CSC 480/580 Principles of Machine Learning
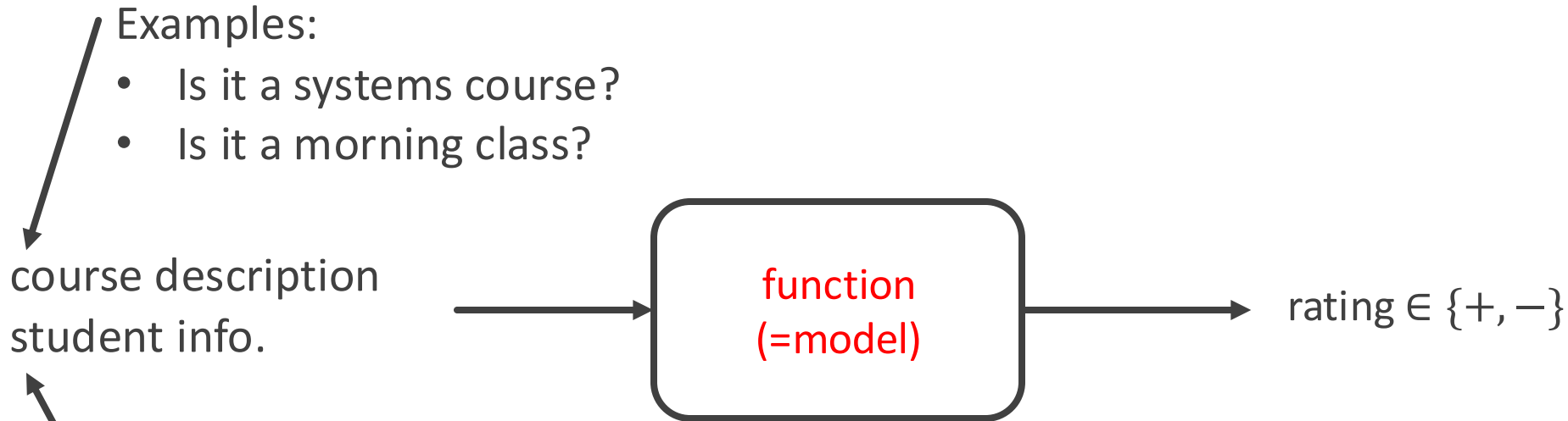
# 01 Decision Trees

**Jason Pacheco**

THE UNIVERSITY OF ARIZONA

# Example: course recommendation

- Build a software: given a student, recommend a set of courses that s/he would like

Examples:
- Is it a systems course?
- Is it a morning class?

course description
student info.

function
(=model)

rating $\in \{+, -\}$

Examples:
- which courses has taken before?
- likes morning class?

I'll explain
1. What kind of functions we'll be using
2. How to train one from data

# Model: Decision Tree
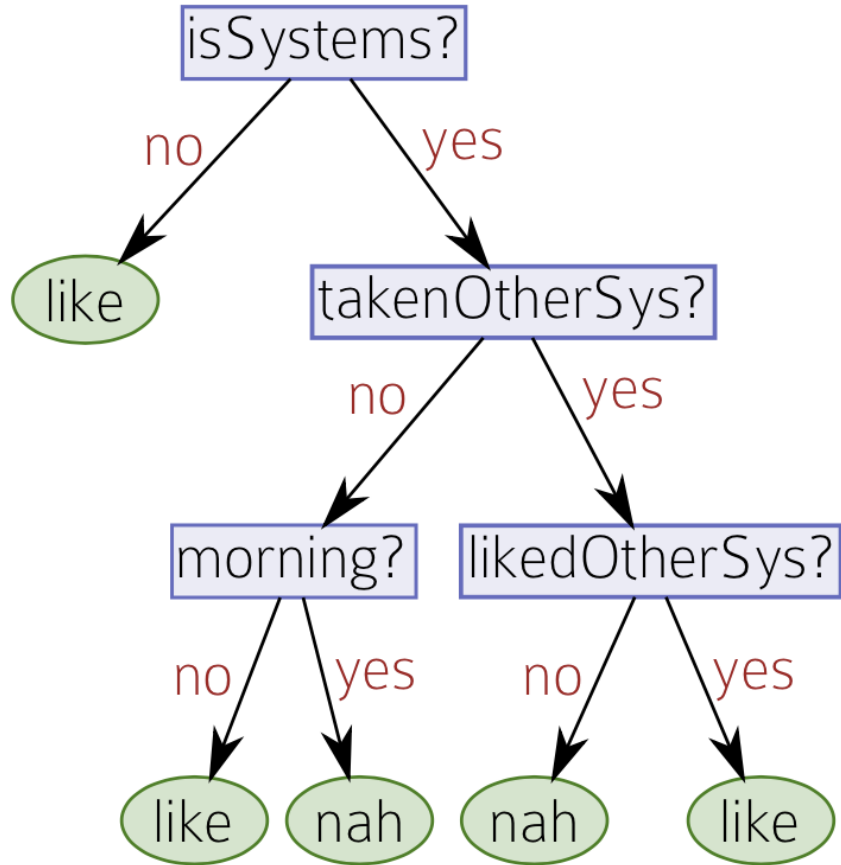
# Model: Decision Tree: Example



Figure 1.2: A decision tree for a course recommender system, from which the in-text "dialog" is drawn.

**Input**: the course & student info

Use questions to arrive at a conclusion.

**Terminology:**

- (Question, Answer) → (Feature, Feature Value)
- "Like" / "Nah" → Label
- {(A set of (Question & Answer)'s, Label)} → Train Data
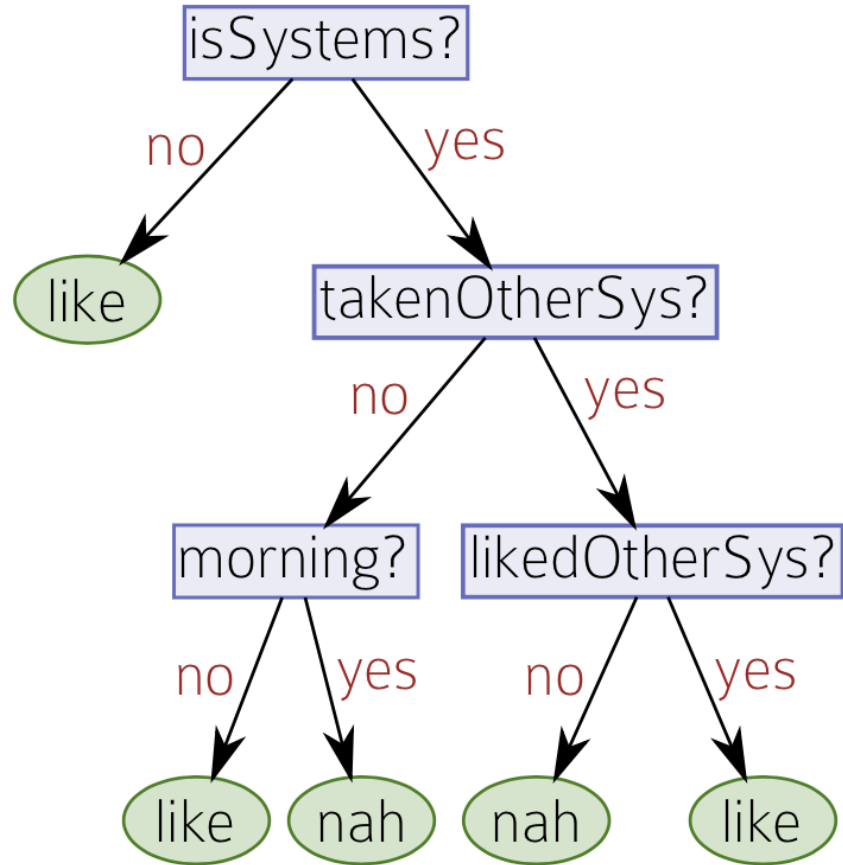
# Basic tree terminology



Figure 1.2: A decision tree for a course recommender system, from which the in-text "dialog" is drawn.

node                    parent
root node               children
leaf node               ancestor
internal node           subtree
                        depth

Q: How many nodes are there?

Q: What's the depth of this tree?

- Key advantage of decision trees: *intepretability*

- Useful in consequential settings, e.g. medical treatment, loan approval, etc.

nodes organized in a tree-based structure, leading to a prediction (Fig. 1). The interpretability of decision trees allows physicians to understand why a prediction or stratification is being made, providing an account of the reasons behind the decision to subsequently accept or override the model's output. This interaction between humans and algorithms can provide

# Prediction using a decision tree

- Test: predict using a decision tree:

test point: the data point to be classified
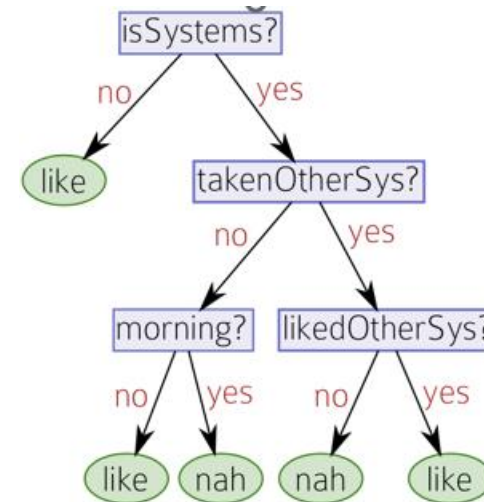(vs train point: data point to be used for training)

**Algorithm 2** DecisionTreeTest(*tree, test point*)

1: **if** *tree* is of the form Leaf(*guess*) **then**
2:     **return** *guess*
3: **else if** *tree* is of the form Node(*f, left, right*) **then**
4:     **if** $f = no$ in *test point* **then**
5:         **return** DecisionTreeTest(*left, test point*)
6:     **else**
7:         **return** DecisionTreeTest(*right, test point*)
8:     **end if**
9: **end if**

guess = prediction

left = no
right = yes

- Training: how to design a learning algorithm $\mathcal{A}$ that can build trees $f$ from training data?

6

# How to train

# Train dataset

feature vector $\in \mathbb{R}^d$    label $\in \{+,-\}$

*Define the labeled train data* $S = \{(x_i, y_i)\}_{i=1}^{n}$

Features can be a function of the <u>user</u> being recommended; e.g., are you a morning person?

**Features**

**Feature Values**

To make this a <u>binary</u> classification, we map

$\{+2,+1,0\} \Rightarrow$ "Liked" (+)
$\{-1,-2\} \Rightarrow$ "Nah" (-)

**Labels**

**(Labeled) Data Point**

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|--------|-------|-----|------|------|----------|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

# Background: Majority vote classifier

The most basic classifier you can think of.

$\in \mathbb{R}^d$

$\in \{1, \dots, C\}$

How to train:

- Given: A (training) dataset with **n** data points $\{(x_i, y_i)\}_{i=1}^n$ with **C** classes.

- Compute the most common class $c^*$ in the dataset.

$$c^* := \arg \max_{c \in \{1, \dots, C\}} \sum_{i=1}^n \mathbf{I}\{y_i = c\}$$

(break ties arbitrarily)

$$\mathbf{I}\{A\} := \begin{cases} 1 & \text{if } A \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

"indicator function"

- Output a classifier $f(x) = c^*$.

Stupid enough classifier! Always try to beat this classifier.

Often, state-of-the-art ML algorithms perform barely better than the majority vote classifier..
$\Rightarrow$ happens when there is no association between features and labels in the dataset

# Background: Train set accuracy/error

- Suppose the ML algorithm has trained a function $f$ using the dataset $D = \{(x_i, y_i)\}_{i=1}^{n}$

- Train set accuracy:

$$\widehat{acc}(f) := \frac{1}{n}\sum_{i=1}^{n} \mathbf{I}\{f(x_i) = y_i\}$$

- Train set error: $\widehat{err}(f) = \frac{1}{n}\sum_{i=1}^{n} \mathbf{I}\{f(x_i) \neq y_i\} = 1 - \widehat{acc}(f)$

- Q: We have 100 train set (images) consisting of 5 cats, 80 dogs, and 15 lions. What is the train set accuracy of the majority vote classifier? What is the error?

# Training: The ideal criterion

- The training data $D = \{(x_i, y_i)\}_{i=1}^n$ $\qquad x_i \in \{y, n\}^d \qquad y_i \in \{+, -\}$

$$\hat{f} := \arg \max_{f \in \text{DecisionTrees}} \frac{1}{n} \sum_{i=1}^n \mathbf{I}\{f(x_i) = y_i\}$$

The main principle governing most of the ML algorithms.

It's called "**empirical risk minimization (ERM)**" ( empirical risk = training set error )
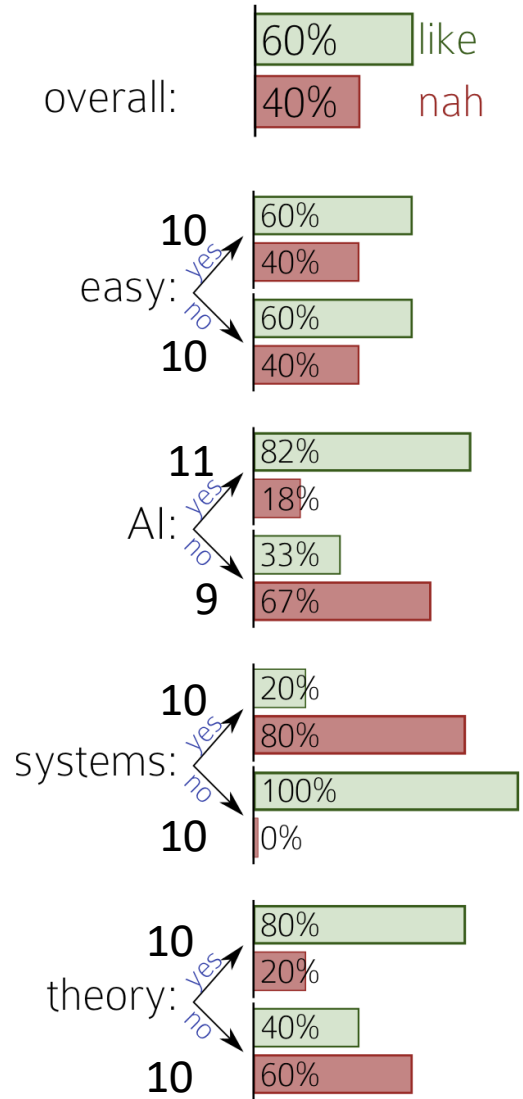
**The issue**:
- Naïve search: $O(d^d)$ time complexity
- It's <u>NP-Hard</u> -- don't expect to have an efficient algorithm.

**Solution**:
- Perform greedy approximation!
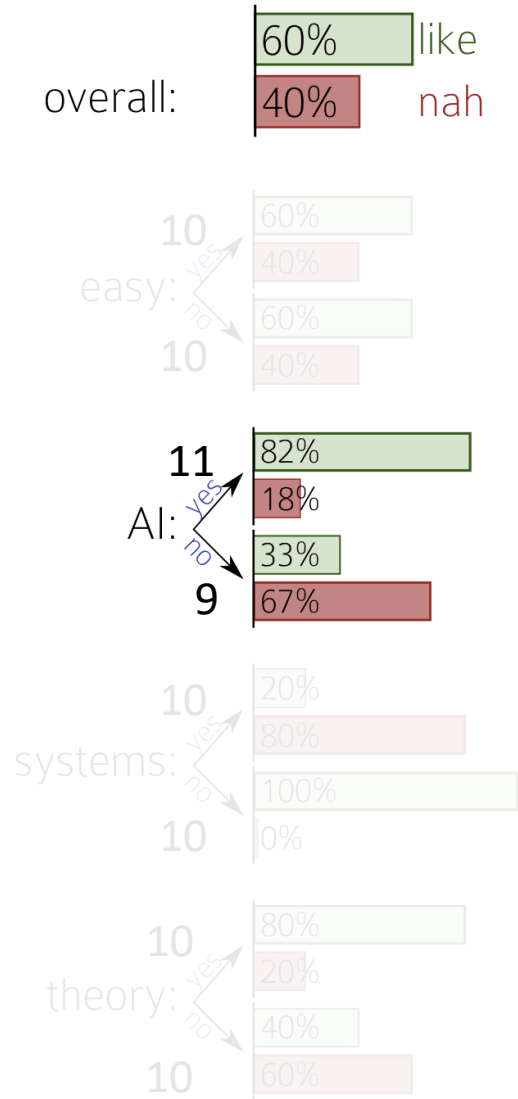
# How to train a decision tree

overall:
60% like
40% nah

easy:
10 — 60% / 40%
10 — 60% / 40%

AI:
11 — 82% / 18%
9 — 33% / 67%

systems:
10 — 20% / 80%
10 — 100% / 0%

theory:
10 — 80% / 20%
10 — 40% / 60%

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|--------|-------|-----|------|------|----------|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

Baseline: 'majority vote' classifier

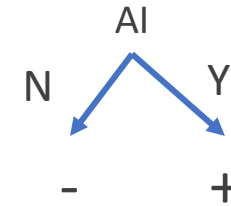Q: What is the train set accuracy?  0.60

# How to train a decision tree

**overall:** 60% like / 40% nah

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|---|---|---|---|---|---|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| o | n | n | n | n | y |
| o | y | n | n | y | y |
| o | n | y | n | y | n |
| o | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

easy: 10 → 60% / 40%, 10 → 60% / 40%

AI: 11 → 82% / 18%, 9 → 33% / 67%

systems: 10 → 20% / 80%, 10 → 100% / 0%

theory: 10 → 80% / 20%, 10 → 40% / 60%

Baseline: 'majority vote' classifier

Q: What is the train set accuracy? **0.60**

Suppose we place the node AI at the root. Let us set the prediction at each leaf node as the majority vote.

AI
N → -
Y → +

What is the train set accuracy now? Use weighted average.
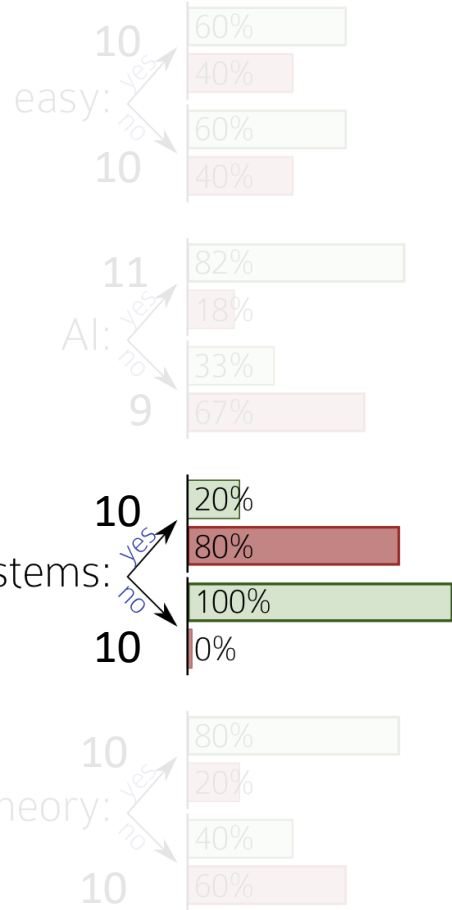
$$\frac{9}{20}\cdot\frac{6}{9} + \frac{11}{20}\cdot\frac{9}{11} = \frac{15}{20} = 0.75 \quad \text{improved!}$$

AI=N        AI=Y
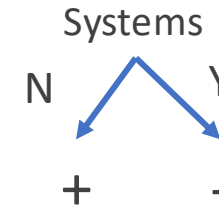
# How to train a decision tree

overall:
60% like
40% nah

easy:
10 — 60% / 40%
10 — 60% / 40%

AI:
11 — 82% / 18%
9 — 33% / 67%

systems:
10 yes — 20% / 80%
10 no — 100% / 0%

theory:
10 yes — 80% / 20%
10 no — 40% / 60%

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|---|---|---|---|---|---|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

Suppose placing the node Systems at the root.
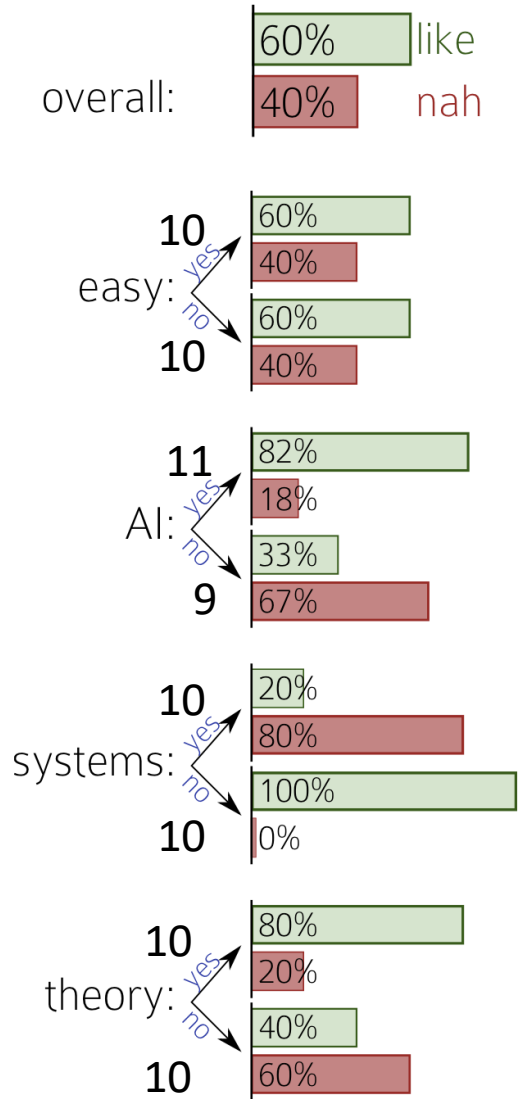
Systems
N / Y
+    -

What is the train set accuracy now?

$$\frac{10}{20} \cdot \frac{10}{10} + \frac{10}{20} \cdot \frac{8}{10} = \frac{18}{20} = 0.9 \quad \text{even better!}$$
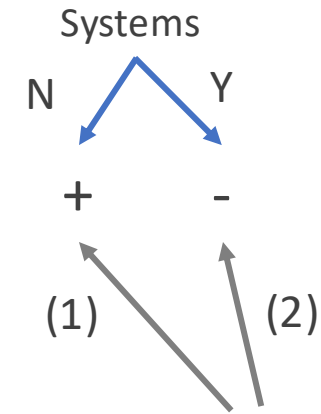
What would you do to build a depth-1 tree?

try out each feature and choose the one that leads to the largest accuracy!
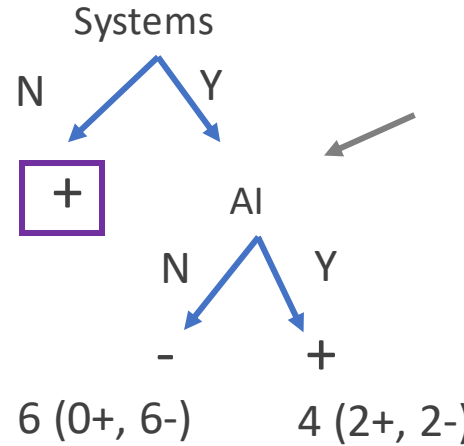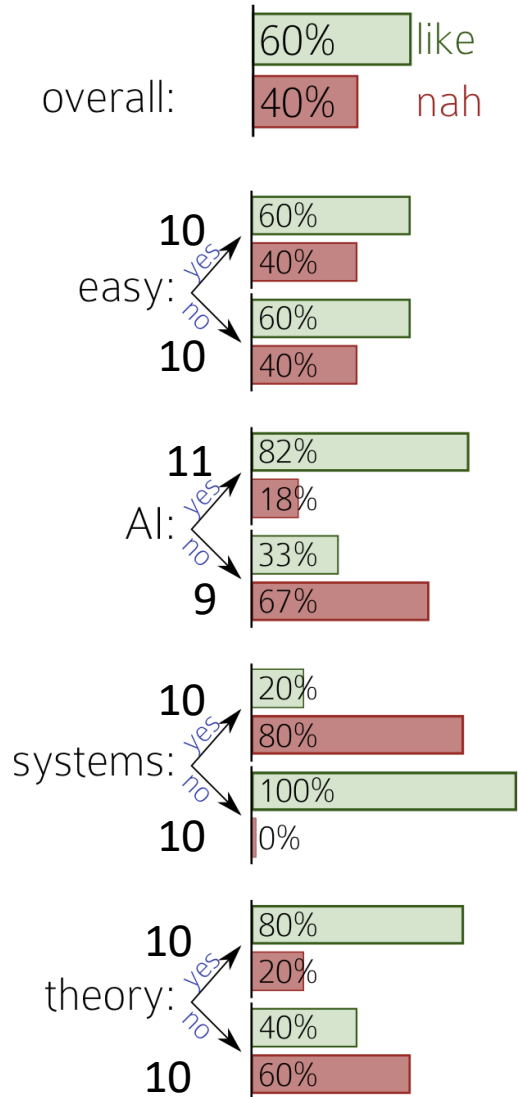
# How to train a decision tree

overall:
- 60% like
- 40% nah

easy:
- 10 → yes: 60% / 40%
- 10 → no: 60% / 40%

AI:
- 11 → yes: 82% / 18%
- 9 → no: 33% / 67%

systems:
- 10 → yes: 20% / 80%
- 10 → no: 100% / 0%

theory:
- 10 → yes: 80% / 20%
- 10 → no: 40% / 60%

What about depth 2?

Systems

N          Y

+          -

(1)        (2)

Which nodes to put at each leaf node?

Focus on (2). Try placing AI

# How to train a decision tree

overall:
60% like
40% nah

easy:
10 → yes: 60% / 40%
10 → no: 60% / 40%

AI:
11 → yes: 82% / 18%
9 → no: 33% / 67%

systems:
10 → yes: 20% / 80%
10 → no: 100% / 0%

theory:
10 → yes: 80% / 20%
10 → no: 40% / 60%

Systems
N    Y

+

AI
N    Y

-         +

6 (0+, 6-)    4 (2+, 2-)

Q: How many training data points fall here?   10

Q: How many training data points arrive at these two leaves? How many for each label?

Q: what prediction should we use for each leaf?

Q: What is the train set accuracy, conditioning on Systems=Y?

'local' train set accuracy

$$\frac{6}{10} \cdot \frac{6}{6} + \frac{4}{10} \cdot \frac{2}{4} = \frac{8}{10}$$

Try all the other nodes and pick the one with the largest (local) acc.!
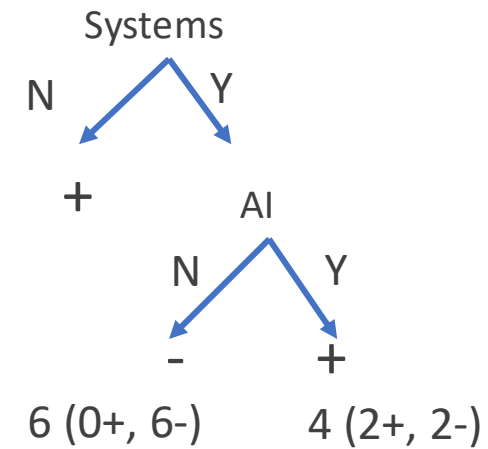
Then, repeat the same for Systems=N branch!

⇒ But this has 1.0 local train set acc. No need to expand anymore!

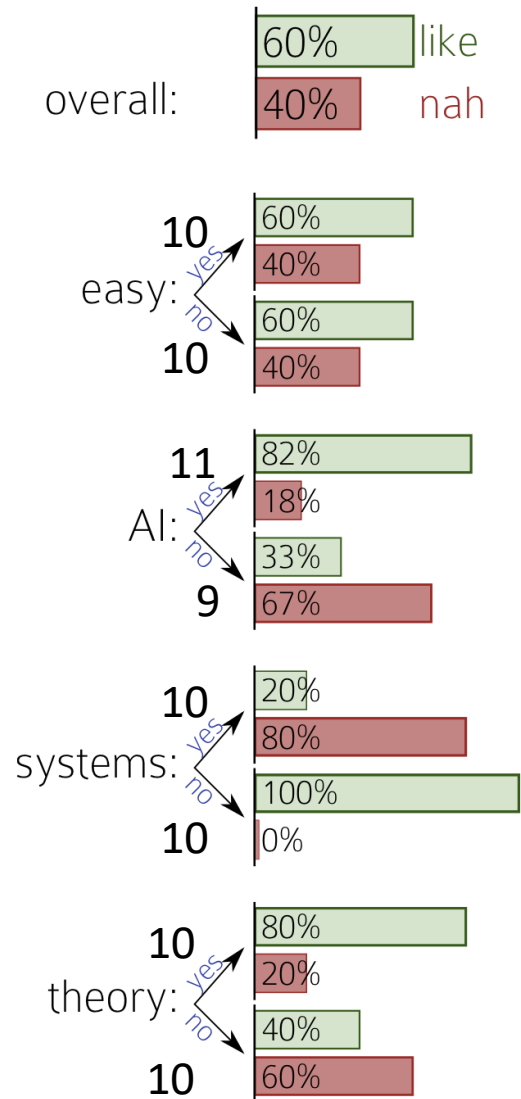Move onto expanding nodes at depth 2!

# How to train a decision tree

| Rating | Easy? | AI? | Sys? | Thy? | Morning? |
|--------|-------|-----|------|------|----------|
| +2 | y | y | n | y | n |
| +2 | y | y | n | y | n |
| +2 | n | y | n | n | n |
| +2 | n | n | n | y | n |
| +2 | n | y | y | n | y |
| +1 | y | y | n | n | n |
| +1 | y | y | n | y | n |
| +1 | n | y | n | y | n |
| 0 | n | n | n | n | y |
| 0 | y | n | n | y | y |
| 0 | n | y | n | y | n |
| 0 | y | y | y | y | y |
| -1 | y | y | y | n | y |
| -1 | n | n | y | y | n |
| -1 | n | n | y | n | y |
| -1 | y | n | y | n | y |
| -2 | n | n | y | y | n |
| -2 | n | y | y | n | y |
| -2 | y | n | y | n | n |
| -2 | y | n | y | n | y |

Systems

N          Y

+          AI

N          Y

-          +

6 (0+, 6-)        4 (2+, 2-)

# How to train a decision tree

Overall idea:
1. Set the root node as a leaf node.
2. Grab a leaf node for whose 'local' train accuracy is not 1.0.
3. Loop through features to find a feature $f^*$ that maximizes the 'local' train accuracy and replace the leaf node with a node with feature $f^*$; add its leaf nodes and set their predictions by majority vote.
   (note: skip the features used by an ancestor)
4. Repeat 2-3 until there is no more 'expandable' leaf node.

(i)   local train acc. is not 1.0 <u>and</u>
(ii)  ancestors did not use all the features yet

**Algorithm 1** DECISIONTREETRAIN(*data, remaining features*)

1: *guess* ← most frequent answer in *data*                    // default answer for this data      guess=majority vote
2: **if** the labels in *data* are unambiguous **then**
3:     **return** LEAF(*guess*)                              // base case: no need to split further      unambiguous
4: **else if** *remaining features* is empty **then**                                          = achieves 100% local acc. when
5:     **return** LEAF(*guess*)                              // base case: cannot split further      using the majority vote
6: **else**                                      // we need to query more features
7:     **for all** $f \in$ *remaining features* **do**
8:         $NO$ ← the subset of *data* on which $f$=*no*
9:         $YES$ ← the subset of *data* on which $f$=*yes*
10:        *score*[$f$] ← # of majority vote answers in $NO$            has the same role as computing
11:                    + # of majority vote answers in $YES$
                                                    $$\frac{|NO|}{|YES| + |NO|} \, \widehat{acc}\,(NO) + \frac{|YES|}{|YES| + |NO|} \, \widehat{acc}\,(YES)$$
                    // the accuracy we would get if we only queried on *f*
12:    **end for**
13:    $f$ ← the feature with maximal *score(f)*
14:    $NO$ ← the subset of *data* on which $f$=*no*
15:    $YES$ ← the subset of *data* on which $f$=*yes*
16:    *left* ← DECISIONTREETRAIN($NO$, *remaining features* \ {$f$})
17:    *right* ← DECISIONTREETRAIN($YES$, *remaining features* \ {$f$})
18:    **return** NODE($f$, *left*, *right*)
19: **end if**

19

# Type of features

- Binary

- Categorical: values in $\{1, \ldots, C\}$    e.g., occupation, blood type
  - Option 1: Instead of 2 children, have C children.
  - Option 2: Derive C features of the form "feature=c?" for every $c \in C$.
    ↑ binary features!

    Q: How about features of the form "feature$\in D$" for every $D \subset C$?    computational complexity ↑

- Real value   e.g., weight, age, price
  - Sort the values.
  - Find the **breakpoints**: For every two adjacent points with opposite labels, compute the midpoint.
  - Derive features like "weight $\leq$ breakpoint"

# Types of labels

- Binary
  - Accuracy is not sensitive to node purity...we will look at alternatives


- Multiclass: What changes do we need to make?
  - Almost none! Just extend the definition of accuracy to multiclass.

$$\widehat{\mathrm{acc}}(f) := \frac{1}{n} \sum_{i=1}^{n} \mathbf{I}\{f(x_i) = y_i\}$$


- Real Value
  - This is a regression problem...we will get back to this

# Variations: binary case

## Notions of uncertainty: binary case $(\mathcal{Y} = \{0, 1\})$

Suppose in a set of examples $S \subseteq \mathcal{X} \times \{0, 1\}$, a $p$ fraction are labeled as 1
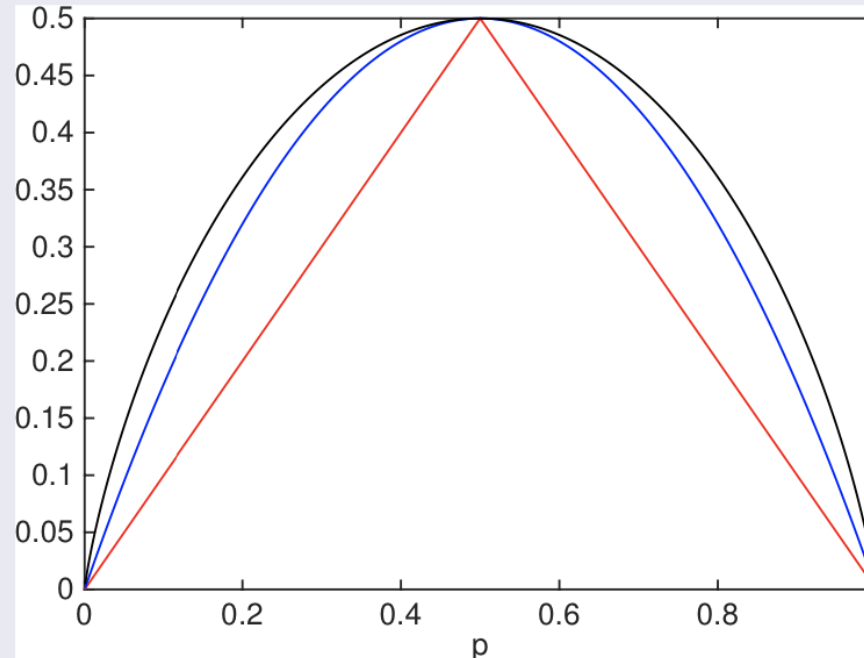
**❶ Classification error**: (red)

$$u(S) := \min\{p, \, 1-p\}$$

**❷ Gini index**: (blue)

$$u(S) := 2p(1-p)$$

**❸ Entropy**: (black)

$$u(S) := p \log \frac{1}{p} + (1-p) \log \frac{1}{1-p}$$



Gini index and entropy (after some rescaling) are concave upper-bounds on classification error

"label uncertainty"

classification error here
$= 1 -$ accuracy
(verify yourself)

Let $q$ is the fraction of data points with feature=Y.

**Modification**:
Set score[f] as
$$q \cdot \big(-u(YES)\big) + (1-q) \cdot \big(-u(NO)\big)$$

# If the number of classes is >2

## Notions of uncertainty: general case

Suppose in $S \subseteq \mathcal{X} \times \mathcal{Y}$, a $p_k$ fraction are labeled as $k$ (for each $k \in \mathcal{Y}$).

1. **Classification error**:
$$u(S) := 1 - \max_{k \in \mathcal{Y}} p_k$$

2. **Gini index**:
$$u(S) := 1 - \sum_{k \in \mathcal{Y}} p_k^2$$

3. **Entropy**:
$$u(S) := \sum_{k \in \mathcal{Y}} p_k \log \frac{1}{p_k}$$

Each is *maximized* when $p_k = 1/|\mathcal{Y}|$ for all $k \in \mathcal{Y}$
(i.e., equal numbers of each label in $S$)

Each is *minimized* when $p_k = 1$ for a single label $k \in \mathcal{Y}$
(so $S$ is **pure** in label)

# Regression

- Classification vs Regression
  - Both supervised learning
  - Regression has <u>real-valued</u> labels.

- Examples: Price prediction. Property value prediction.

- Standard measure of performance: mean squared error: $\frac{1}{n}\sum_{i=1}^{n}(f(x_i)-y_i)^2$

  Q: why are we using <u>squared error</u> (f-y)^2 rather than <u>absolute error</u> |f-y|?   my opinion: convenience & tradition

- Changes needed:
  - How to make predictions at the leaf node?

    Average labels of the data at the leaf; denote by $\bar{y}_{YES}$ and $\bar{y}_{NO}$.

  - How to adjust score[f]?

    Use negative squared error

$$\frac{|YES|}{|YES|+|NO|}\cdot\left(-\frac{1}{|YES|}\sum_{i\in YES}(\bar{y}_{YES}-y_i)^2\right)+\frac{|NO|}{|YES|+|NO|}\left(-\frac{1}{|NO|}\sum_{i\in NO}(\bar{y}_{NO}-y_i)^2\right)$$
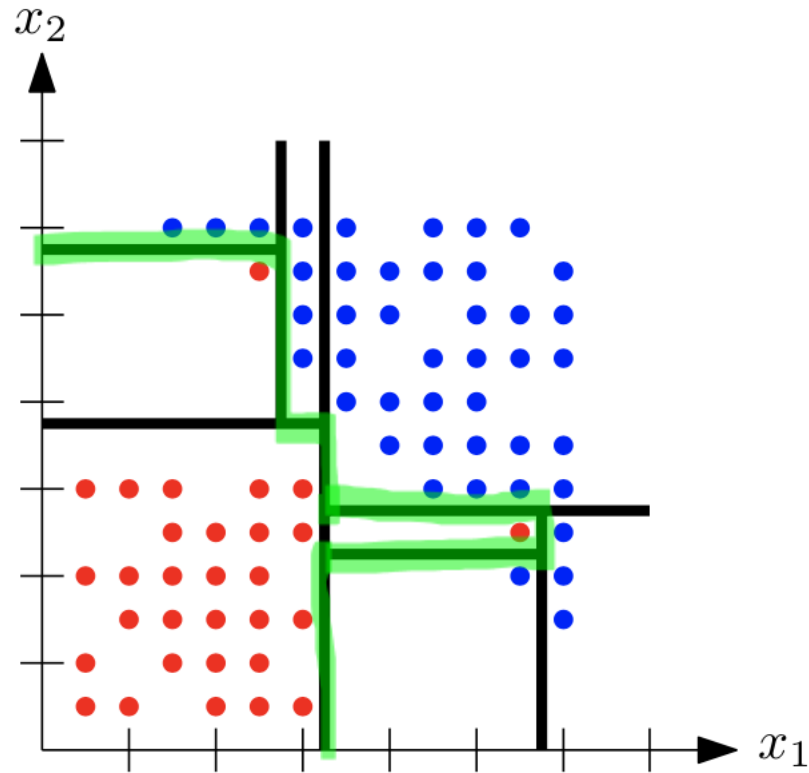
    (notations from the decision tree pseudocode)

Comparison: For classification

$$\widehat{err}(f)=\frac{1}{n}\sum_{i=1}^{n}\mathbf{I}\{f(x_i)\neq y_i\}$$
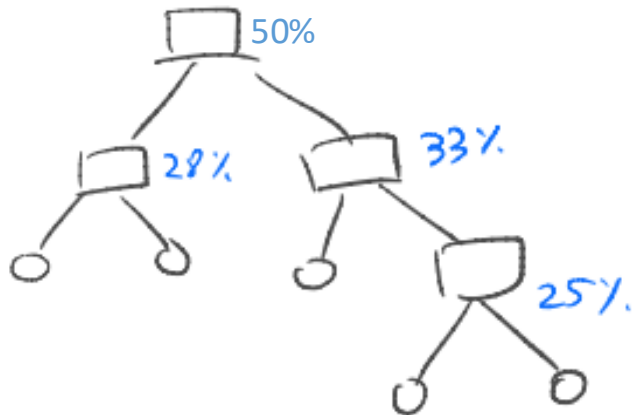
# "Spurious" patterns can be learned



note axis-parallel decision boundaries

# Unlearn spurious patterns by <u>pruning</u>

Split the given data into **train set** and **validation set**

- Build a decision tree based on the **train set**

- min_error ← compute the **validation set** error

- While true
  - For each <u>non-leaf node</u>, pretend that it is a leaf node and then compute the validation set error (but do not make it a leaf node yet)
  - current_error ← the smallest validation set error above.
  - If current_error ≥ min_error
    - Break
  - Else
    - **Prune** the one that reduces the validation set error the most
    - min_error ← current_error



original validation set error: 35%

- ► Spam dataset
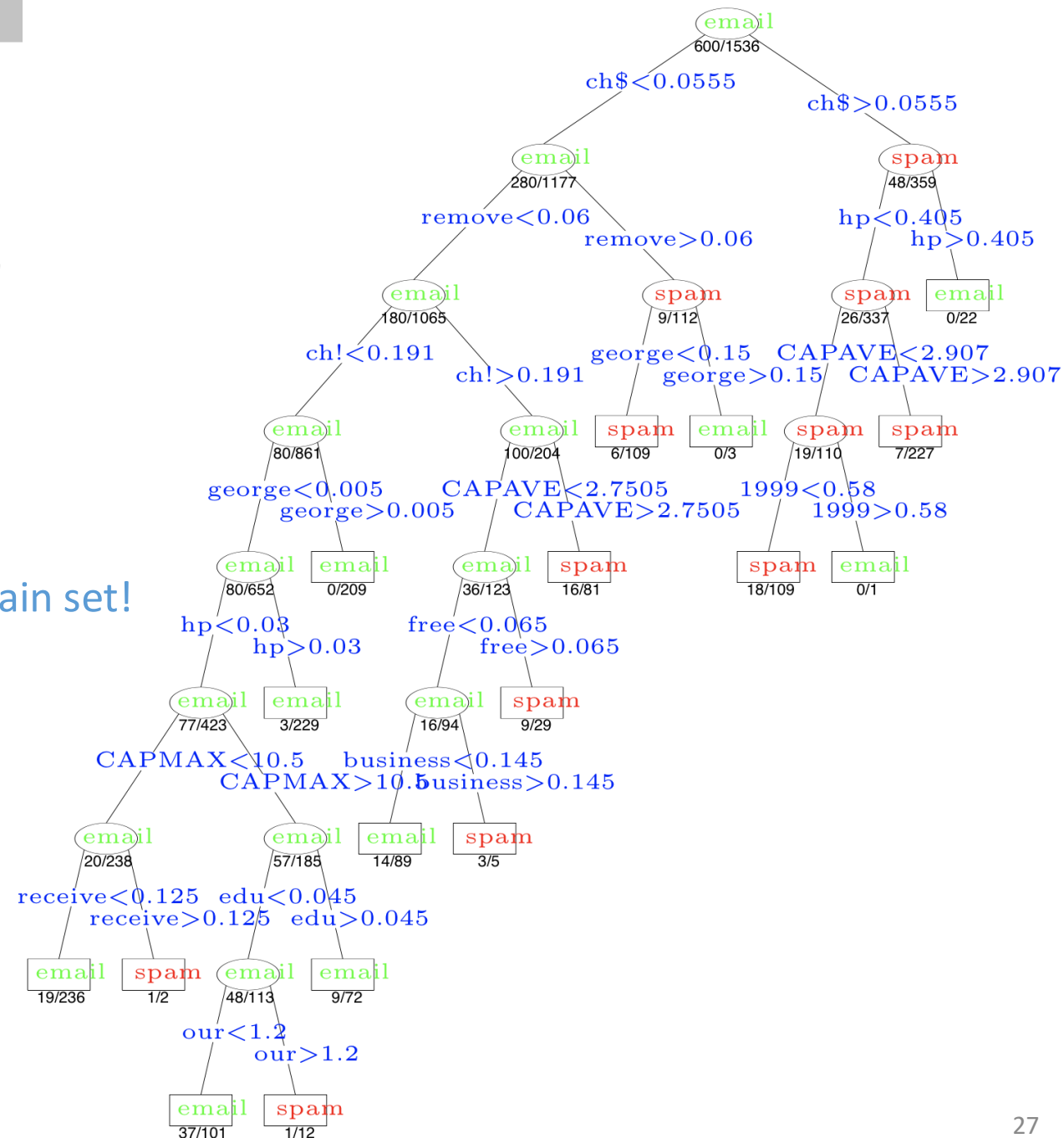- ► 4601 email messages, about 39% are spam
- ► Classify message by spam and not-spam
- ► 57 features
  - ► 48 are of the form "percentage of email words that is (WORD)"
  - ► 6 are of the form "percentage of email characters is (CHAR)"
  - ► 3 other features (e.g., "longest sequence of all-caps")
- ► Final tree after pruning has 17 leaves, 9.3% test error rate

error rate computed on test set data
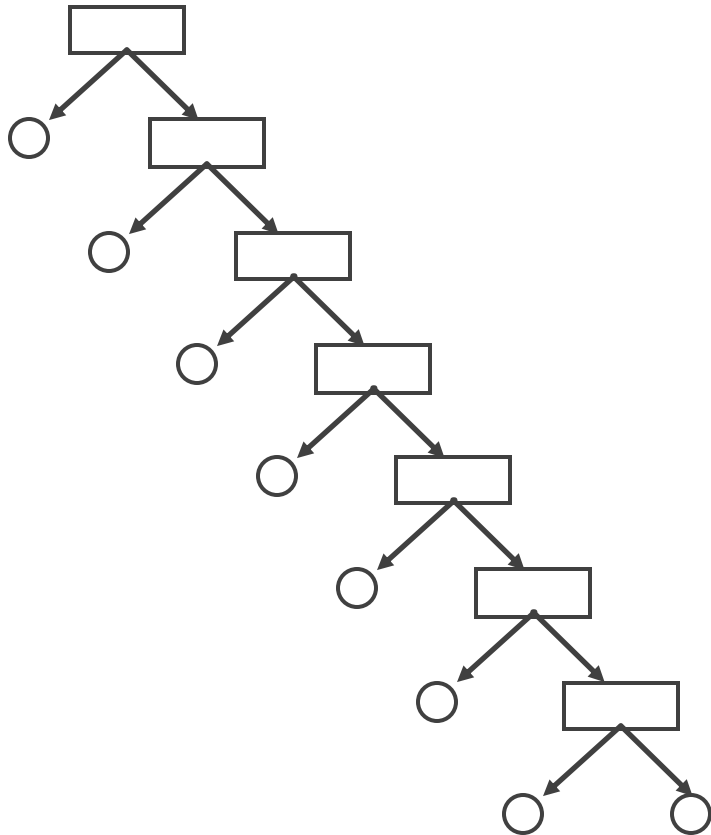⇒ test set data should not have been part of the train set!

Q: what would be the majority vote accuracy?

email 600/1536
ch$<0.0555   ch$>0.0555

email 280/1177   spam 48/359
remove<0.06   remove>0.06   hp<0.405   hp>0.405

email 180/1065   spam 9/112   spam 26/337   email 0/22
ch!<0.191   ch!>0.191   george<0.15   george>0.15   CAPAVE<2.907   CAPAVE>2.907

email 80/861   email 100/204   spam 6/109   email 0/3   spam 19/110   spam 7/227
george<0.005   george>0.005   CAPAVE<2.7505   CAPAVE>2.7505   1999<0.58   1999>0.58

email 80/652   email 0/209   email 36/123   spam 16/81   spam 18/109   email 0/1
hp<0.03   hp>0.03   free<0.065   free>0.065

email 77/423   email 3/229   email 16/94   spam 9/29
CAPMAX<10.5   CAPMAX>10.5   business<0.145   business>0.145

email 20/238   email 57/185   email 14/89   spam 3/5
receive<0.125   receive>0.125   edu<0.045   edu>0.045

email 19/236   spam 1/2   email 48/113   email 9/72
our<1.2   our>1.2

email 37/101   spam 1/12

# Time complexity

- $d$: number of binary features, $m$: the number of data points

The worst-case configuration has $O(m)$ leaf nodes    $\Rightarrow$ O(m) internal node



$\Rightarrow$ Each internal node pays O(dm) for choosing which feature

$\Rightarrow$ Total: $O(dm^2)$