

CHAPTER 9:

DECISION TREES

SECTIONS 9.1 ~ 9.2.1

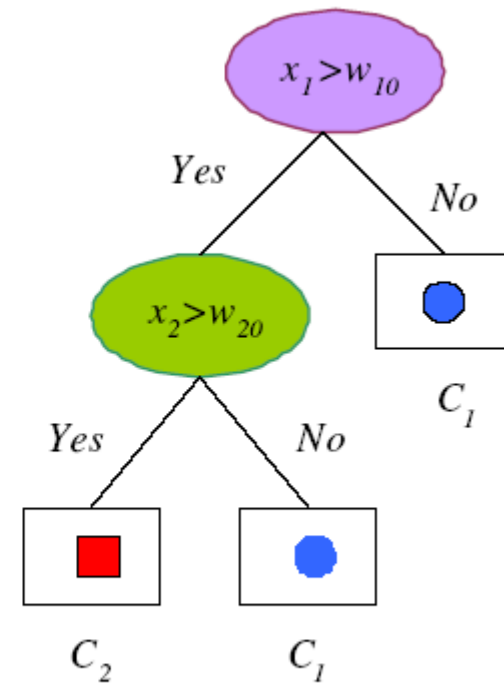
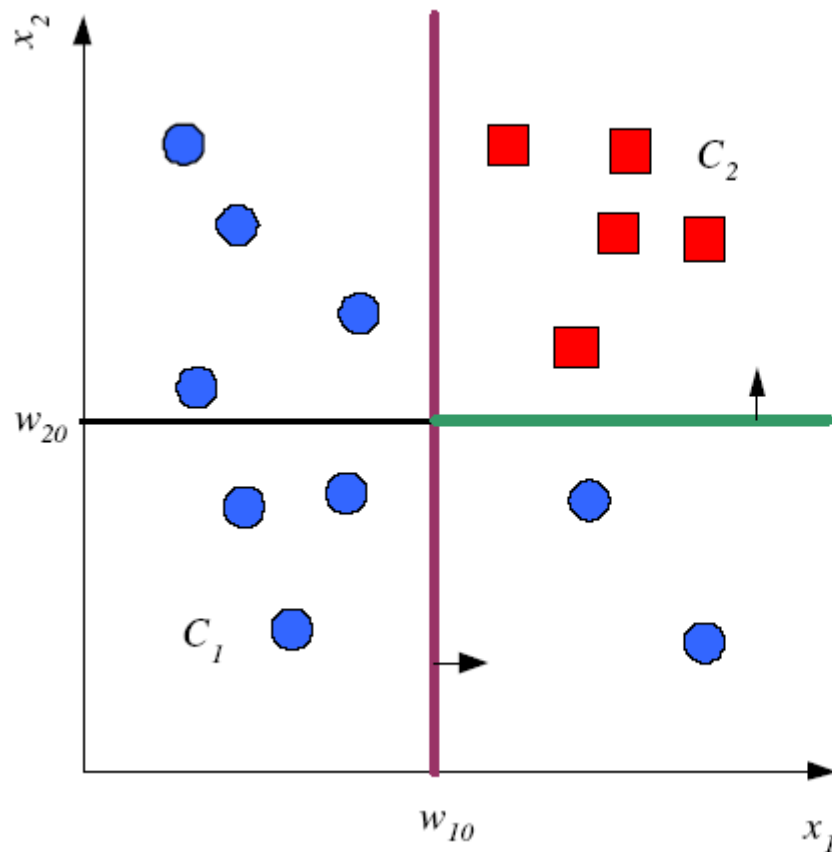
Decision tree

2

- Motivation
 - ▣ Explainable model with rules (if-then-else)
 - ▣ Usually not as accurate as other models
 - ▣ Classification, regression, etc
- How to learn a decision tree from data
 - ▣ Recursive
 - ▣ Greedily split the dataset into subgroups
 - ▣ Multiple hyperparameters and fine calibration

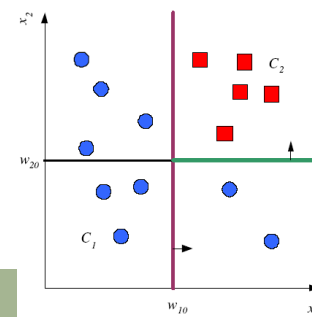
Tree Uses Nodes and Leaves

3



Two Types of Nodes

4



□ Terminal nodes (leafs)

- ▣ Provide the class (prediction)

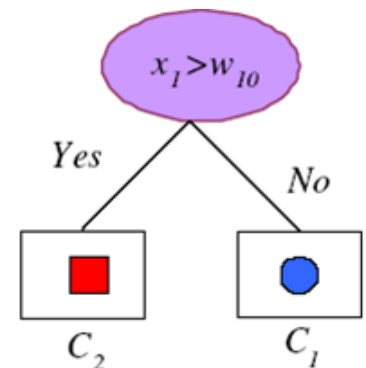
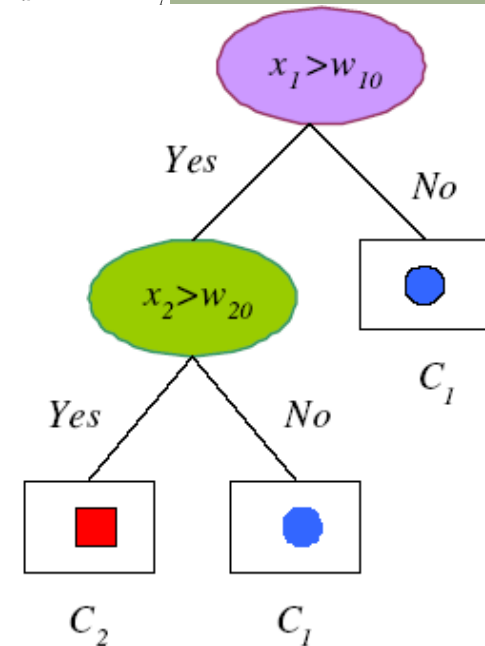
□ Internal nodes

- ▣ Test the value of one or more features
- ▣ Branch (split the data) based on the test
- ▣ Root
- ▣ Can have two or more children (not 1)
- ▣ Variables can be numeric/discrete

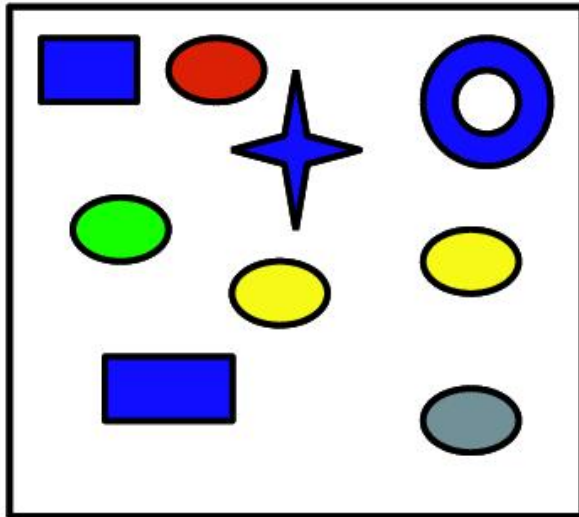
□ Terminal node does not have to be pure

□ How to learn a decision tree from data

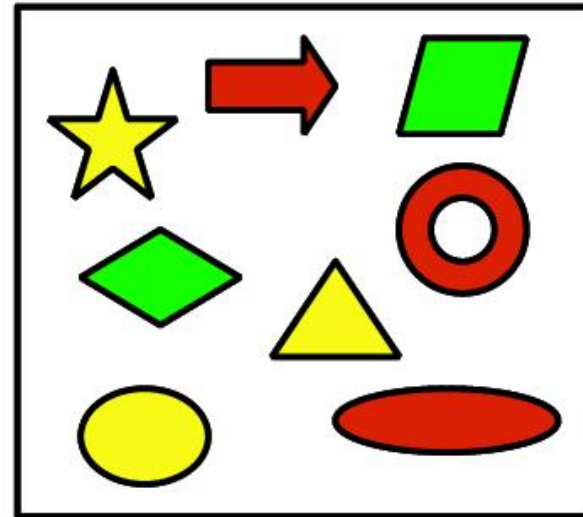
- ▣ Choose (feature, branch) for internal nodes



pos



neg



D features (attributes)



N cases

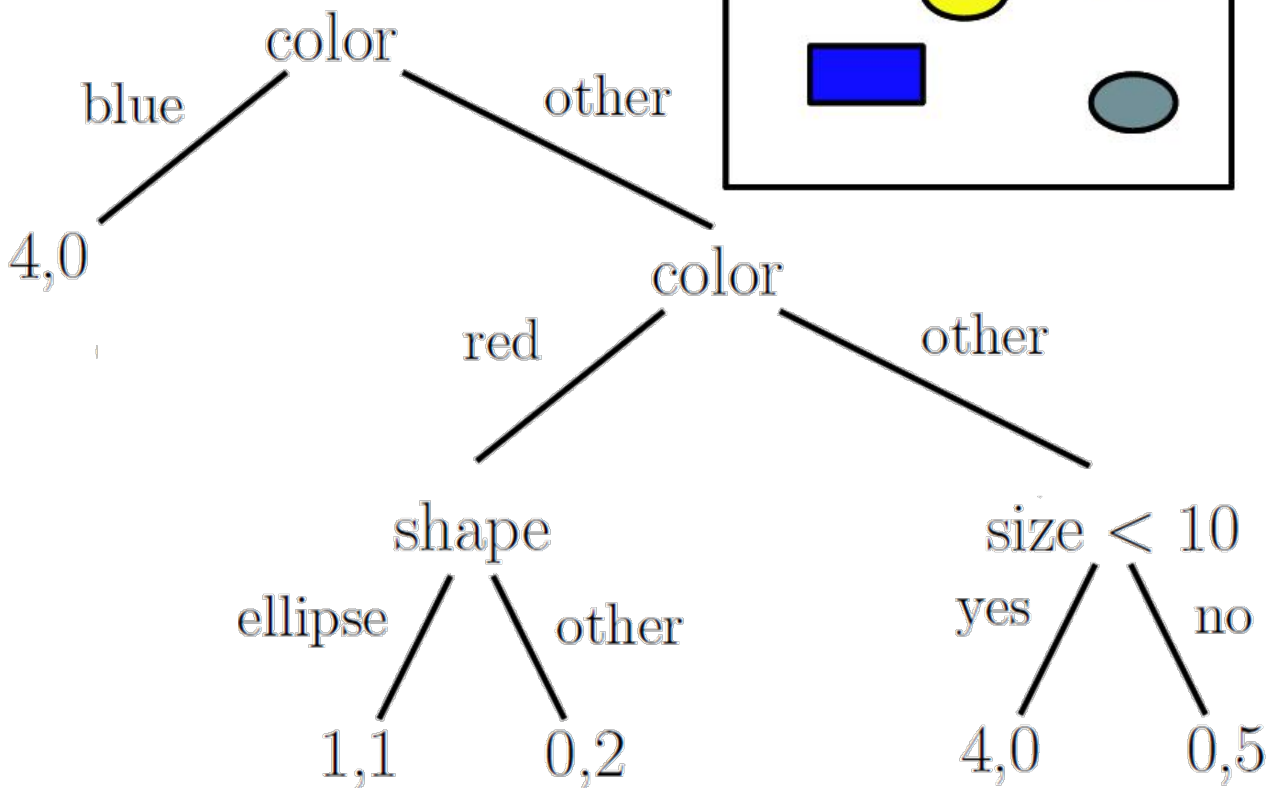
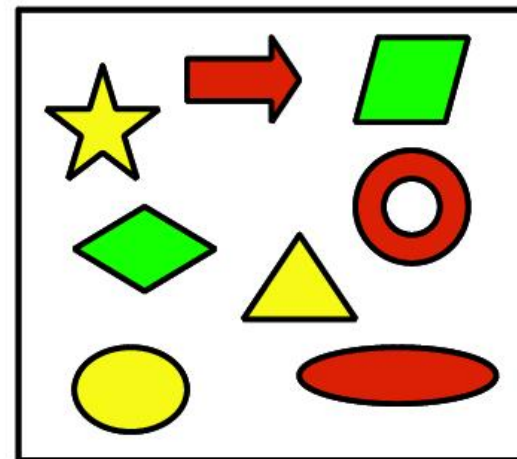
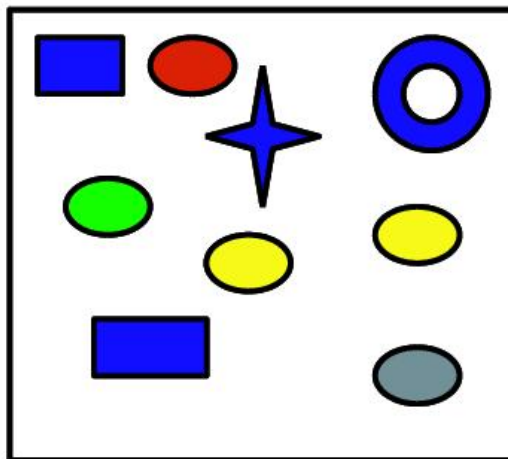
Color	Shape	Size (cm)
Blue	Square	10
Red	Ellipse	2.4
Red	Ellipse	20.7

Label
1
1
0

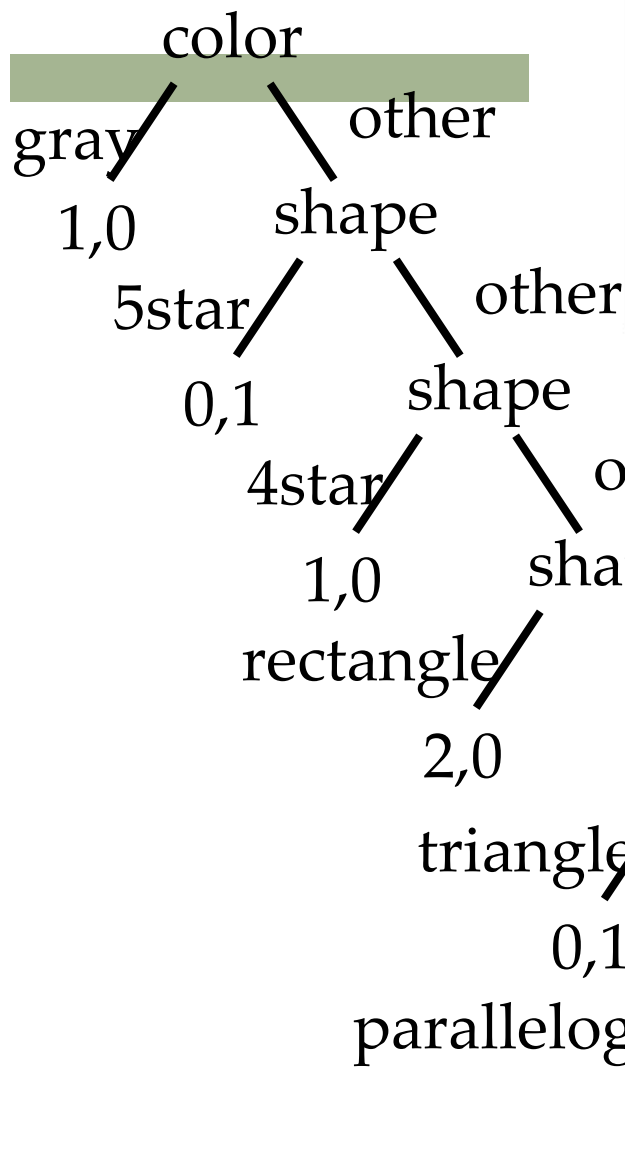
6

pos

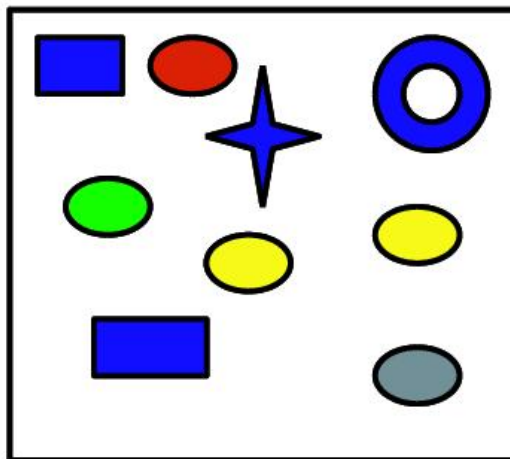
neg



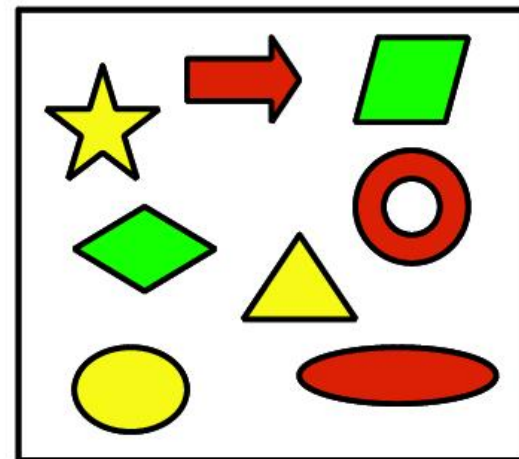
not pure



pos



neg



we can be more specific



diamond

shape

0,1

other

color

green

other

1,0

color

blue

other

1,0

shape

donut

other

0,1

shape

arrow

other

0,1

size < 2cm

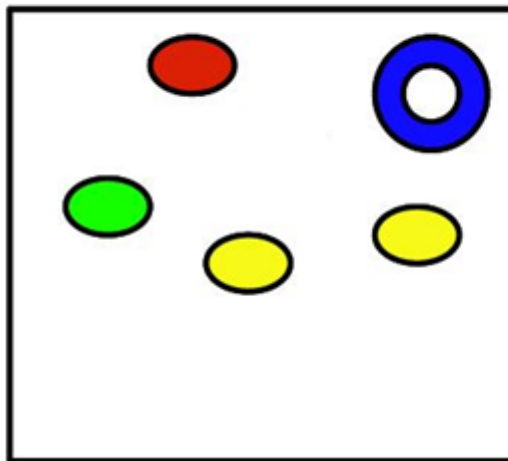
yes

no

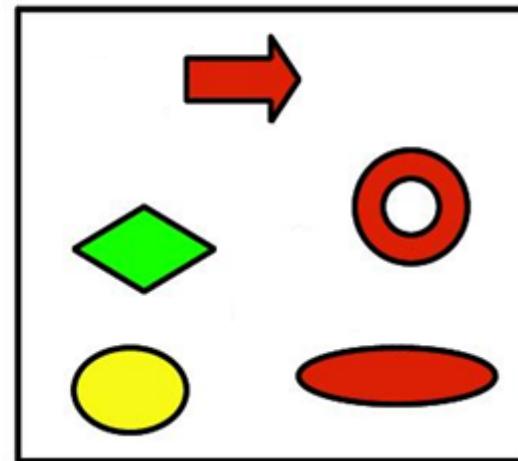
3,0

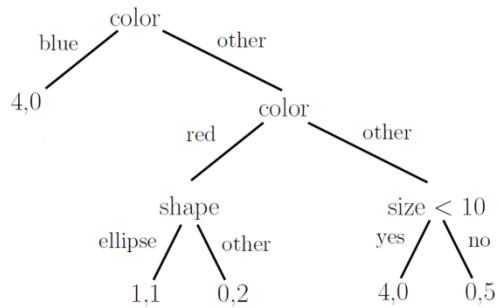
0,2

pos



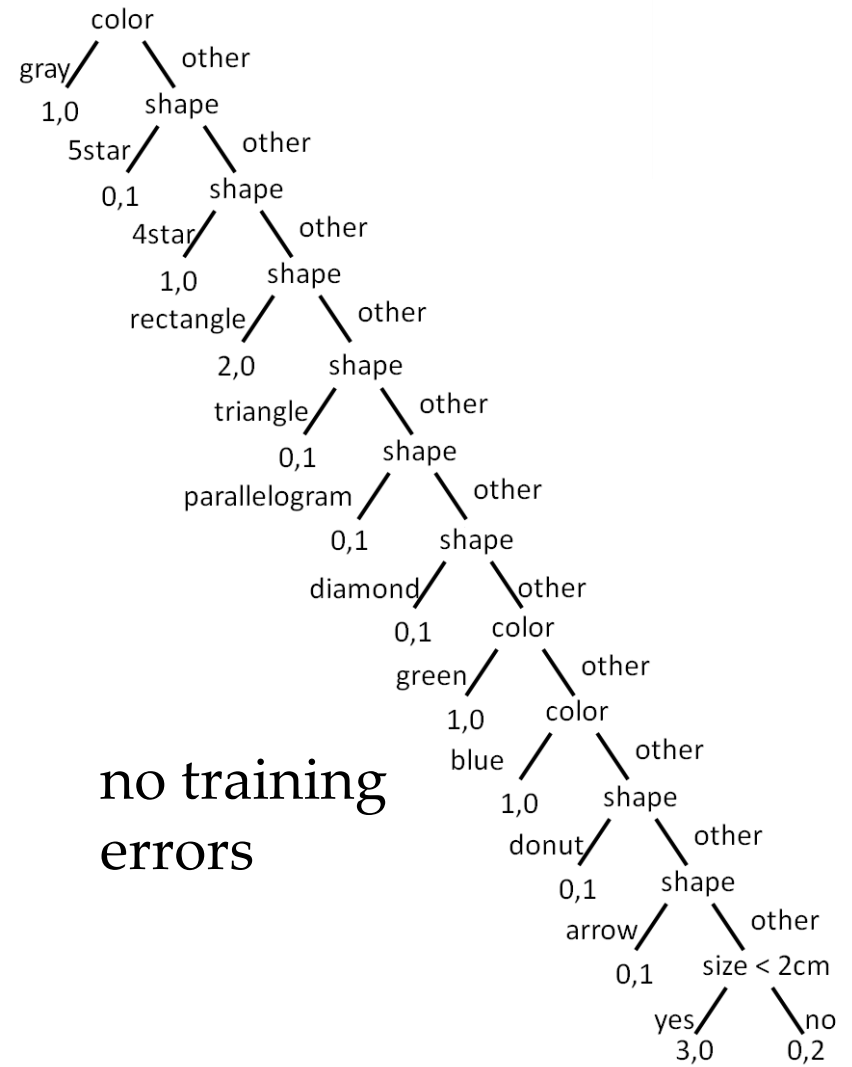
neg





1 training error

or



no training errors

Decision tree: divide and conquer

10

- Supervised learning: classification and regression
- Hierarchical, nonparametric (no a priori structure)
- Internal decision nodes implements a test function
 - ▣ Univariate: Uses a single attribute, x_i
 - Numeric x_i : Binary split : $x_i > w_m$
 - Discrete x_i : n -way split for n possible values or binary
 - ▣ Multivariate: Uses multiple/all attributes, \mathbf{x}
- Leaves
 - ▣ Classification: Class labels, or proportions
 - ▣ Regression: Numeric; r average, or local fit
- Highly interpretable (compare with Naïve Bayes)
- Learning is greedy; find the best split recursively (Breiman et al, 1984; Quinlan, 1986, 1993)

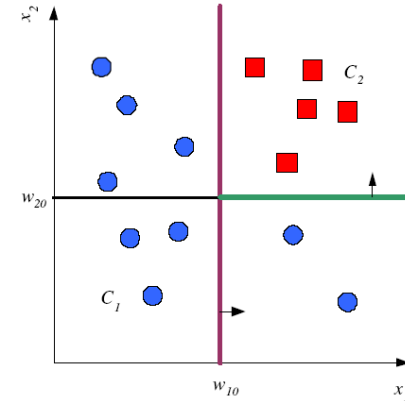
Classification Trees (ID3,CART,C4.5)

11

- Algorithms differ in branching models

- Pick a feature x_j
- Discrete case (with n values): split into n branches
- Numeric case: discretize into **two** by thresholding

$$\underline{f}_m(\mathbf{x}) : x_j > w_{m0} \text{ (threshold)}$$



- **Goal:** find the smallest tree that has low/zero training error

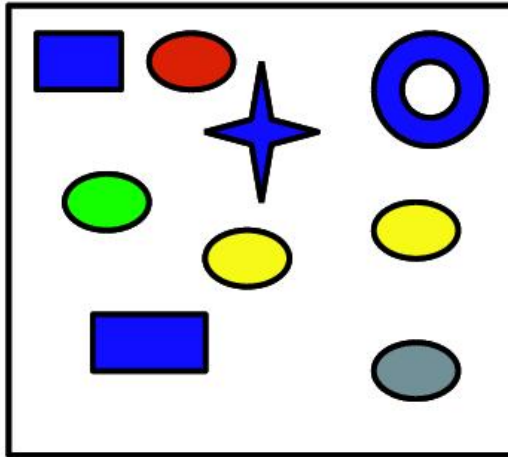
- Small means depth, number node, breath, branching factor, etc
- NP-complete, forced to use local search based on heuristic

- Greedy: each step we look for the best split

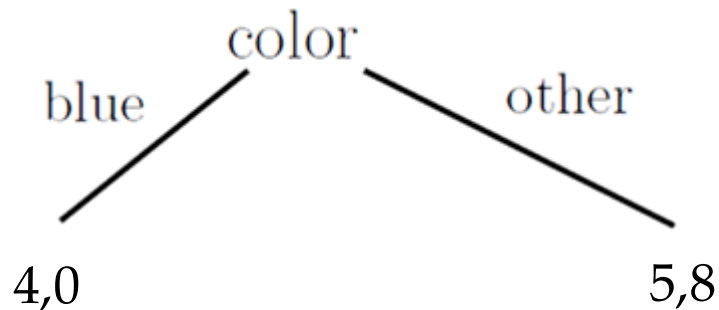
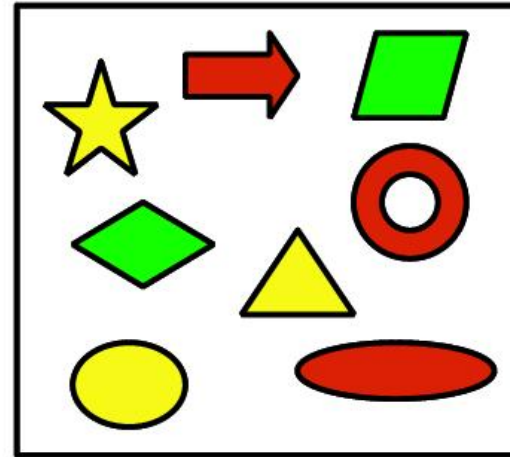
Score(D_1, D_2, \dots, D_k) measuring “goodness” of splitting the data into k subsets

- Continue recursively until no more split (leaf node)

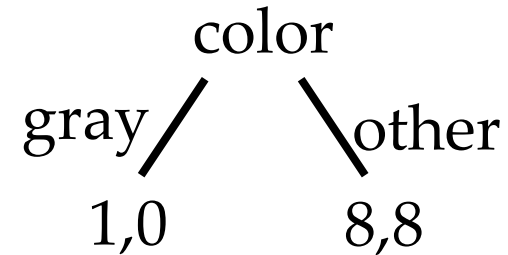
pos



neg



Error rate: 5/17

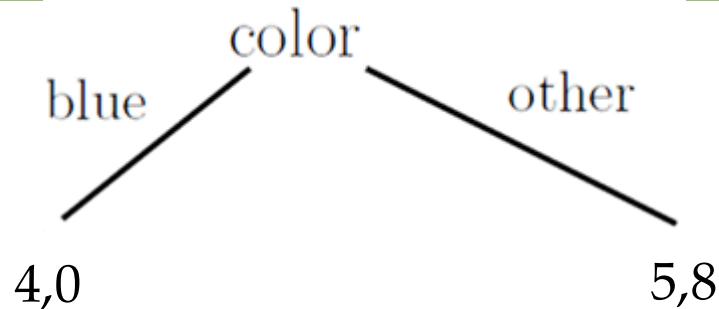


Error rate: 8/17

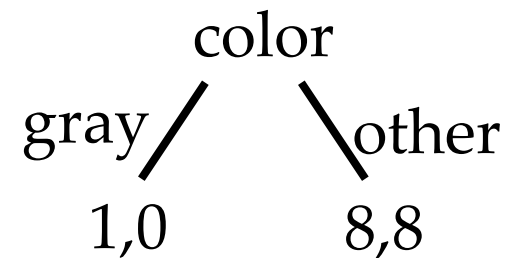
Use majority label at the leaf, then compute error rate
Also equal to the weighted average of error on both groups.

Selection of feature and splitting

13



Error rate: 5/17



Error rate: 8/17

Select the feature and splitting value that “**progresses most**” towards a lower error.

best_loss = infinity

for feat in all possible features, i.e., {shape, color, size, ...}

for v in all possible values of feat (e.g., {red, blue, green yellow} for color)

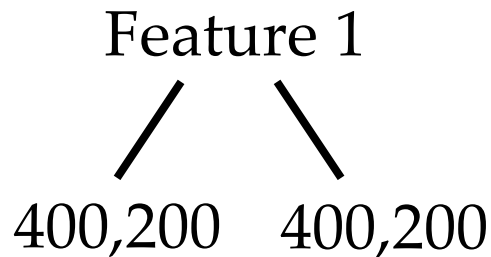
 score = **error of splitting the current dataset along (feat, v)**

if score < best_loss

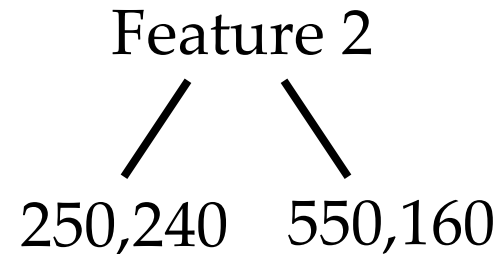
 best_loss = score, and record feat and v

Accuracy score pitfall

14



Error rate: $(200+200)/1200$



Error rate: $(240+160)/1200$

Both have the same error rate!

Which one is “progressing more” towards a better solution?

```
best_loss = infinity
for feat in all possible features, i.e., {shape, color, size, ...}
    for v in all possible values of feat (e.g., {red, blue, green yellow} for color)
        score = some smart cost of splitting the current dataset along (feat, v)
        if score < best_loss
            best_loss = score, and record feat and v
```

Best split in classification

15

- For node m , N_m instances reach m ,
 N_m^i belong to C_i , then $\hat{p}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m}$
- Node m is pure if $p_m^i = 1$ for a certain i
- If node m is pure, generate a leaf and stop, otherwise split and continue recursively
- Measure of impurity is entropy:

$$I_m = - \sum_{i=1}^K p_m^i \log_2 p_m^i$$

Entropy as an impurity measure

16

Measure of (degree of) uncertainty

The more clueless I am about the answer initially, the more information is contained in the answer

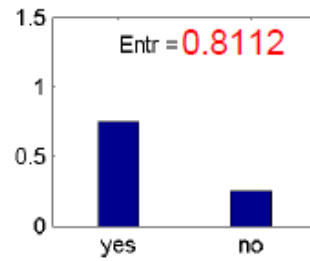
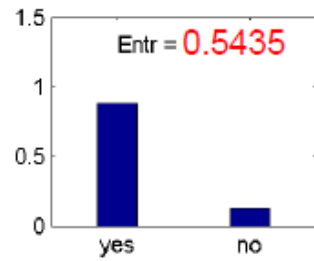
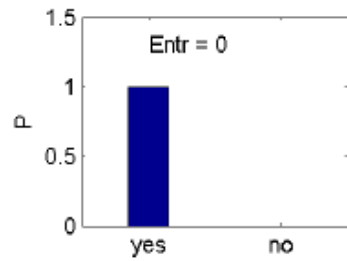
Information in an answer when prior is (P_1, \dots, P_n)

$$\sum_{i=1}^n -P_i \log_2 P_i$$

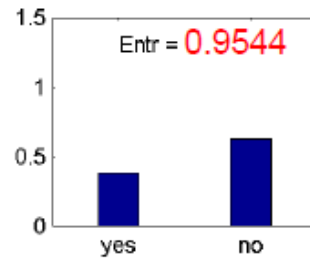
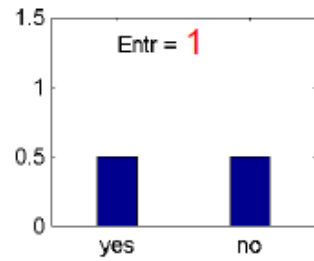
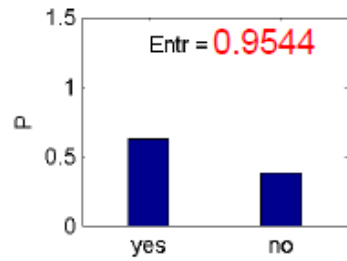
Scale: 1 bit = answer to Boolean question with equal prior

Roll of a 4-sided die has 2 bits of information.

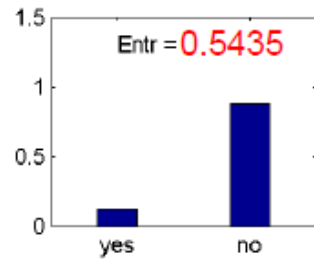
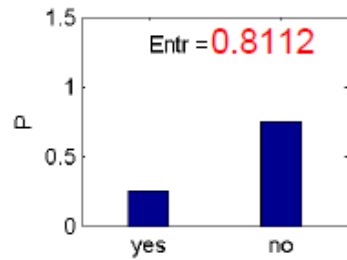
Acquisition of information leads to reduction in entropy



The entropy is maximal when all possibilities are equally likely.

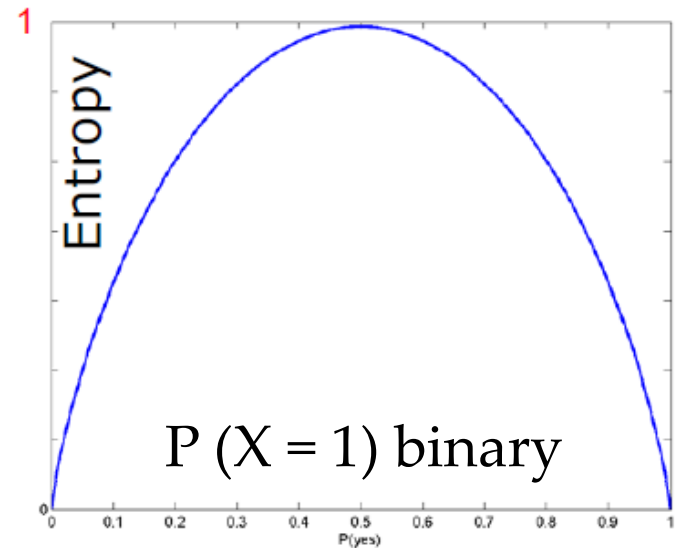


The goal of the decision tree is to decrease the entropy in each node.



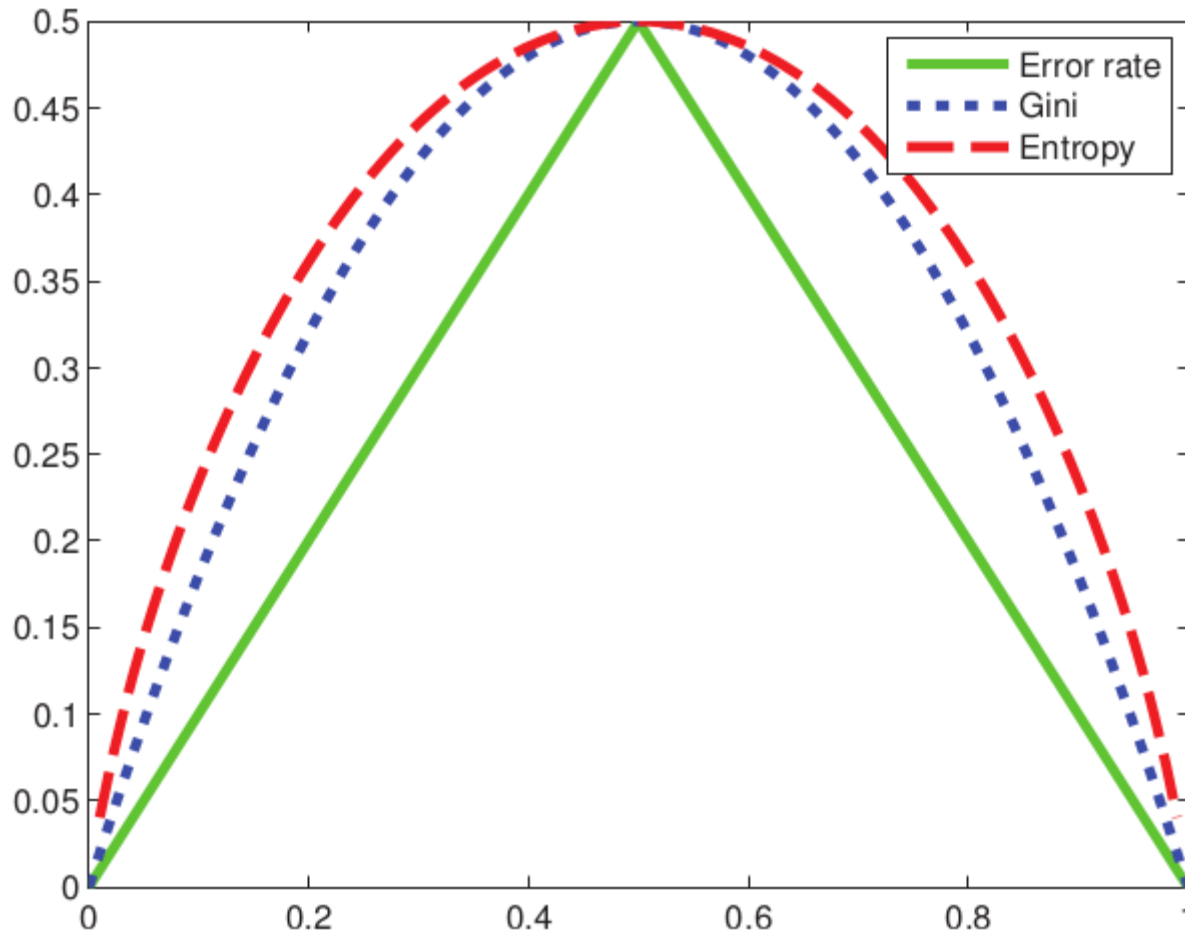
Entropy is zero in a pure "yes" node (or pure "no" node).

Entropy is a measure of "uncertainty" of a random variable.



Impurity measures

18



P (X=1) binary

Entropy:

$$\sum_{i=1}^n -P_i \log_2 P_i$$

= (in the binary case)

$$-p \log p - (1-p) \log p$$

The curve for entropy is rescaled (halved) on the left.

Gini index:

$$1 - \sum_{i=1}^n P_i^2$$

Error rate:

$$1 - \max_{i=1 \dots n} P_i$$

Best split in classification

19

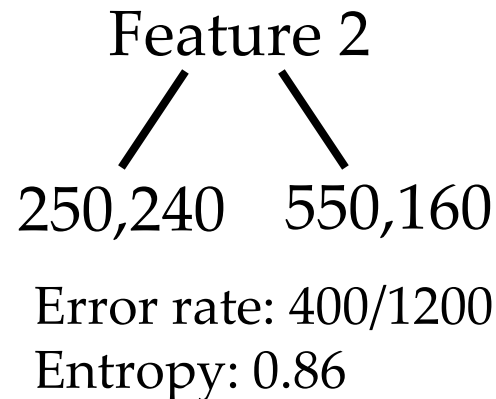
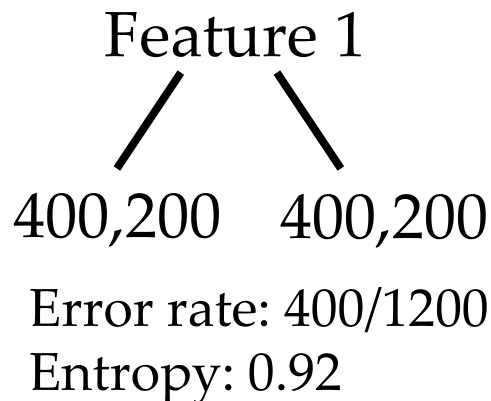
- If node m is pure, generate a leaf and stop, otherwise split and continue recursively

- Impurity after split: N_{mj} of N_m take branch j . N_{mj}^i belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$

$$J'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

- Find the variable and split that best reduces impurity (among all variables -- and split positions for numeric variables)



Best split in classification

20

- If node m is pure, generate a leaf and stop, otherwise split and continue recursively
- Impurity after split: N_{mj} of N_m take branch j . N_{mj}^i belong to C_i
$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}} \quad \mathcal{I}'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$
- Find the variable and split that best reduces impurity (among all variables -- and split positions for numeric variables)

best_loss = infinity

for feat in all possible features, i.e., {shape, color, size, ...}

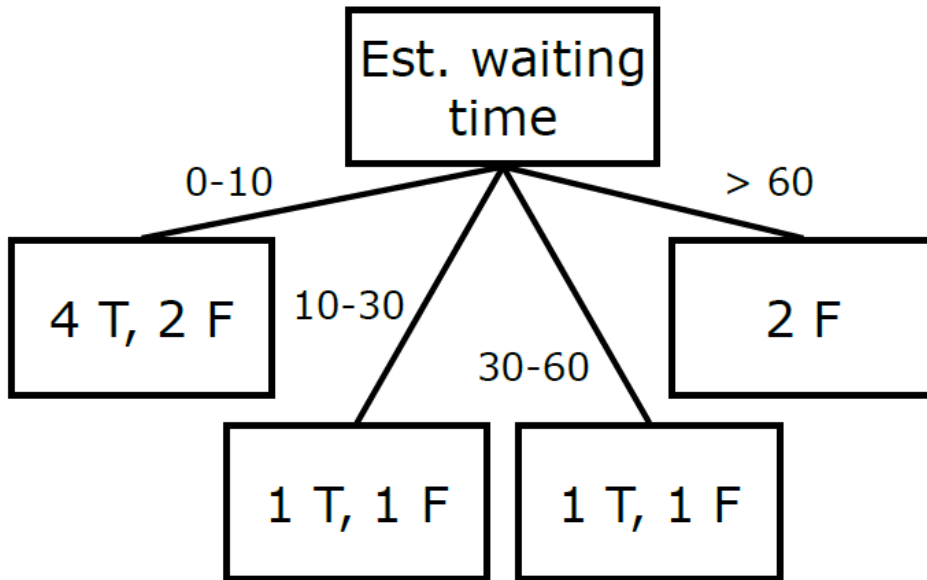
for \mathbf{v} in all possible values of feat (e.g., {red, blue, green yellow} for color)

 score = **impurity after splitting along (feat, v)**

if score < best_loss

 best_loss = score, and record feat and \mathbf{v}

Information Gain Example



Example	Attributes										Target
	Alt	Bar	Fri	Hun	Pat	Price	Rain	Res	Type	Est	WillWait
X ₁	T	F	F	T	Some	\$\$\$	F	T	French	0-10	T
X ₂	T	F	F	T	Full	\$	F	F	Thai	30-60	F
X ₃	F	T	F	F	Some	\$	F	F	Burger	0-10	T
X ₄	T	F	T	T	Full	\$	F	F	Thai	10-30	T
X ₅	T	F	T	F	Full	\$\$\$	F	T	French	>60	F
X ₆	F	T	F	T	Some	\$\$	T	T	Italian	0-10	T
X ₇	F	T	F	F	None	\$	T	F	Burger	0-10	F
X ₈	F	F	F	T	Some	\$\$	T	T	Thai	0-10	T
X ₉	F	T	T	F	Full	\$	T	F	Burger	>60	F
X ₁₀	T	T	T	T	Full	\$\$\$	F	T	Italian	10-30	F
X ₁₁	F	F	F	F	None	\$	F	F	Thai	0-10	F
X ₁₂	T	T	T	T	Full	\$	F	F	Burger	30-60	T

**More details at
Decision_tree_
example.pdf in
the slides
section**

$$\text{Remainder}(\text{Wait}) = \frac{6}{12} \left[-\left(\frac{4}{6}\right) \log_2 \left(\frac{4}{6}\right) - \left(\frac{2}{6}\right) \log_2 \left(\frac{2}{6}\right) \right]$$

$$+ \frac{2}{12} \left[-\left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) \right]$$

$$+ \frac{2}{12} \left[-\left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) - \left(\frac{1}{2}\right) \log_2 \left(\frac{1}{2}\right) \right] + \frac{2}{12} \left[-\left(\frac{0}{2}\right) \log_2 \left(\frac{0}{2}\right) - \left(\frac{2}{2}\right) \log_2 \left(\frac{2}{2}\right) \right] = 0.7925$$

$$\text{Gain}(\text{Wait}) = 1 - 0.7925 = 0.2075$$

GenerateTree(\mathcal{X})

If NodeEntropy(\mathcal{X}) $< \theta_I$ /* eq. 9.3

Create leaf labelled by majority class in \mathcal{X}

Return

$i \leftarrow \text{SplitAttribute}(\mathcal{X})$

For each branch of \mathbf{x}_i

Find \mathcal{X}_i falling in branch

GenerateTree(\mathcal{X}_i)

SplitAttribute(\mathcal{X})

MinEnt \leftarrow MAX

For all attributes $i = 1, \dots, d$

If \mathbf{x}_i is discrete with n values

Split \mathcal{X} into $\mathcal{X}_1, \dots, \mathcal{X}_n$ by \mathbf{x}_i

$e \leftarrow \text{SplitEntropy}(\mathcal{X}_1, \dots, \mathcal{X}_n)$

$$I'_m = - \sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

If $e < \text{MinEnt}$ MinEnt $\leftarrow e$; bestf $\leftarrow i$

Else /* \mathbf{x}_i is numeric */

For all possible splits

Split \mathcal{X} into $\mathcal{X}_1, \mathcal{X}_2$ on \mathbf{x}_i

$e \leftarrow \text{SplitEntropy}(\mathcal{X}_1, \mathcal{X}_2)$

If $e < \text{MinEnt}$ MinEnt $\leftarrow e$; bestf $\leftarrow i$

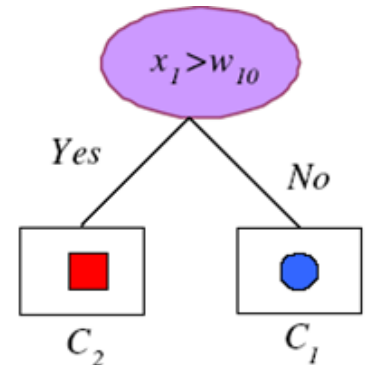
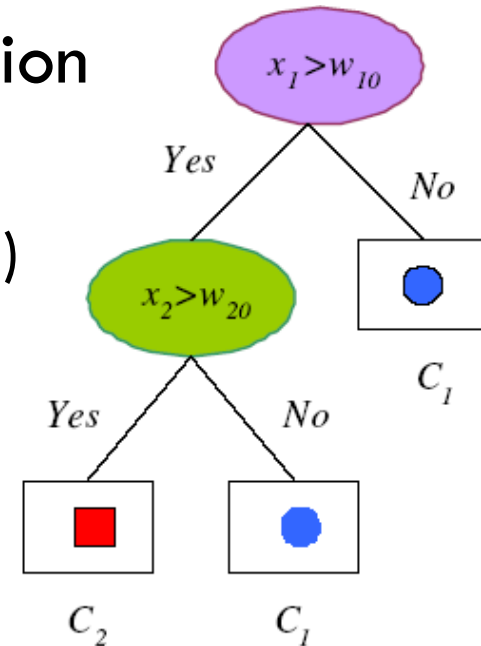
Return bestf

computational
cost?

Pruning Trees

23

- Remove subtrees for better generalization (decrease variance)
 - ▣ Prepruning: Early stopping (e.g. $< 5\%$ points)
 - ▣ Postpruning: Grow the whole tree then prune subtrees that overfit on the pruning set
 - ▣ Set aside a subset of data for pruning
- Prepruning is faster, postpruning is more accurate (requires a separate pruning set)



Decision tree

24

- Motivation
 - ▣ Explainable model with rules (if-then-else)
 - ▣ Usually not as accurate as other models
 - ▣ Classification, regression, etc
- How to learn a decision tree from data
 - ▣ Recursive
 - ▣ Greedily split the dataset into subgroups
 - ▣ Various impurity measures
 - ▣ Pruning