# Decision Trees

CSC 461: Machine Learning

Fall 2022

Prof. Marco Alvarez
University of Rhode Island

# Introduction

# Supervised learning setup

‣ Data instance

  ✓ in general, $x \in \mathbb{R}^d$ is a **feature vector** of <u>discrete values</u>, but continuous values can also be handled
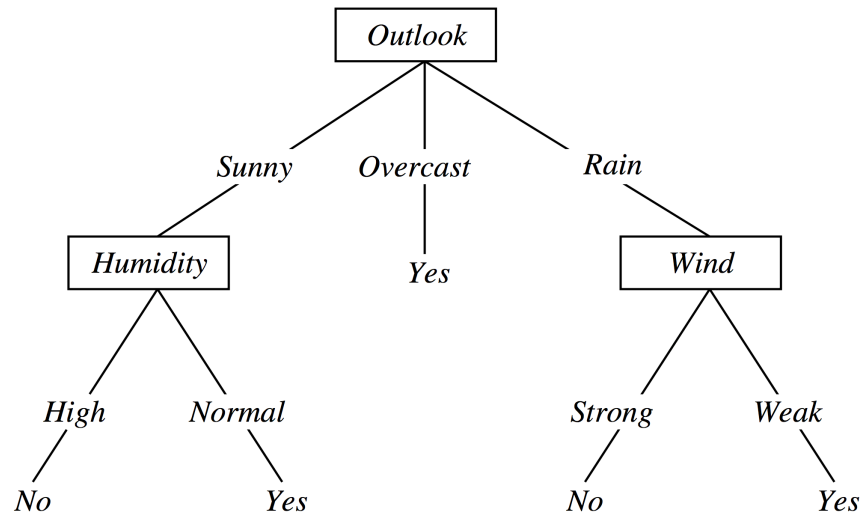
  ✓ $y \in \{1, 2, \ldots, k\}$

‣ Hypothesis

  ✓ each hypothesis **g** is a **decision tree**

$$g : \mathcal{X} \mapsto \mathcal{Y}, g \in \mathcal{H}$$

# Tennis dataset

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| D1 | Sunny | Hot | High | Weak | No |
| D2 | Sunny | Hot | High | Strong | No |
| D3 | Overcast | Hot | High | Weak | Yes |
| D4 | Rain | Mild | High | Weak | Yes |
| D5 | Rain | Cool | Normal | Weak | Yes |
| D6 | Rain | Cool | Normal | Strong | No |
| D7 | Overcast | Cool | Normal | Strong | Yes |
| D8 | Sunny | Mild | High | Weak | No |
| D9 | Sunny | Cool | Normal | Weak | Yes |
| D10 | Rain | Mild | Normal | Weak | Yes |
| D11 | Sunny | Mild | Normal | Strong | Yes |
| D12 | Overcast | Mild | High | Strong | Yes |
| D13 | Overcast | Hot | Normal | Weak | Yes |
| D14 | Rain | Mild | High | Strong | No |

Machine Learning, Tom Mitchell, McGraw Hill, 1997

## Example (Decision Tree)



Outlook
- Sunny → Humidity
  - High → No
  - Normal → Yes
- Overcast → Yes
- Rain → Wind
  - Strong → No
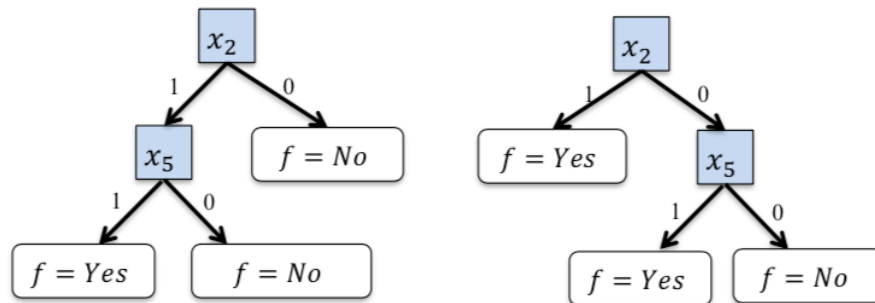  - Weak → Yes

Machine Learning, Tom Mitchell, McGraw Hill, 1997

## Representation

‣ Nodes test features/attributes

‣ Branches represent possible values for a feature

‣ Leaves represent outputs (classes)

‣ Assuming boolean inputs/outputs, draw the trees:

$A \wedge B$
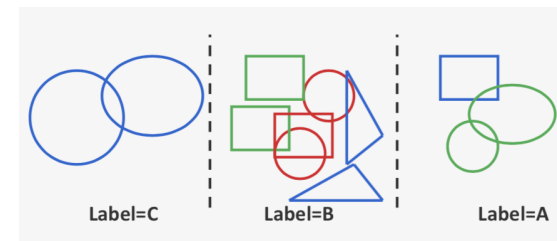$A \vee B$
$(A \wedge B) \vee (C \wedge \neg D \wedge E)$

## What functions are represented?

## Build/test your own tree

‣ Assume features: color and shape



Label=C     Label=B     Label=A

‣ What are the labels for a red triangle and a green triangle?

## Extracting rules from the tree
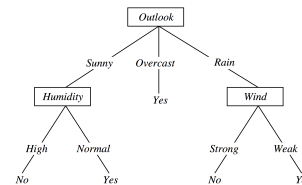
Look at the paths



Machine Learning, Tom Mitchell, McGraw Hill, 1997

## Disjunction of conjunctions

$$\dots \vee (\dots \wedge \dots) \vee (\dots \wedge \dots) \vee \dots$$

If … $(Outlook = Sunny \wedge Humidity = Normal) \vee$
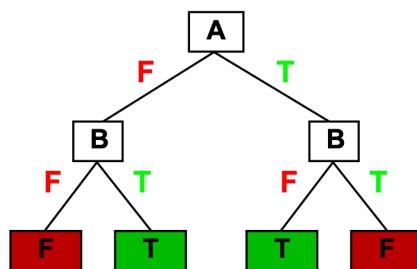$(Outlook = Overcast) \vee$
$(Outlook = Rain \wedge Wind = Weak)$

then it belongs to class YES



## Expressiveness

‣ A decision tree can represent any boolean/ discrete function (discrete input/discrete output)

| A | B | A xor B |
|---|---|---------|
| F | F | F |
| F | T | T |
| T | F | T |
| T | T | F |



http://aima.eecs.berkeley.edu/slides-pdf/chapter18.pdf

## Hypothesis space

How many distinct decision trees can be created with d=5 boolean features?

| x | | | | | y |
|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | T/F |
| 0 | 0 | 0 | 0 | 1 | T/F |
| 0 | 0 | 0 | 1 | 0 | T/F |
| 0 | 0 | 0 | 1 | 1 | T/F |
| 0 | 0 | 1 | 0 | 0 | T/F |
| … | | | | | T/F |
| 1 | 1 | 1 | 1 | 1 | T/F |

$2^5 = 32$ entries

how many boolean functions with 5 features are there, given that entries can be T/F?

$$2^{2^5}$$

Try d == 10

# Hypothesis space

‣ More expressive hypothesis space …

  ✓ allows learning complex target functions

  ✓ increases number of consistent hypotheses

  ✓ may not **generalize**, due to **overfitting**

‣ DT learning goal

  ✓ find a **small tree** <u>consistent</u> **with the training data**

  ✓ **NP-complete** (polynomial algorithm may not exist)

# Consistent hypotheses

‣ A hypothesis **g** is consistent with a set of training examples **D** if and only if **g(x) = y** for all pairs **(x, y)** in **D**

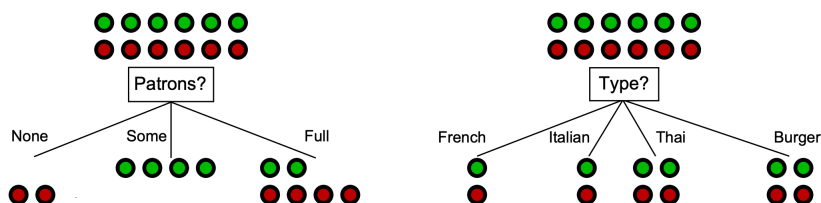  ✓ **"hope"**: if **g** is consistent with training data, then it would be accurate on new instances

‣ There is a tree consistent with any training set (just list all paths) — **it may not generalize well**
‣ Preferably we want more **compact** trees that can **generalize** better

# Learning a Decision Tree

# Induction of a Decision Tree

‣ Build the tree using a **top-down** approach

  ✓ selecting one feature to split at a time

‣ **Greedy** algorithm

  ✓ makes the **optimal** choice at each step **(which feature to split)**

  ✓ the greedy nature of the algorithm cannot guarantee **optimality (smallest tree consistent with the data)**

‣ **NP-complete** problem

  ✓ "Although a solution to an NP-complete problem can be verified "quickly", there is no known way to find a solution quickly" [wikipedia]

# Which feature is better? Why?



Patrons?

None    Some    Full

Type?

French    Italian    Thai    Burger

- Which feature is more **informative**?

- Which provides the minimum **0/1 loss** if we use the majority vote for classifying new instances?

---

# Which feature is better?

$[29+,35-]$    A1=?

t    f

$[21+,5-]$    $[8+,30-]$

$[29+,35-]$    A2=?

t    f

$[18+,33-]$    $[11+,2-]$

---

# How to choose the splitting feature?

‣ Information Gain

  ✓ used in **ID3**

‣ Gain Ratio

  ✓ used in C4.5

‣ Gini Measure
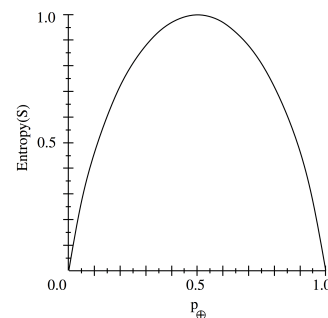
  ✓ used in CART

‣ …

**ID3 was invented by Ross Quinlan**



---

# Entropy

‣ Assume a set S of positive/negative instances

  ✓ entropy measures the **impurity** of S

w.r.t. a binary variable

$$E(S) = -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

# Entropy

‣ Assuming **k** possible values, each with different probabilities, then:

$$E(S) = -\sum_{i=1}^{k} p_i \log_2 p_i$$

> What is the entropy if all instances belong to the same category?
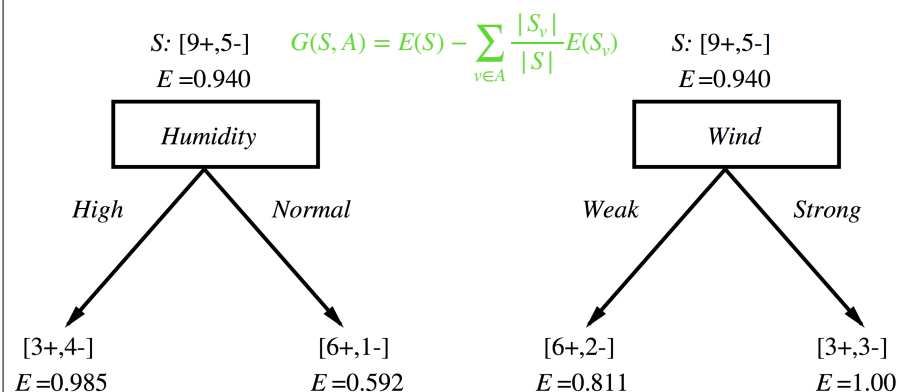
# Information Gain

‣ Expected reduction in **Entropy** after splitting

$$G(S,A) = E(S) - \sum_{v \in A} \frac{|S_v|}{|S|} E(S_v)$$

‣ **Information gain** increases for low entropy values

# Calculate the Information Gain

$$S: [9+,5-] \qquad G(S,A) = E(S) - \sum_{v \in A} \frac{|S_v|}{|S|} E(S_v) \qquad S: [9+,5-]$$

$E = 0.940$                                                 $E = 0.940$

|  Humidity  |          |  Wind  |          |
|------------|----------|--------|----------|
| High       | Normal   | Weak   | Strong   |
| [3+,4-]    | [6+,1-]  | [6+,2-]| [3+,3-]  |
| E = 0.985  | E = 0.592| E = 0.811 | E = 1.00 |

# Induction of a Decision Tree

---
**Algorithm** GrowTree($D, F$)

---
**Input** : data $D$; set of features $F$.
**Output** : feature tree $T$ with labelled leaves.
**if** Homogeneous($D$) **then return** Label($D$) ;    // Homogeneous, Label: see text
$S \leftarrow$ BestSplit($D, F$) ;              // e.g., BestSplit-Class (Algorithm 5.2)
split $D$ into subsets $D_i$ according to the literals in $S$;
**for** each $i$ **do**
    **if** $D_i \neq \emptyset$ **then** $T_i \leftarrow$ GrowTree($D_i, F$) **else** $T_i$ is a leaf labelled with Label($D$);
**end**
**return** a tree whose root is labelled with $S$ and whose children are $T_i$

---

## Induction of a Decision Tree

---

**Algorithm** BestSplit-Class($D, F$) – find the best split for a decision tree.

---

**Input** : data $D$; set of features $F$.

**Output** : feature $f$ to split on.

$I_{min} \leftarrow 1$;

**for** each $f \in F$ **do**

    split $D$ into subsets $D_1, \ldots, D_l$ according to the values $v_j$ of $f$;

    **if** $\text{Imp}(\{D_1, \ldots, D_l\}) < I_{min}$ **then**

        $I_{min} \leftarrow \text{Imp}(\{D_1, \ldots, D_l\})$;
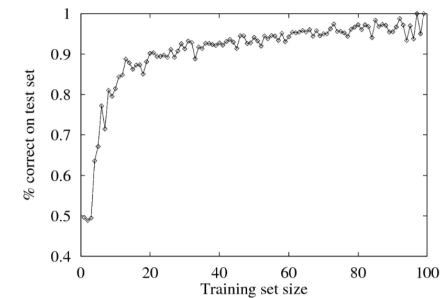
        $f_{best} \leftarrow f$;
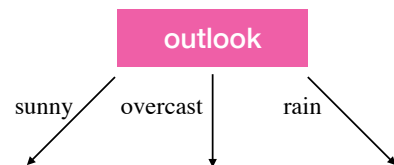
    **end**

**end**

**return** $f_{best}$

---

## Resulting tree

‣ Tree is expected to be small and consistent with training examples

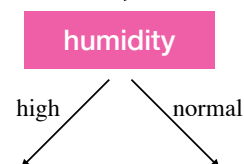‣ Tree does not necessarily agree with the correct function (bigger training datasets help)



## Example

| Day | Outlook | Temperature | Humidity | Wind | Play |
|-----|---------|-------------|----------|------|------|
| D1 | sunny | hot | high | weak | no |
| D2 | sunny | hot | high | strong | no |
| D3 | overcast | hot | high | weak | yes |
| D4 | rain | mild | high | weak | yes |
| D5 | rain | cool | normal | weak | yes |
| D6 | rain | cool | normal | strong | no |
| D7 | overcast | cool | normal | strong | yes |
| D8 | sunny | mild | high | weak | no |
| D9 | sunny | cool | normal | weak | yes |
| D10 | rain | mild | normal | weak | yes |
| D11 | sunny | mild | normal | strong | yes |
| D12 | overcast | mild | high | strong | yes |
| D13 | overcast | hot | normal | weak | yes |
| D14 | rain | mild | high | strong | no |

**Panel 1**

outlook — sunny, overcast, rain

| Day | Outlook | Temperature | Humidity | Wind | Play |
|-----|---------|-------------|----------|------|------|
| D1 | sunny | hot | high | weak | no |
| D2 | sunny | hot | high | strong | no |
| D3 | overcast | hot | high | weak | yes |
| D4 | rain | mild | high | weak | yes |
| D5 | rain | cool | normal | weak | yes |
| D6 | rain | cool | normal | strong | no |
| D7 | overcast | cool | normal | strong | yes |
| D8 | sunny | mild | high | weak | no |
| D9 | sunny | cool | normal | weak | yes |
| D10 | rain | mild | normal | weak | yes |
| D11 | sunny | mild | normal | strong | yes |
| D12 | overcast | mild | high | strong | yes |
| D13 | overcast | hot | normal | weak | yes |
| D14 | rain | mild | high | strong | no |

calculate IG and select highest

**Panel 2**

outlook — sunny, overcast, rain
sunny → humidity — high, normal

| Day | Temperature | Humidity | Wind | Play |
|-----|-------------|----------|------|------|
| D1 | hot | high | weak | no |
| D2 | hot | high | strong | no |
| D8 | mild | high | weak | no |
| D9 | cool | normal | weak | yes |
| D11 | mild | normal | strong | yes |

calculate IG and select highest

**Panel 3**

outlook — sunny, overcast, rain
sunny → humidity — high → NO, normal

| Day | Temperature | Wind | Play |
|-----|-------------|------|------|
| D1 | hot | weak | no |
| D2 | hot | strong | no |
| D8 | mild | weak | no |

calculate IG and select highest

**Panel 4**

outlook — sunny, overcast, rain
sunny → humidity — high → NO, normal → YES

| Day | Temperature | Wind | Play |
|-----|-------------|------|------|
| D9 | cool | weak | yes |
| D11 | mild | strong | yes |

calculate IG and select highest

**Panel 1 (top-left) — tree: outlook → {sunny → humidity → high: NO, normal: YES; overcast → YES; rain → }**

| Day | Temperature | Humidity | Wind | Play |
|-----|-------------|----------|------|------|
| D3 | hot | high | weak | yes |
| D7 | cool | normal | strong | yes |
| D12 | mild | high | strong | yes |
| D13 | hot | normal | weak | yes |

calculate IG and select highest

**Panel 2 (top-right) — tree: outlook → {sunny → humidity → high: NO, normal: YES; overcast → YES; rain → wind → weak, strong}**

| Day | Temperature | Humidity | Wind | Play |
|-----|-------------|----------|------|------|
| D4 | mild | high | weak | yes |
| D5 | cool | normal | weak | yes |
| D6 | cool | normal | strong | no |
| D10 | mild | normal | weak | yes |
| D14 | mild | high | strong | no |

calculate IG and select highest

**Panel 3 (bottom-left) — tree: outlook → {sunny → humidity → high: NO, normal: YES; overcast → YES; rain → wind → weak: YES, strong}**

| Day | Temperature | Humidity | Play |
|-----|-------------|----------|------|
| D4 | mild | high | yes |
| D5 | cool | normal | yes |
| D10 | mild | normal | yes |

calculate IG and select highest

**Panel 4 (bottom-right) — tree: outlook → {sunny → humidity → high: NO, normal: YES; overcast → YES; rain → wind → weak: YES, strong: NO}**

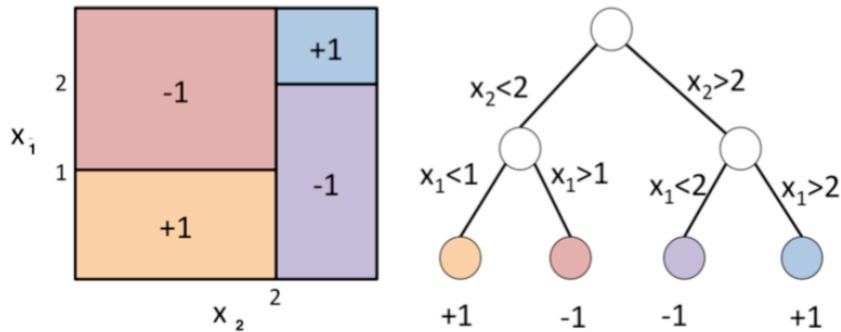| Day | Temperature | Humidity | Play |
|-----|-------------|----------|------|
| D6 | cool | normal | no |
| D14 | mild | high | no |

calculate IG and select highest
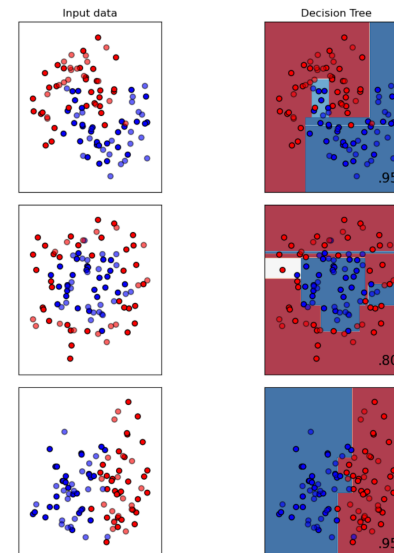
# Final Remarks

# Continuous features

‣ Transform continuous into discrete features

  ✓ use thresholds defined by domain experts or automatically calculated from training data

‣ For example:

  ✓ sort values (training set)

  ✓ find split points where class changes

| Temperature: | 40 | 48 | 60 | 72 | 80 | 90 |
|---|---|---|---|---|---|---|
| PlayTennis: | No | No | Yes | Yes | Yes | No |

54          85

# Nonlinear Decision Boundary



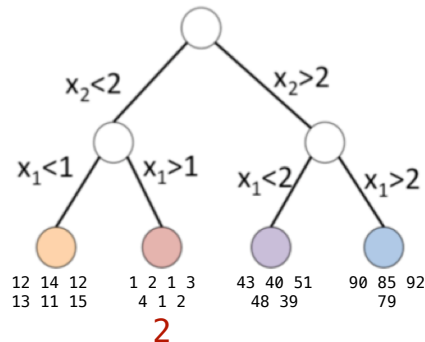from: CS260 Machine Learning Algorithms, Cho-Jui Hsieh, UCLA, 2019
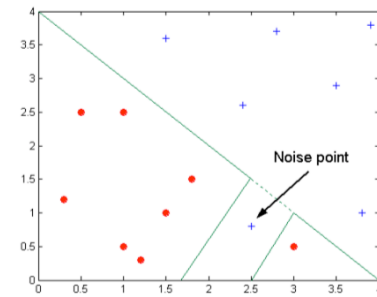
# Nonlinear Decision Boundary

## Continuous outputs

‣ Regression trees

  ✓ can assign a continuous value to a leaf

  ✓ e.g. the **average** of all **y** values that fall into the leaf



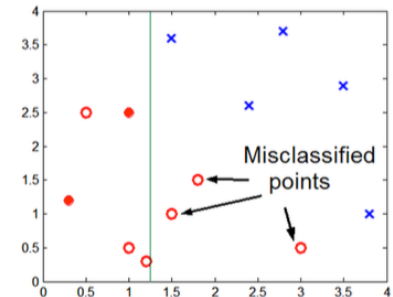|  | $x_2<2$ | | $x_2>2$ | |
| --- | --- | --- | --- | --- |
| | $x_1<1$ | $x_1>1$ | $x_1<2$ | $x_1>2$ |
| | 12 14 12 | 1 2 1 3 | 43 40 51 | 90 85 92 |
| | 13 11 15 | 4 1 2 | 48 39 | 79 |

2

## Model overfitting

A hypothesis **h1** is said to **overfit** the training data if there exists some alternative hypothesis **h2** such that **h1** has smaller error than **h2** over the training examples, but **h2** has a smaller error than **h1** over the entire distribution of instances



due to noise                      due to lack of representative instances

from: Data Mining I, Eirini Ntoutsi, Leibniz University, Summer 2019

## Preventing overfitting (DTs)

‣ Remove irrelevant features

‣ Add more data

‣ **Stop growing branches** during training

  ✓ hard thresholds or statistical measures

‣ Prune the tree post-training

## Additional thoughts on DTs

‣ **Nonlinear** classifiers, which can also provide **interpretability**

‣ Training may be **slow** but inference is **fast**

  ✓ what is the big-O of inference?

‣ Although trees can be small, certain functions will require an exponentially large decision tree

  ✓ e.g. **majority** (1 if n inputs are positive), **parity** (1 if even number of inputs is positive)