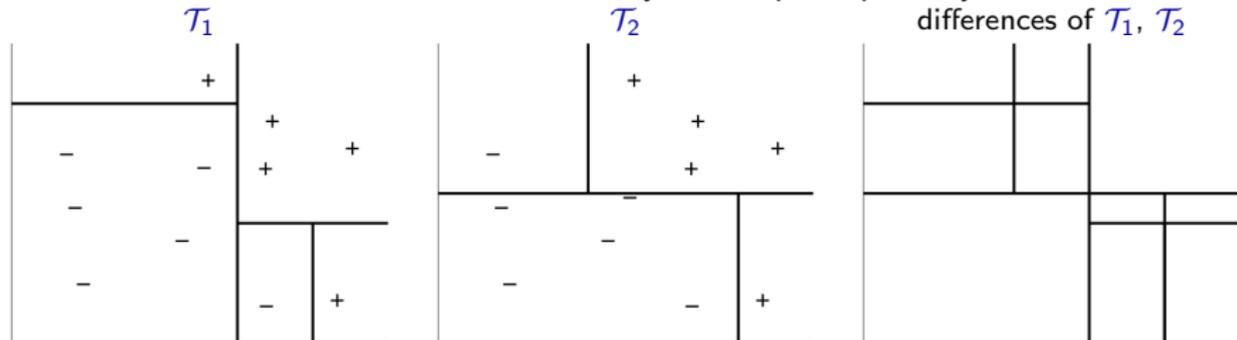
# Predicting with a CART

Given new  $x$

1. Find the unique leaf  $R(x)$  so that  $x \in R(x)$
2. Predict  $\hat{y}$  based on the data in this leaf
  - ▶ **Regression**  
Predict  $\hat{y}(x) = \text{average}\{y^i \text{ with } x^i \in R(x)\}$
  - ▶ **Classification**  
Predict  $\hat{y}(x) = \text{majority}\{y^i \text{ with } x^i \in R(x)\}$

## A decision tree over $\mathcal{D}$ is not unique

Same dataset  $\mathcal{D}$ , two different trees. Both classify the sample  $\mathcal{D}$  perfectly.



But they produce different decision regions.

# Lecture Notes IV – Neural Networks, Part 1

Marina Meilă  
mmp@uwaterloo.ca

With Thanks to Pascal Poupart & Gautam Kamath  
Cheriton School of Computer Science  
University of Waterloo

February 2, 2026

A little history ←

The single “neuron” ←

Two-layer Neural Networks

Hidden layer options

Output layer options

Multi-layer neural networks

**Reading** HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

## A little history

HMM

### First epoch – can computers mimic the brain?

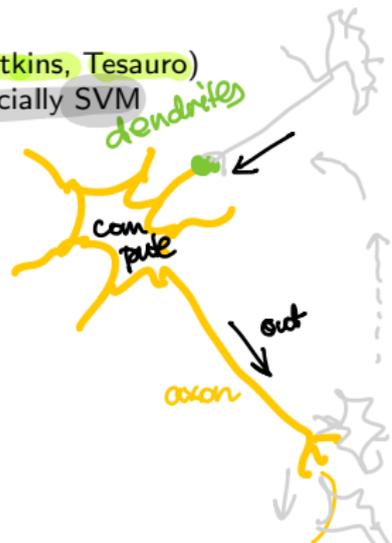
- ▶ 40's-50's the first mathematical models of the neuron (McCullaugh & Pitts, Hebb)
- ▶ '58 the perceptron (Rosenblatt)
- ▶ First winter '69 "Perceptrons" (Minsky and Papert) – they can only classify linearly separable data (and von Neumann computers steal the show)

### Revival in the 80's-'95 – let's use layers!

- ▶ '82 Hopfield associative net (contents adressable memory)
- ▶ '84 Boltzmann Machine (Ackley, Hinton, Sejnowski)
- ▶ autoencoders *Hinton*
- ▶ '86 backpropagation (Sejnowski, Rosenberg)
- ▶ Reinforcement learning (Shannon, Samuel, Barton, Sutton, Watkins, Tesauro)
- ▶ Second winter cca '95 Statistical learning takes the stage, especially SVM

### Current epoch cca 2005 – more data, more layers

- ▶ Deep learning
- ▶ Generative models
- ▶ Attention
- ▶ LLMs



# Brains vs. Computers

## Computer (von Neumann)

- ▶ Electrical binary signals directed by gates
- ▶ Wiring is fixed
- ▶ Sequential and parallel computation
- ▶ Memory is retrieved by address
- ▶ **Fragile** (if a gate stops working, computer crashes). Also, “no friction” – a single instruction can stop the entire machine/program.

## ▶ Brain

- ▶ Electrical signals, units=neurons
- ▶ Wiring is dynamic, changes with brain development, experiences, learning
- ▶ Parallel (and some) sequential computation
- ▶ Memory is distributed
- ▶ **Robust** (when neuron/region dies, brain rewires itself to compensate). **No On/Off master switch**

# Brains vs. Computers

## Computer (von Neumann)

- ▶ Electrical binary signals directed by gates
- ▶ Wiring is fixed
- ▶ Sequential and parallel computation
- ▶ Memory is retrieved by address
- ▶ **Frangible** (if a gate stops working, computer crashes). Also, “no friction” – a single instruction can stop the entire machine/program.

## ▶ Brain

- ▶ Electrical signals, units=neurons
- ▶ Wiring is dynamic, changes with brain development, experiences, learning
- ▶ Parallel (and some) sequential computation
- ▶ Memory is distributed
- ▶ **Robust** (when neuron/region dies, brain rewires itself to compensate). **No On/Off master switch**

## ▶ Neural network

- ▶ Signals are numbers passed between units
- ▶ Network structure is fixed (and dense) but the learned weights allow “rewiring” during training
- ▶ Parallel (in layer) and sequential (feed forward/backward) computation
- ▶ Memory is distributed (in the weights)
- ▶ **Redundant/robust** (no single neuron can influence the output much)

(Artificial) Neural Network (nn) **unit**

$$x \in \mathbb{R}^d$$

$$w_i \in \mathbb{R}^{d+1}$$

- ▶ For each **unit**  $i$

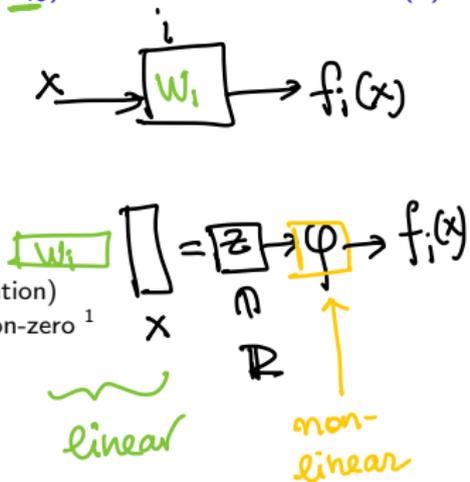
$$y_i \equiv f_i(x) = \phi\left(\sum_j w_{ij}x_j + w_{i0}\right) \quad (1)$$

- ▶ **Weight vector**  $w_i$  **weights**

- ▶  $w_{ij}$  = strength of the link from unit  $i$  to input  $j$
- ▶  $w_{ij} = 0$ : no link
- ▶  $w_{ij}$  can be positive or negative
- ▶ Sometimes we call the input vector  $x = [x_{1:d}]$  **input units**

- ▶ **activation function**  $\phi()$

- ▶ must be non-linear (otherwise the unit is a linear transformation)
- ▶ wanted: monotonically increasing, differentiable, gradient non-zero<sup>1</sup>



if  $\left. \begin{array}{l} \phi = \text{logistic} \\ w_{ij} \leftarrow \beta_j \end{array} \right\} \Rightarrow f_i(x) = P[y=1|x]$   
 in logistic regression

<sup>1</sup>More technically:  $\phi$  can be any continuous, bounded and strictly increasing function on  $\mathbb{R}$ .

## (Artificial) Neural Network (nn) unit

- ▶ For each **unit**  $i$

$$y_i \equiv f_i(x) = \phi\left(\sum_j w_{ij}x_j + w_{i0}\right) \quad (1)$$

- ▶ **Weigth vector**  $w_i$

- ▶  $w_{ij}$  = strength of the link from unit  $i$  to input  $j$
- ▶  $w_{ij} = 0$ : no link
- ▶  $w_{ij}$  can be positive or negative
- ▶ Sometimes we call the input vector  $x = [x_{1:d}]$  **input units**

- ▶ **activation function**  $\phi()$

- ▶ must be non-linear (otherwise the unit is a linear transformation)
- ▶ wanted: monotonically increasing, differentiable, gradient non-zero<sup>1</sup>

### Notation $\phi$ is overloaded

- ▶ When we talk about nn in general:  $\phi$  is any activation function
- ▶ When we do calculations with nn:  $\phi$  is by default the **logistic** function (unless specified otherwise)

$$\text{logistic or sigmoid function } \phi(u) = \frac{1}{1 + e^{-u}} \quad (2)$$

- ▶ When we do statistics or ML (but not nn):  $\phi$  is the logistic function
- ▶ Exercise: compare  $f_i(x)$  from (1) with  $p(x) = Pr[y = 1 | x]$  from logistic regression.

<sup>1</sup>More technically:  $\phi$  can be any continuous, bounded and strictly increasing function on  $\mathbb{R}$ .

# Perceptron

- Single layer feed-forward network

