

ensemble
- bagging
- boosting

CS 412

APR 2ND – DECISION TREES

Final Project

Official write up/description

Groups of up to 3

- Groups of 2, I expect ~50% more
- Groups of 3, ~100% more

Every group will need to meet digitally with me (~15 minutes) between 4/13 and 4/17 to make sure the project is feasible

Project will be 5-10 pages, single-spaced, but with figures

- Problem introduction
- Previous work
- Novel attempt
- Results
- Future work/Conclusion

Team Picked
Project Topic
Next Friday
4/16

Due: May 8th

Rationale for Ensemble Learning

No Free Lunch thm: There is no algorithm that is always the most accurate

Generate a group of base-learners which when combined have higher accuracy

Different learners use different

- Algorithms
- Parameters
- Representations (Modalities)
- Training sets
- Subproblems

(combining)
"weak"
base-learners

bagging: lowers
variance
boosting: lowers
bias

Boosting

Boosting works differently.

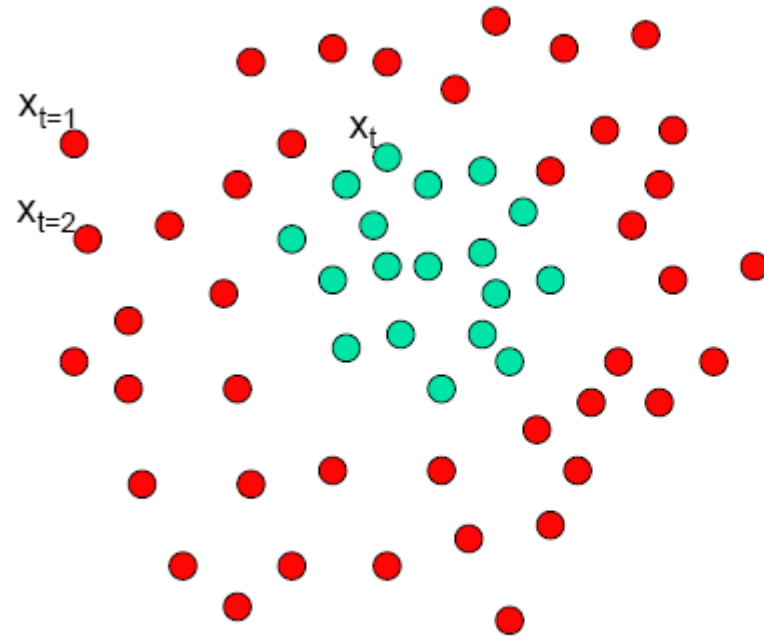
1. Boosting does not involve bootstrap sampling
2. Decision models are grown sequentially: each model is created using information from previously grown trees
- 3. Like bagging, boosting involves combining a large number of models, f^1, \dots, f^B

changing weights of incorrectly classified
points

assign a weight to the model h_t

Boosting - Example

Weak classifier
- linear discriminator



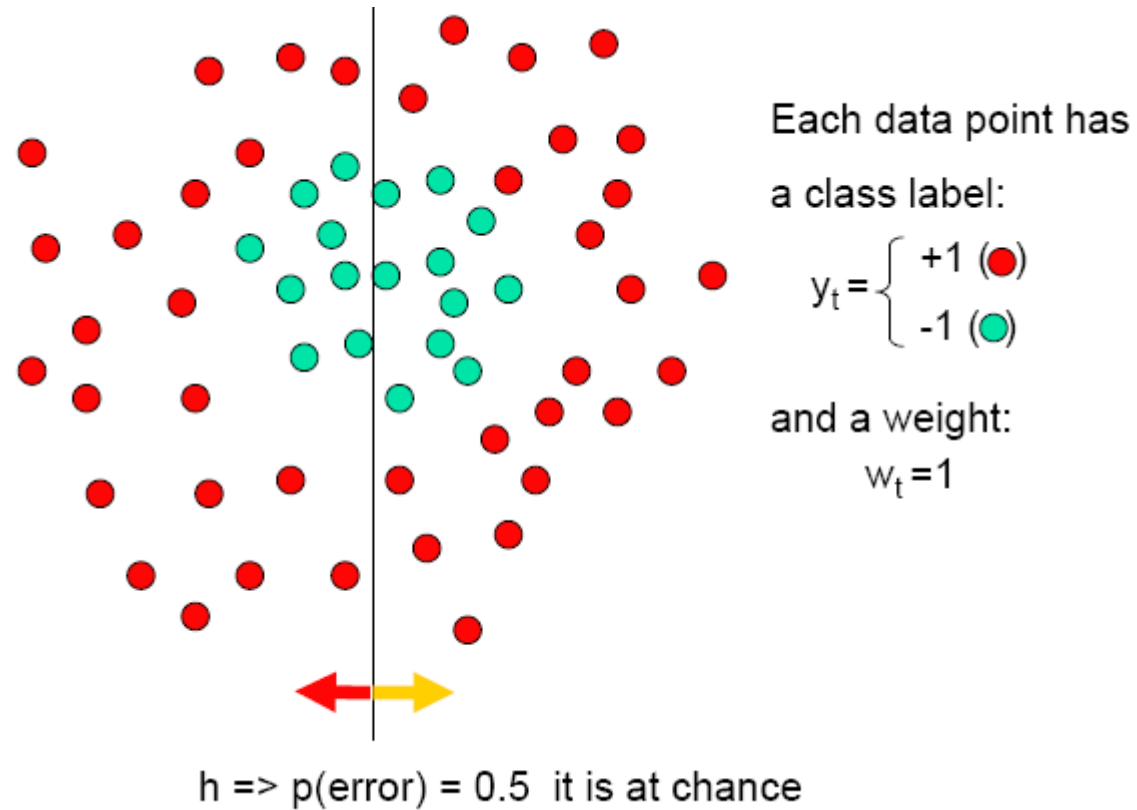
Each data point has
a class label:

$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{cyan circle}) \end{cases}$$

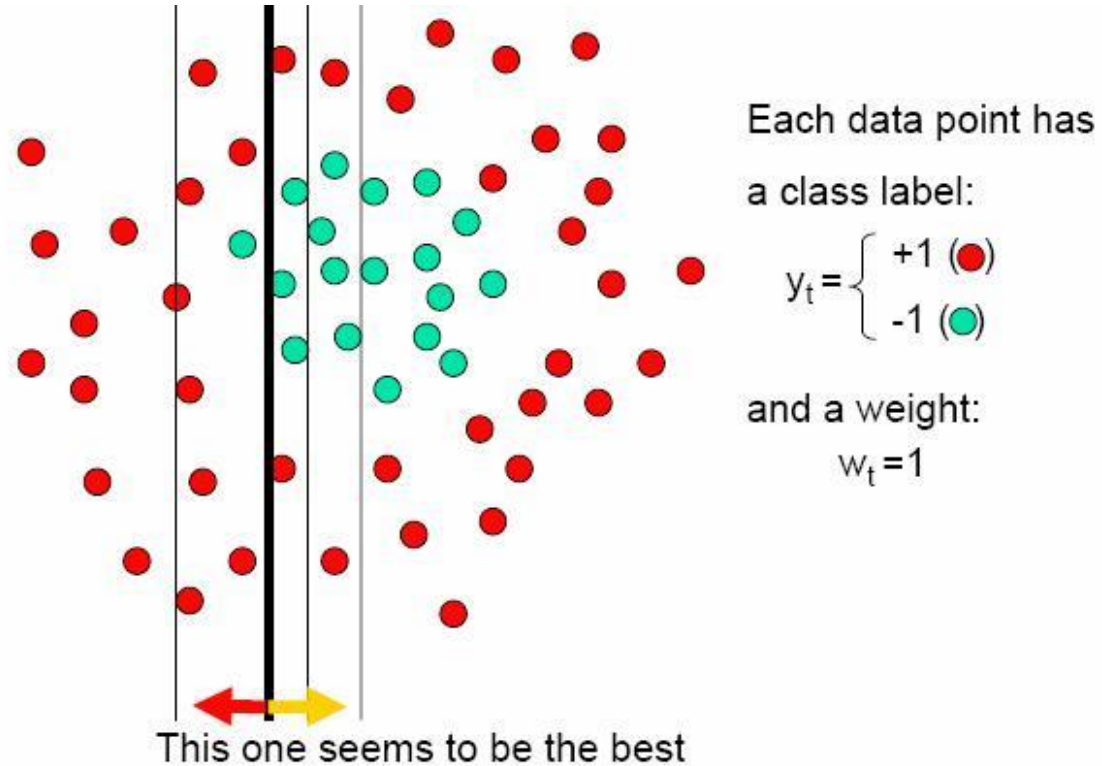
and a weight:

$$w_t = 1$$

Boosting - Example



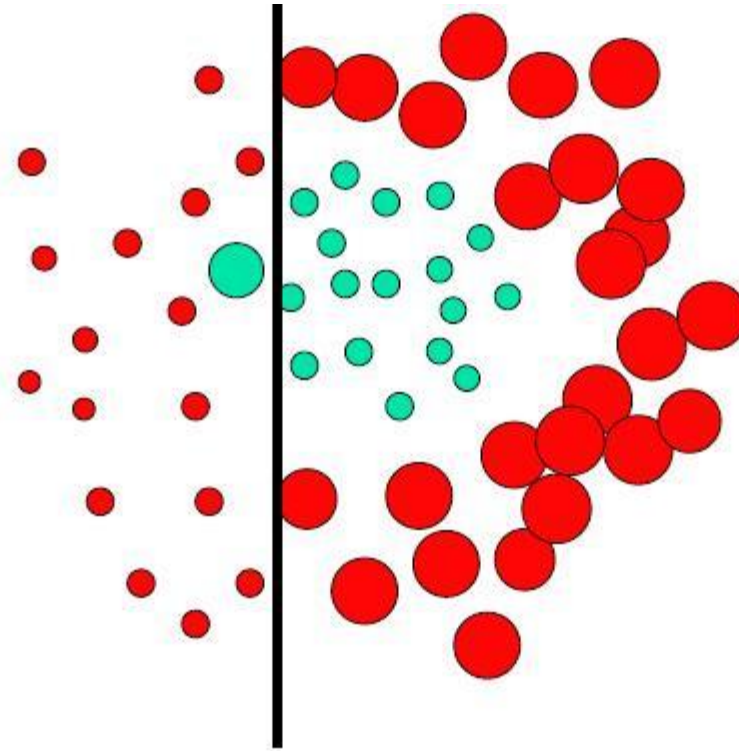
Boosting - Example



This is a '**weak classifier**': It performs slightly better than chance.

Boosting - Example

$$\sum_i D_i = 1$$



Each data point has
a class label:

$$y_t = \begin{cases} +1 & (\text{red}) \\ -1 & (\text{cyan}) \end{cases}$$

We update the weights:

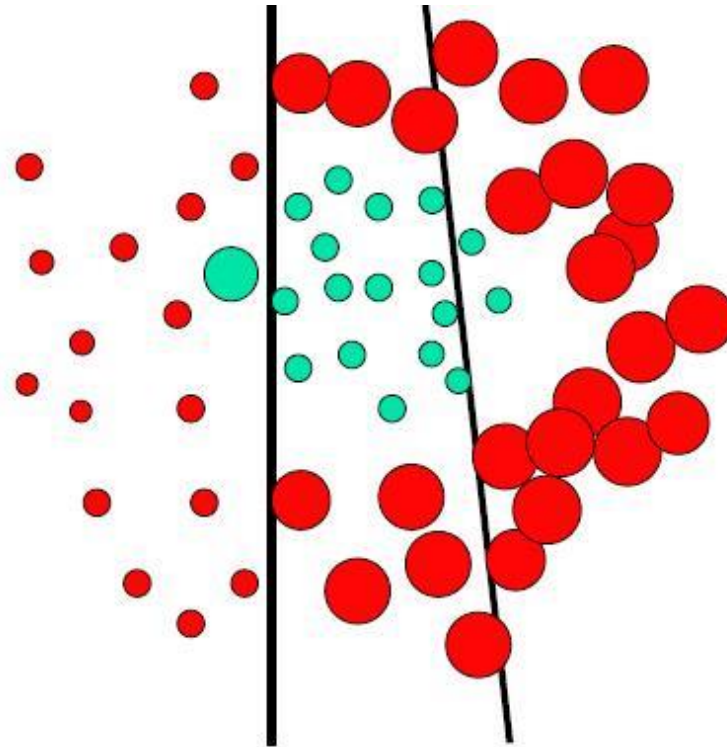
$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

reweight then $\sum D_i = 1$

$$\epsilon_t = \sum_{\text{incorrect}} \frac{1}{n}$$

We set a new problem for which the previous weak classifier performs at chance again

Boosting - Example



Each data point has
a class label:

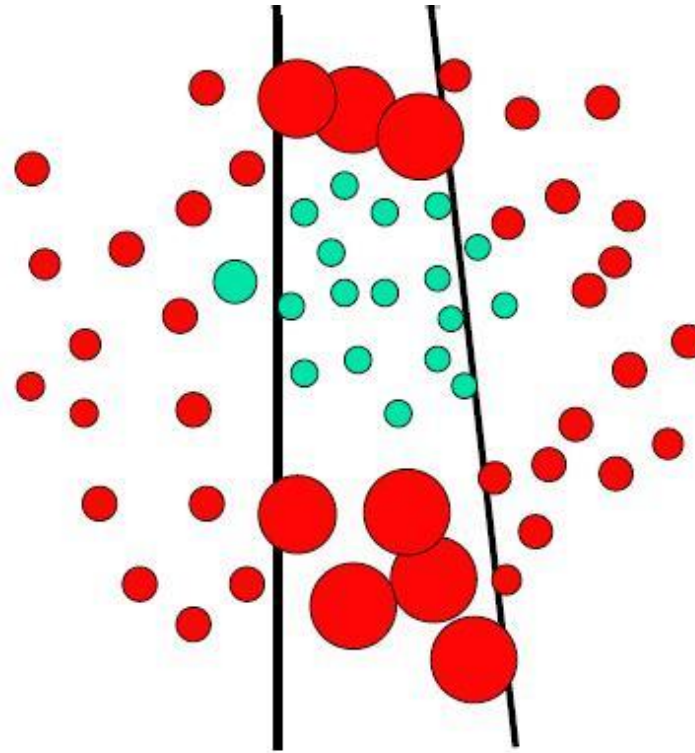
$$y_t = \begin{cases} +1 & (\text{red circle}) \\ -1 & (\text{cyan circle}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

We set a new problem for which the previous weak classifier performs at chance again

Boosting - Example



Each data point has
a class label:

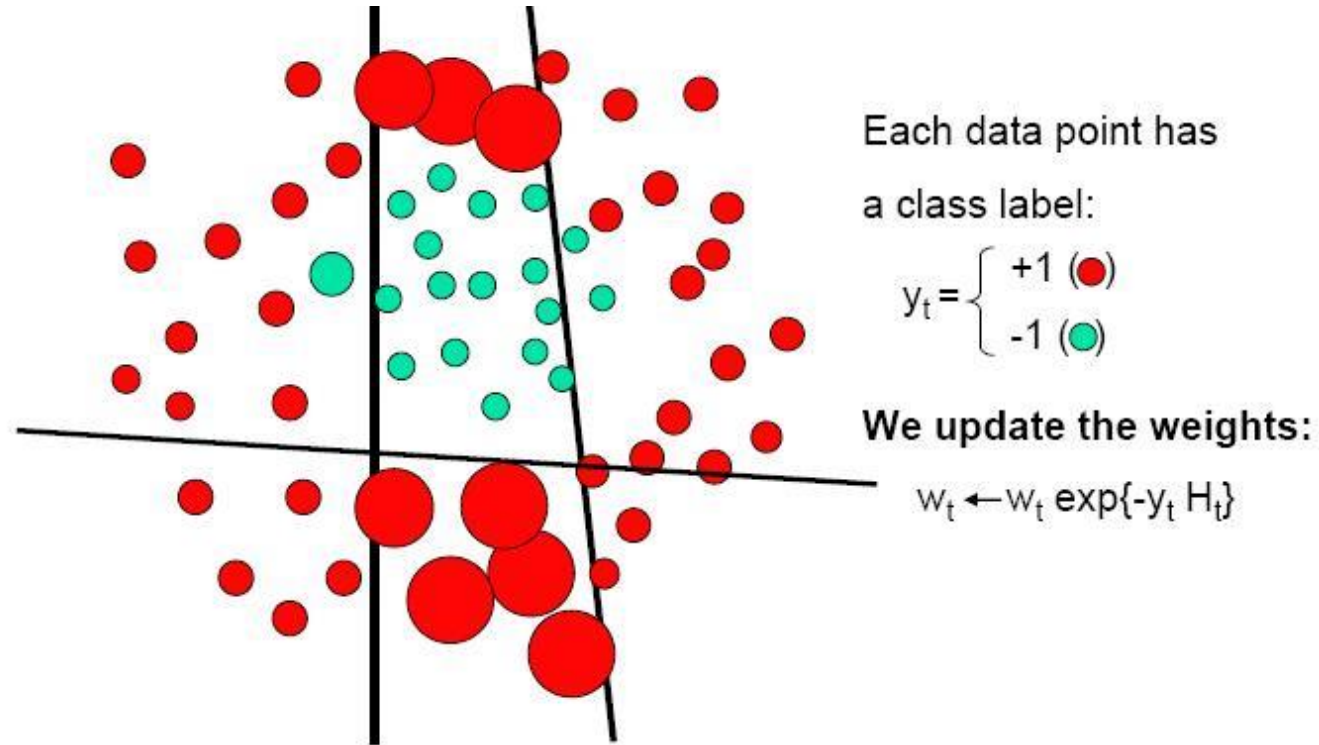
$$y_t = \begin{cases} +1 & (\text{red}) \\ -1 & (\text{cyan}) \end{cases}$$

We update the weights:

$$w_t \leftarrow w_t \exp\{-y_t H_t\}$$

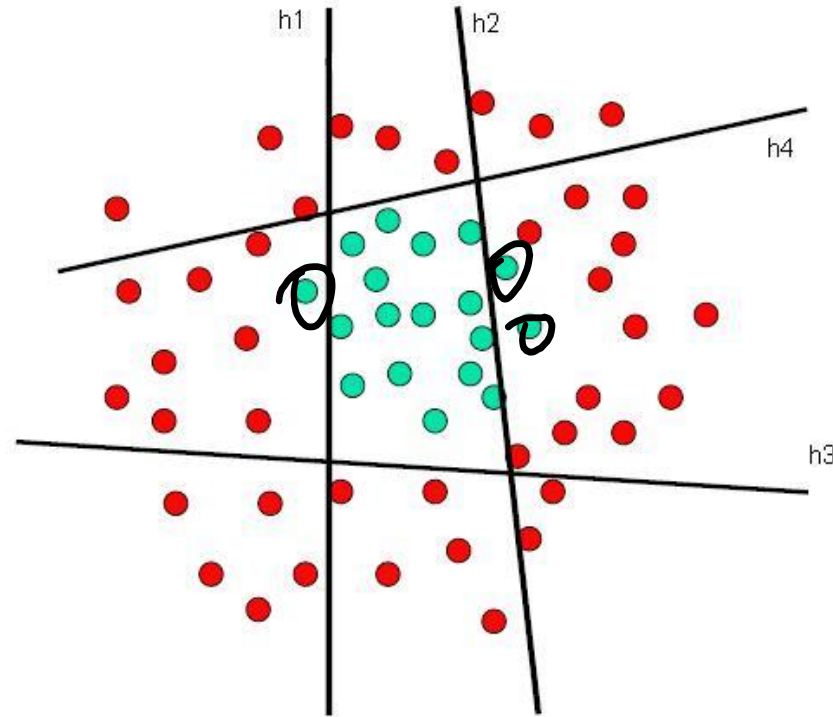
We set a new problem for which the previous weak classifier performs at chance again

Boosting - Example



We set a new problem for which the previous weak classifier performs at chance again

Boosting - Example



The strong (non- linear) classifier is built as the combination of all the weak (linear) classifiers.

iterations
given as a tuning
parameter
↑
Pick optimum
value by
cross-validation

Boosting

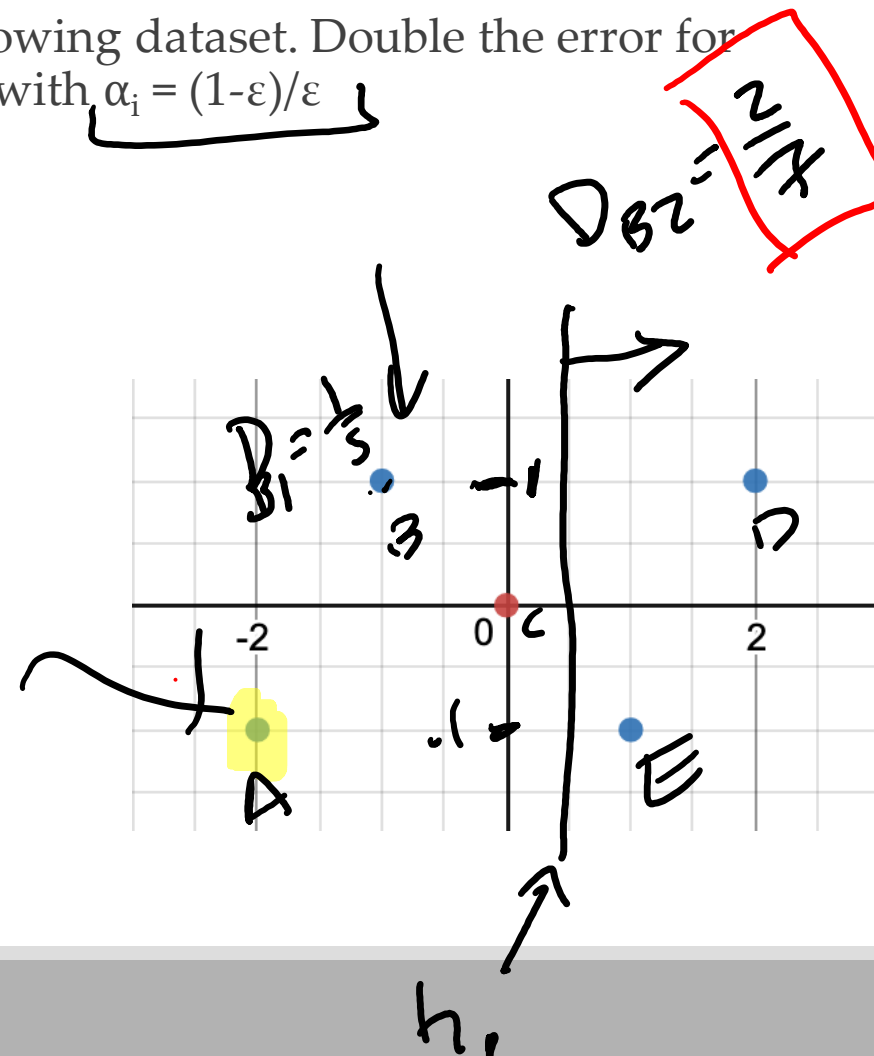
blue
positive

$$\alpha_1 = \frac{1-\epsilon_1}{\epsilon_1}$$

Perform 3 iterations of the boosting algorithm on the following dataset. Double the error for incorrect points at each iteration and weight each model with $\alpha_i = (1-\epsilon_i)/\epsilon_i$

$$\epsilon_1 = \sum_{i \in \text{incorrectly classified}} D_i = \frac{2}{5}$$

$$\alpha_1 = \frac{\frac{3}{5}}{\frac{2}{5}} = \frac{3}{2}$$



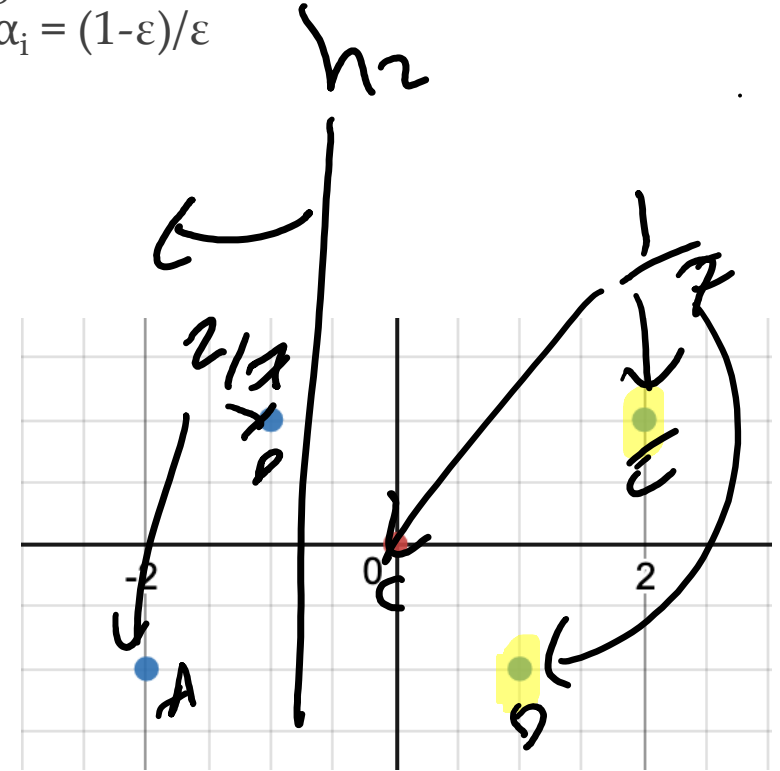
Boosting

Perform 3 iterations of the boosting algorithm on the following dataset. Double the error for incorrect points at each iteration and weight each model with $\alpha_i = (1-\epsilon)/\epsilon$

$$\epsilon_2 = \sum D_i = \frac{1}{7} + \frac{1}{7} = \frac{2}{7}$$

$$\alpha_2 = \frac{\frac{5}{7}}{\frac{2}{7}} = \frac{5}{2}$$

$$\sum \alpha_i h_i =$$



New weights
@ x=3

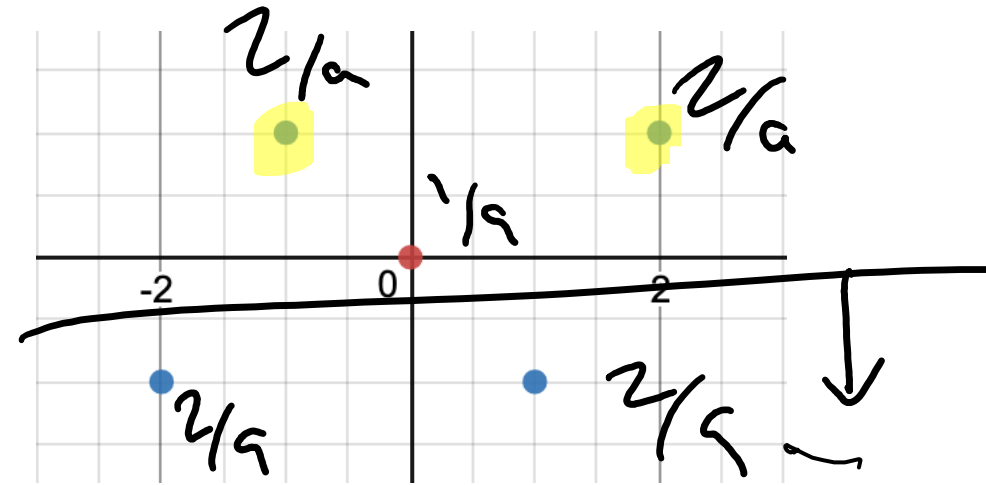
$$\frac{2}{7} + \frac{2}{7} + \frac{1}{7} + \frac{2}{7} + \frac{2}{7} = \frac{9}{7}$$

Boosting

Perform 3 iterations of the boosting algorithm on the following dataset. Double the error for incorrect points at each iteration and weight each model with $\alpha_i = (1-\epsilon)/\epsilon$

$$\epsilon_3 = \frac{2}{9} + \frac{2}{9} = \frac{4}{9}$$

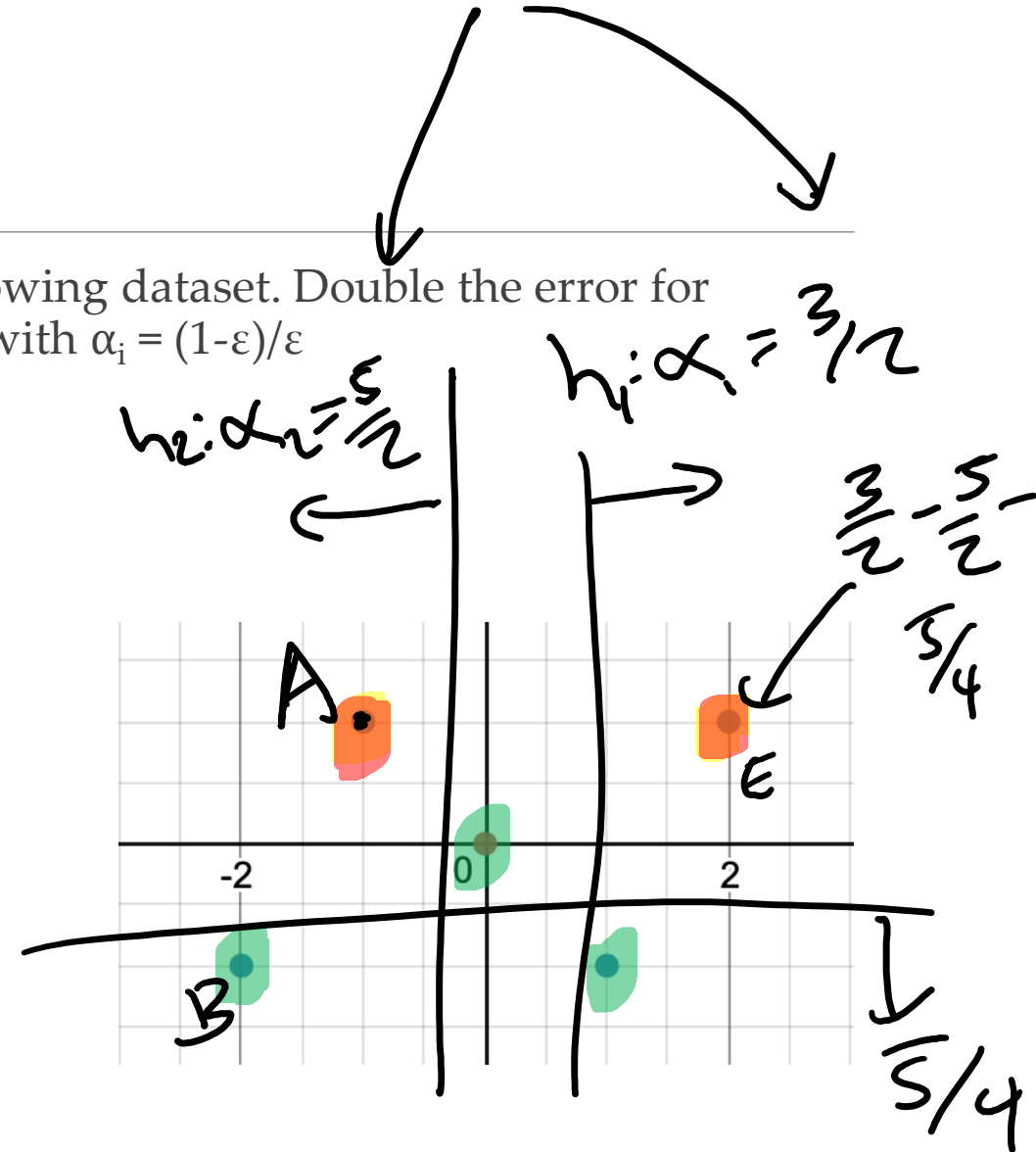
$$\alpha_3 = \frac{\frac{5}{9}}{\frac{4}{9}} = \frac{5}{4}$$



Boosting

Perform 3 iterations of the boosting algorithm on the following dataset. Double the error for incorrect points at each iteration and weight each model with $\alpha_i = (1-\epsilon)/\epsilon$

Class A: $h_1: -\frac{3}{2}$ $\frac{2}{2}$
 $h_2: \frac{5}{2}$
 $h_3: -\frac{5}{4}$
 $\sum: -\frac{1}{4}$



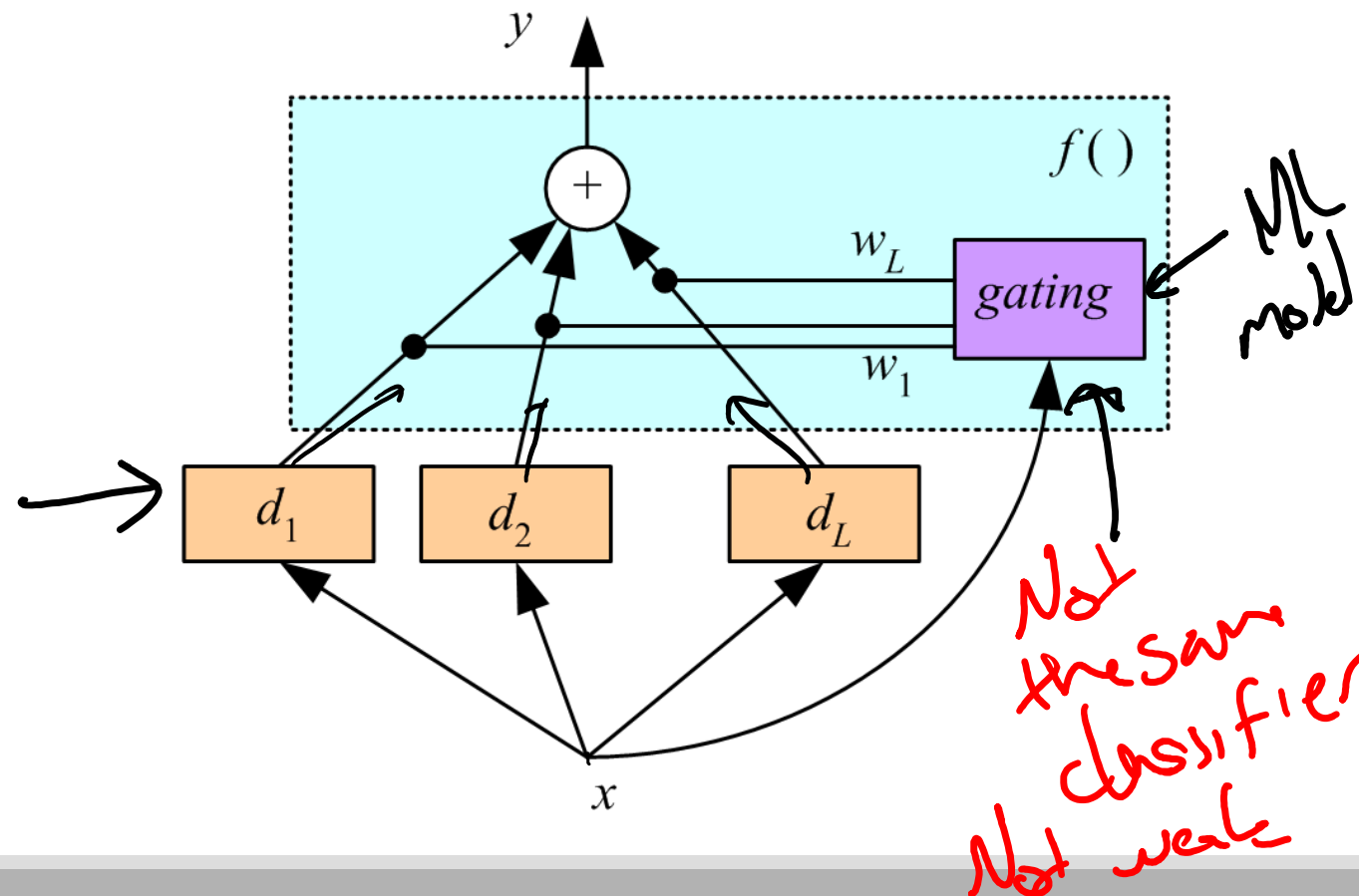
Mixture of Experts: Gating

Voting where weights are input-dependent (gating)

$$y = \sum_{j=1}^L w_j(x) d_j$$

(Jacobs et al., 1991)

Experts or gating can be nonlinear

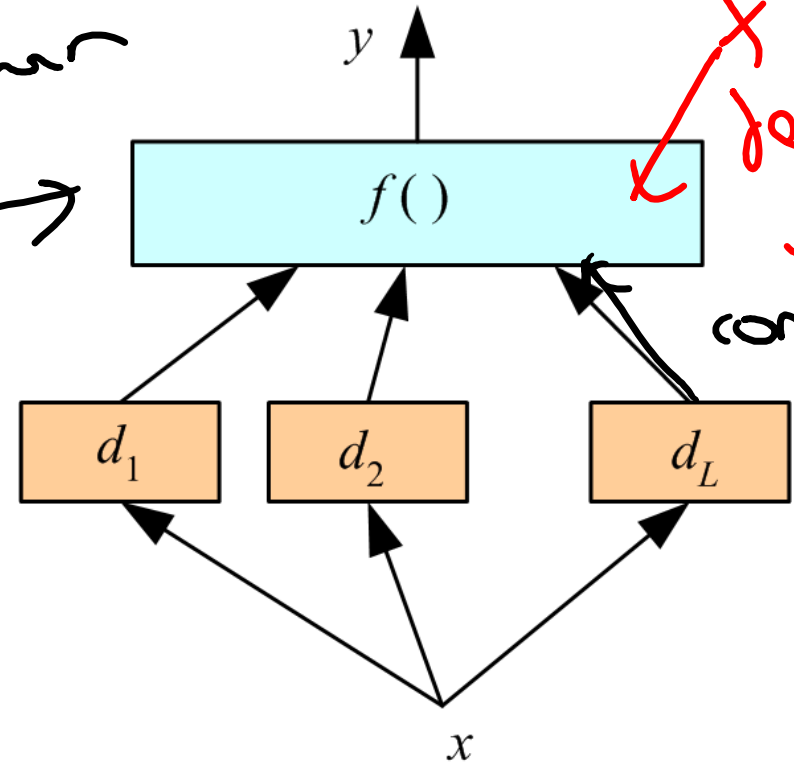


Stacking

Combiner $f()$ is another learner (Wolpert, 1992) that corrects the bias of base learners

Combiner f should be trained on data unused in training the base learners

rather than
just



f learns dependencies w/in confidence learners

high confidence for d_i might impact accuracy for d_k

for cancer

Cascading

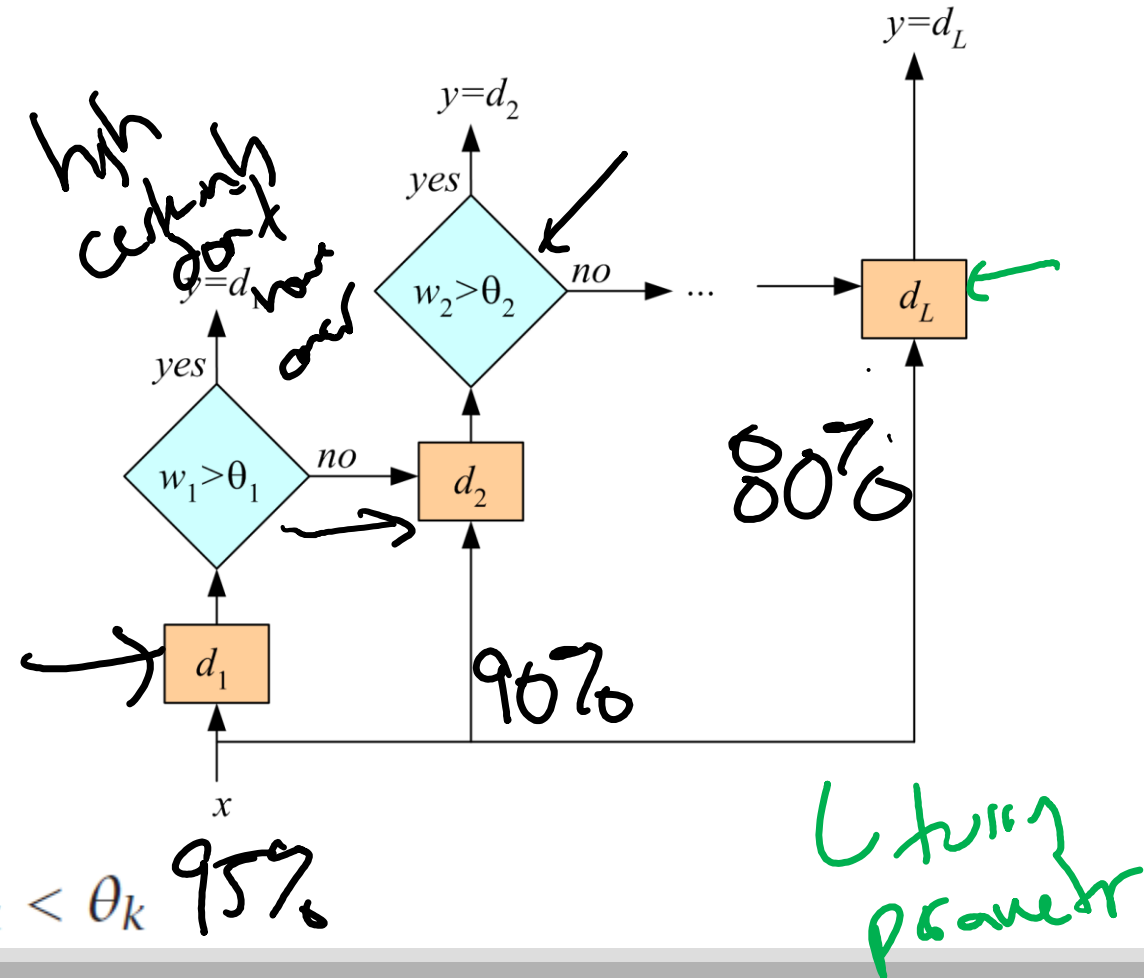
Cascade learners in order of complexity

Use d_j only if preceding ones are not confident (w_i)

asymmetric error

$$1/K < \theta_j \leq \theta_{j+1} < 1$$

$$y_i = d_{ji} \text{ if } w_j > \theta_j \text{ and } \forall k < j, w_k < \theta_k$$



Fine-Tuning an Ensemble

Given an ensemble of dependent classifiers, do not use it as is, try to get independence

- Inaccurate base learners can worsen accuracy (think of majority voting)

bias reduction

not too weak 60-70%

Subset selection: Forward (growing)/Backward (pruning) approaches to improve accuracy/diversity/independence

give each model a different subset of the features

can also select a subset of the learners

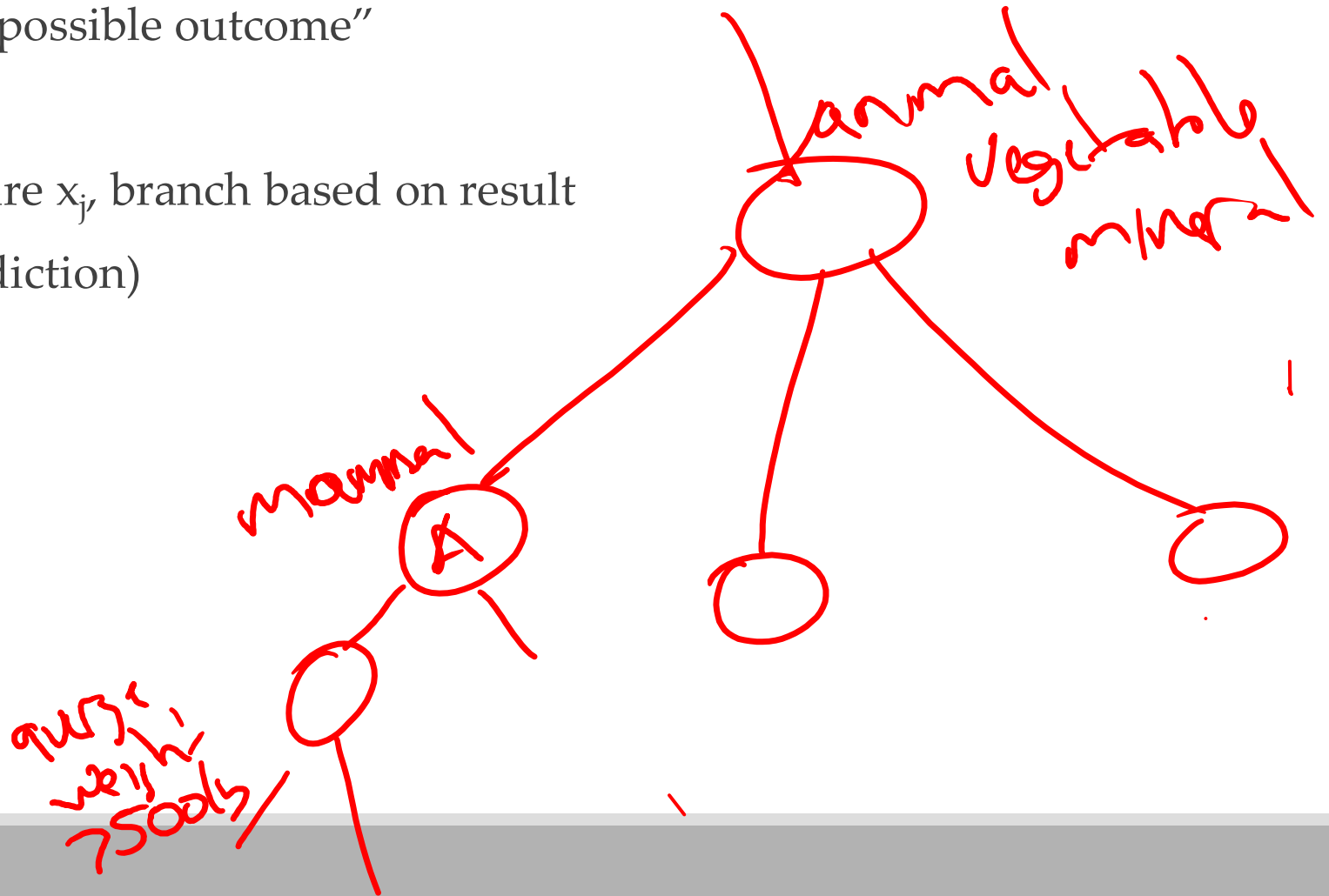


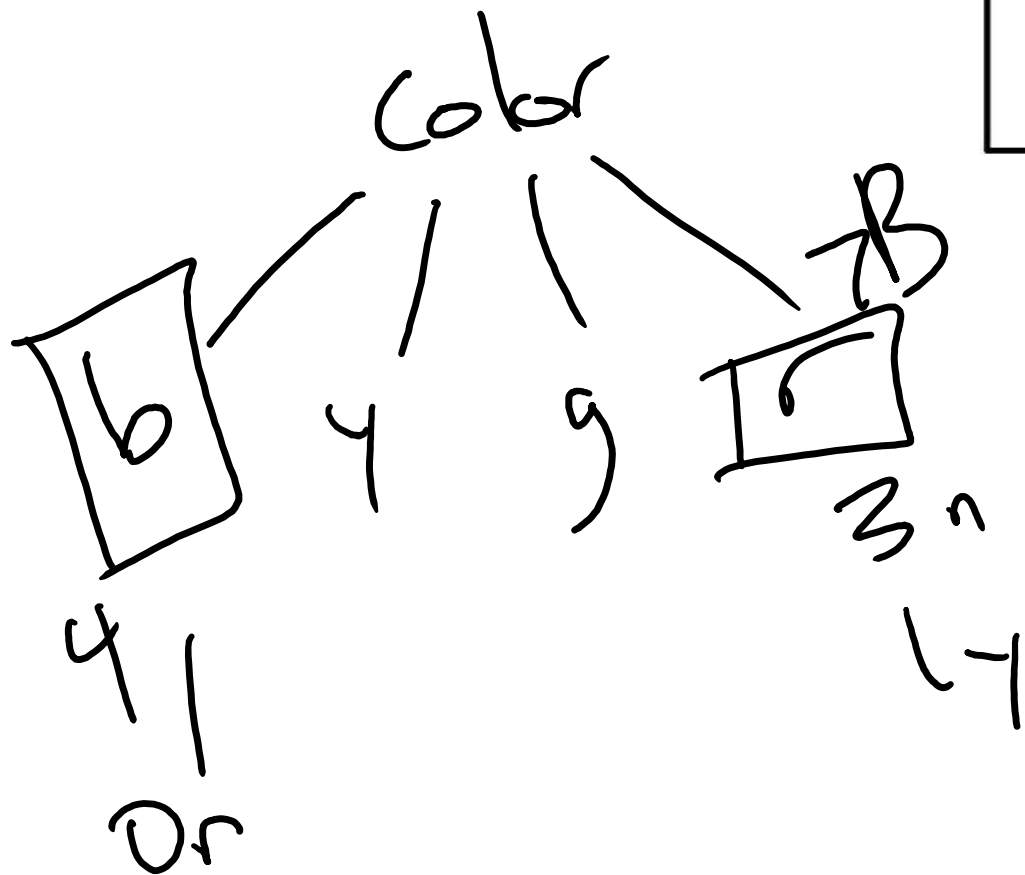
Decision Trees

“20 questions game for each possible outcome”

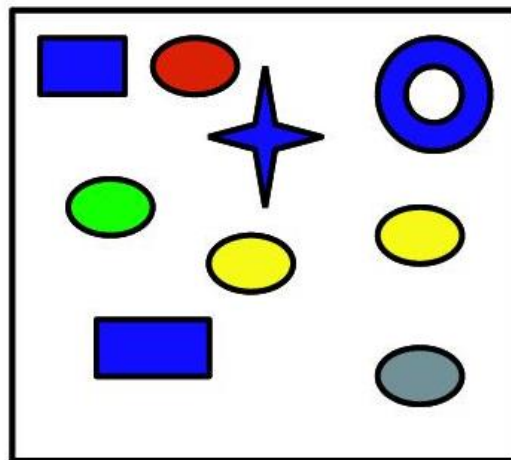
Nodes: test the value of feature x_j , branch based on result

Leafs: provide the class (prediction)

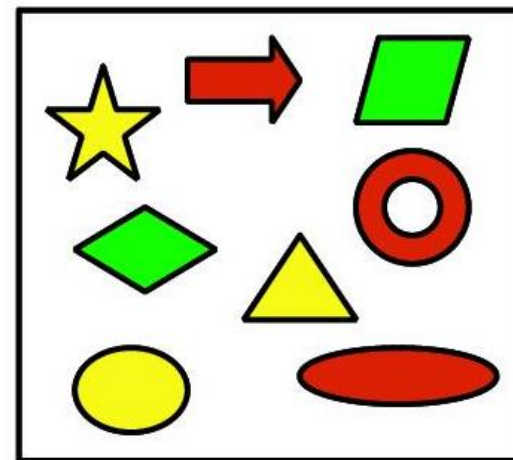




yes

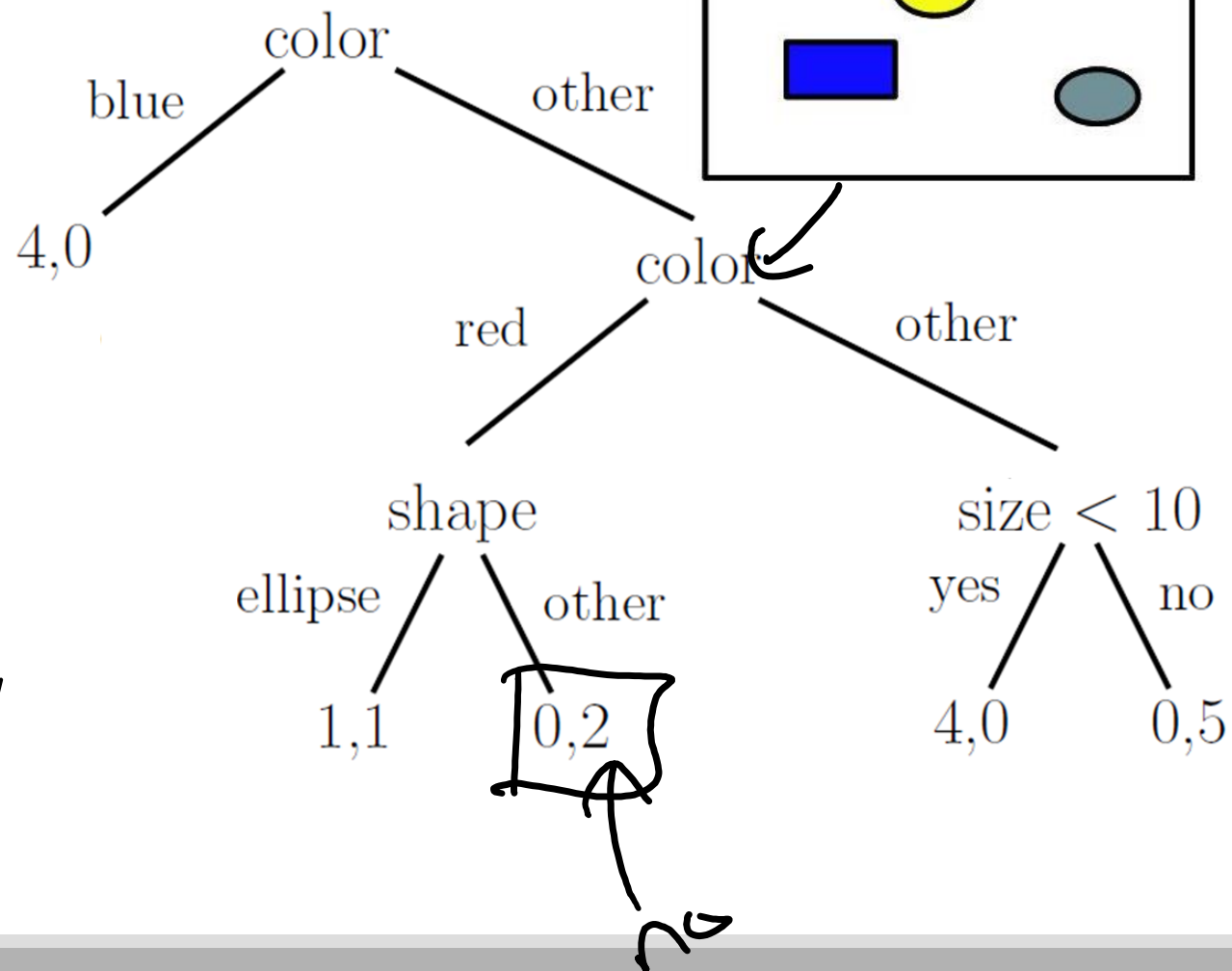
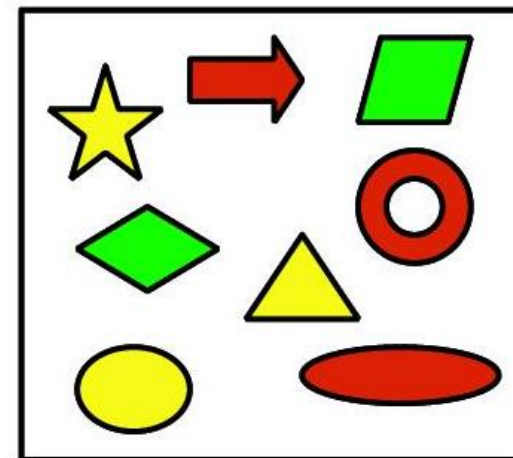
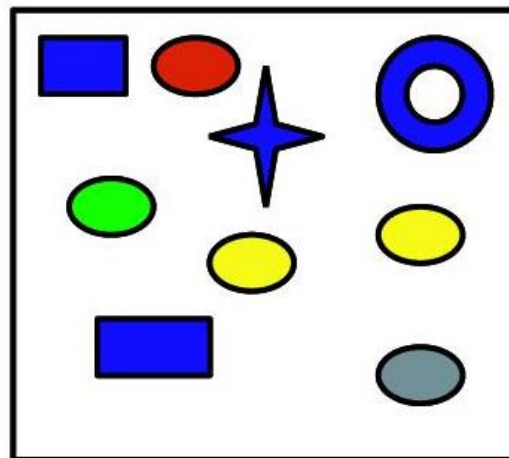


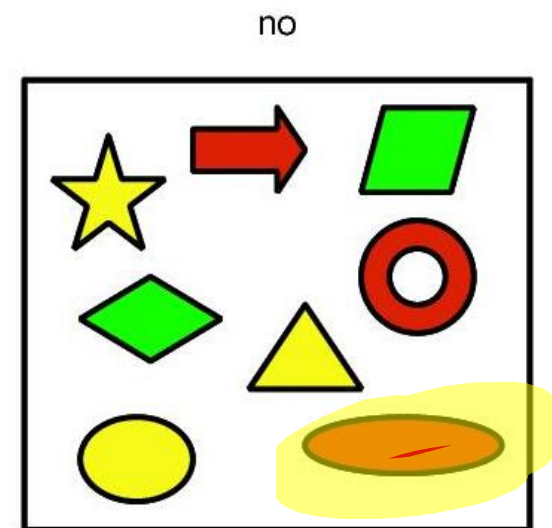
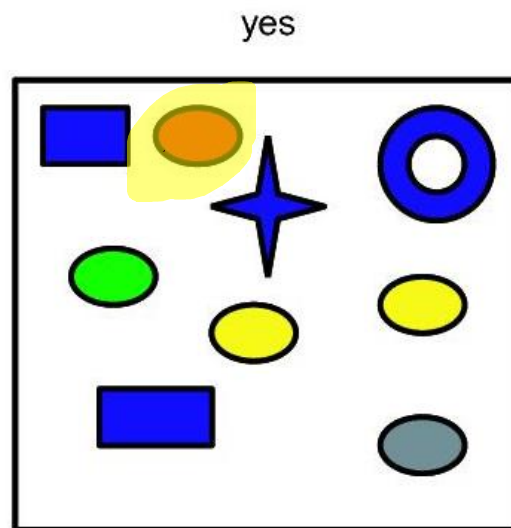
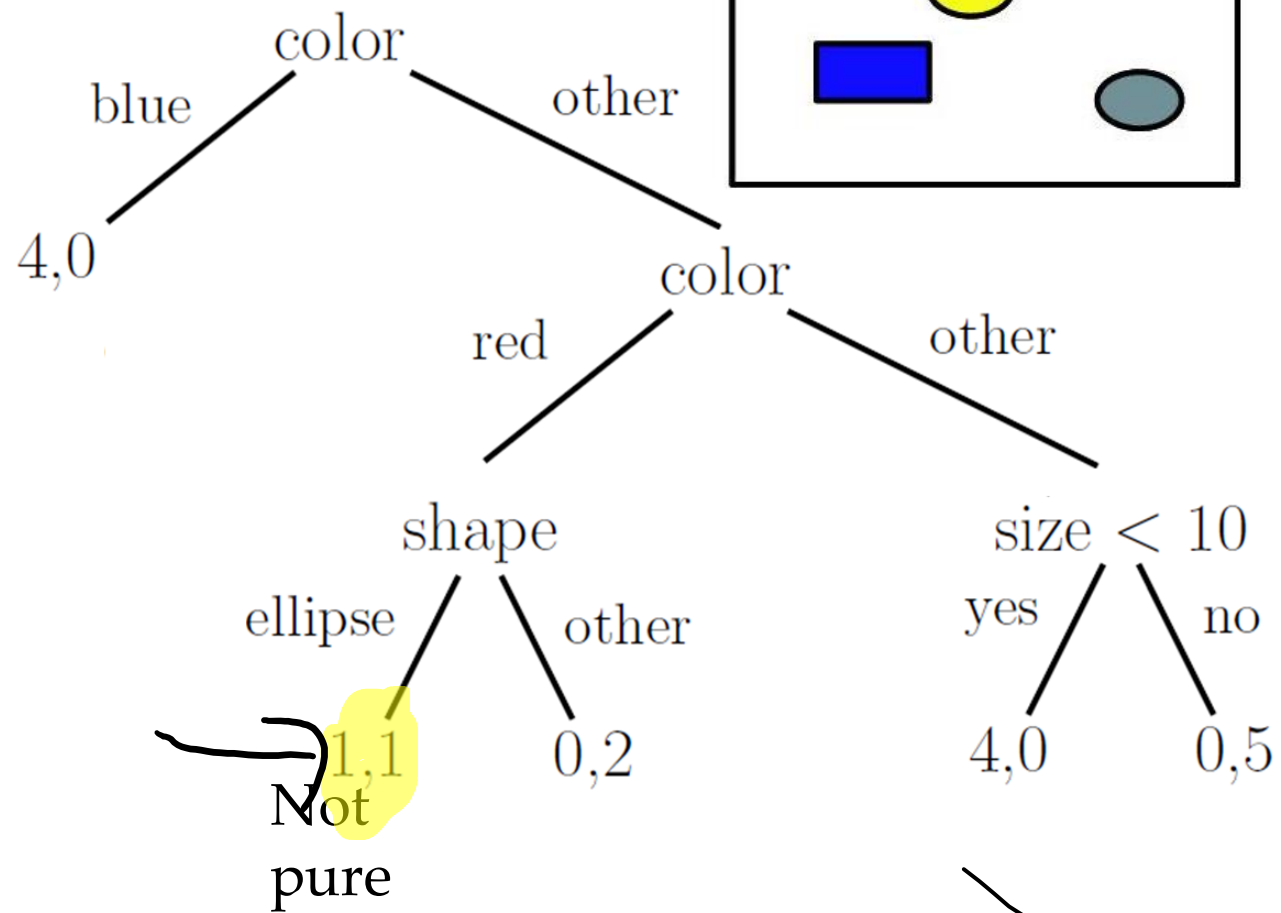
no

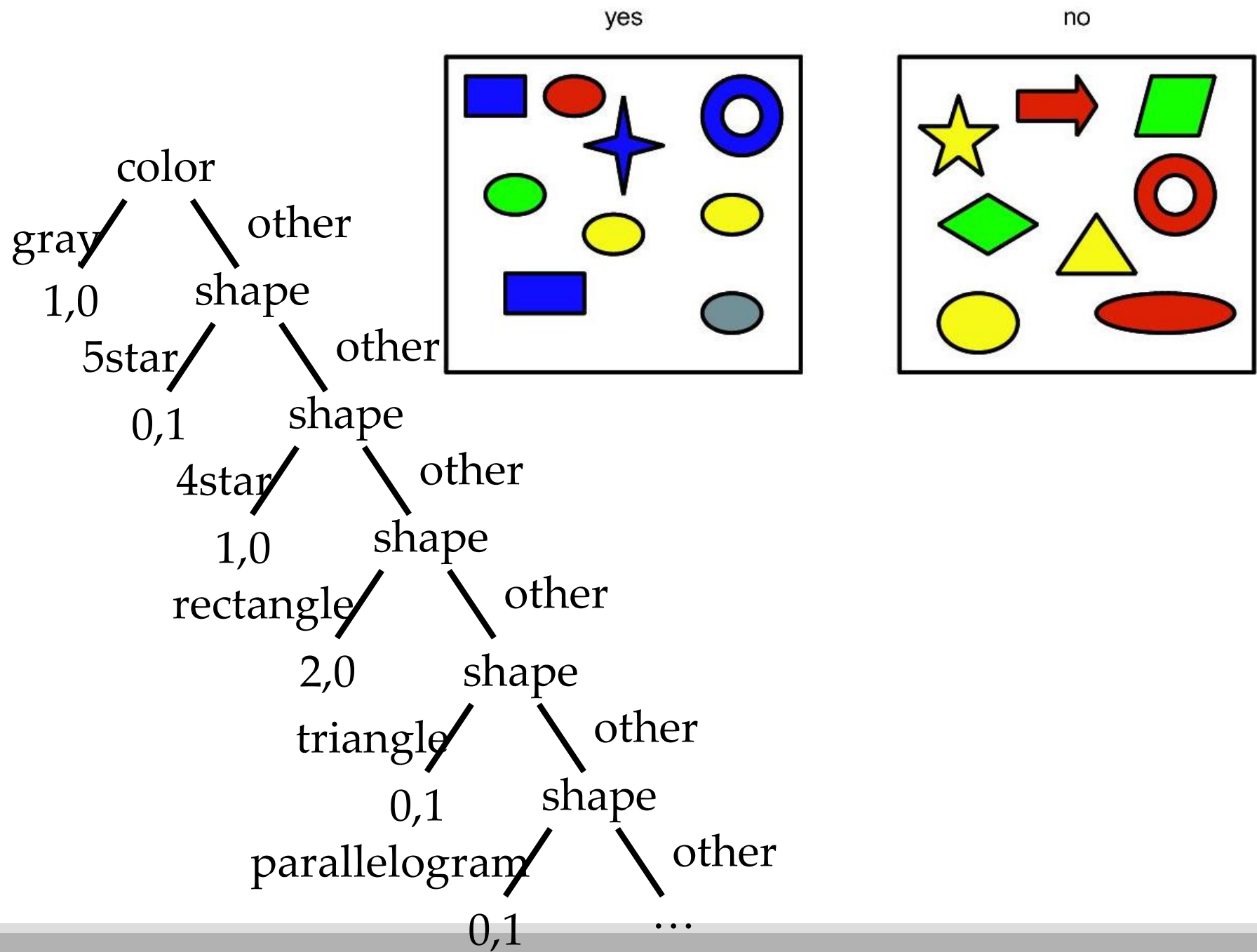


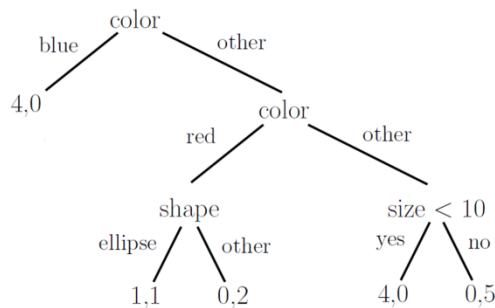
yes

no



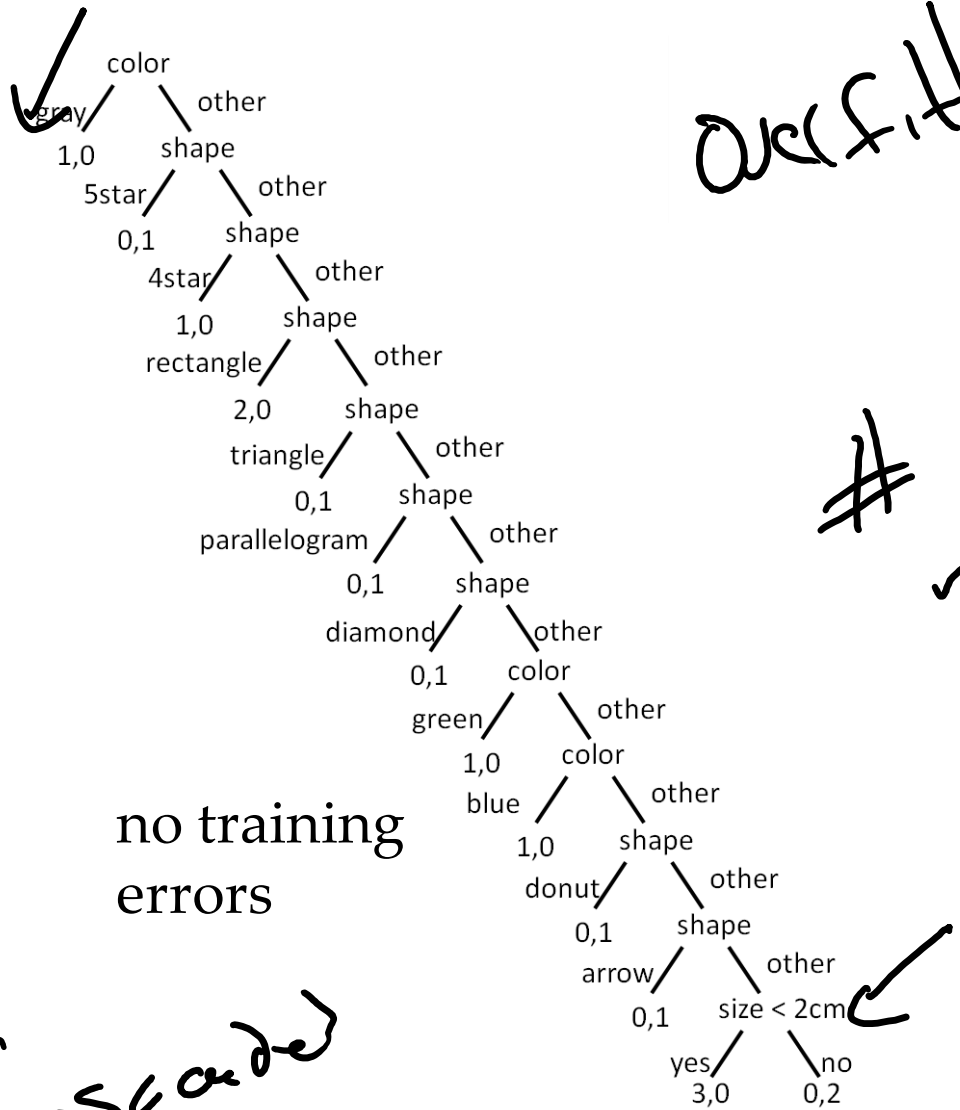






1 training error

or



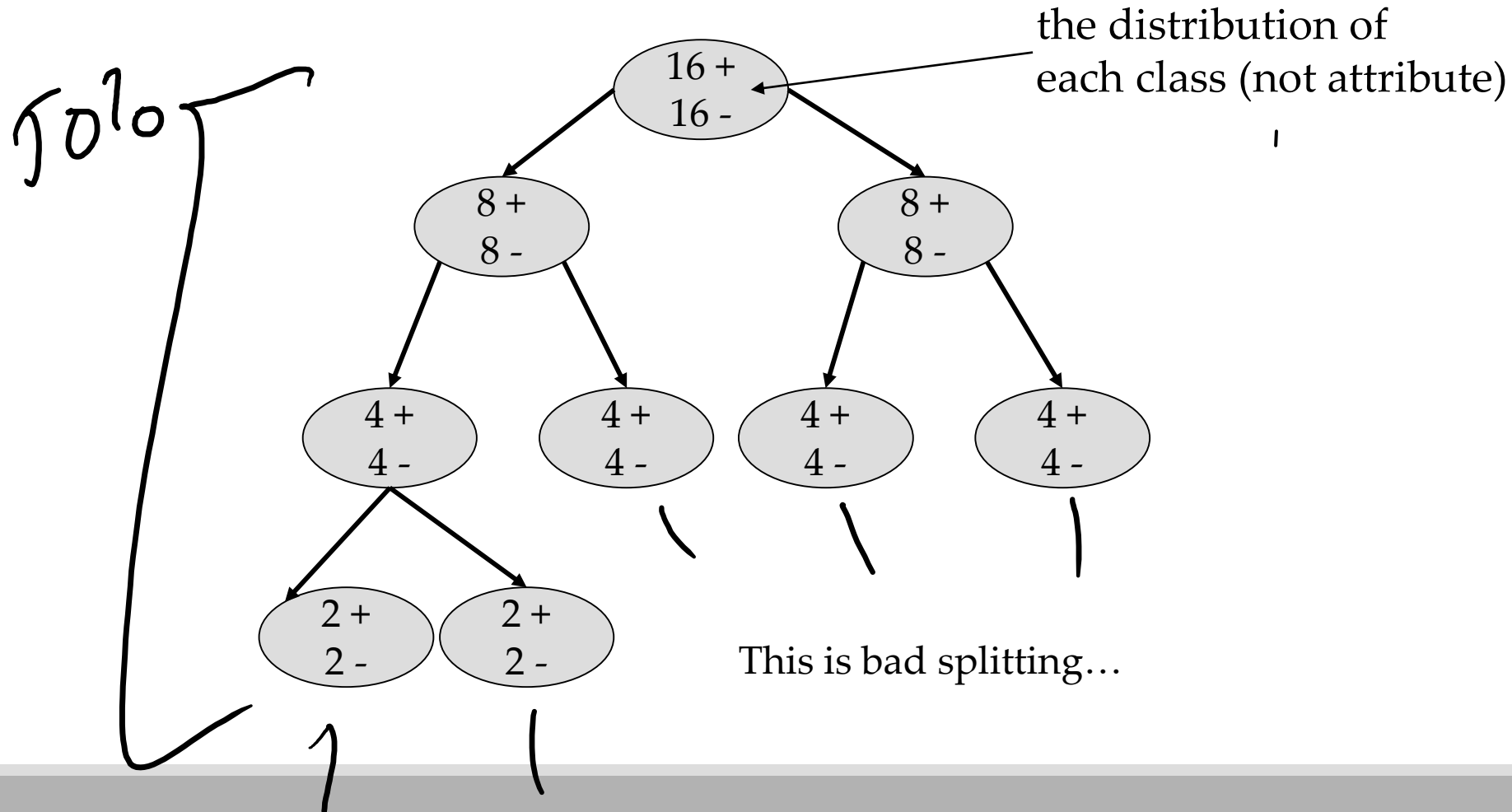
no training errors

overfitting

decision nodes

Cascaded
Decision
stump = tree w/ depth,

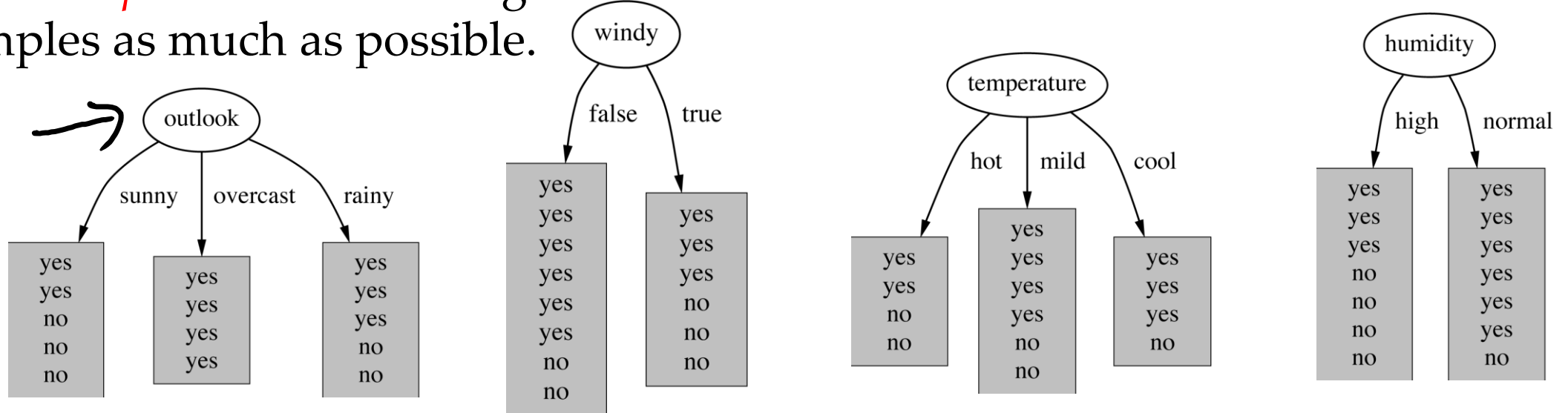
Which attribute to select for splitting?



How do we choose the test ?

Which attribute should be used as the test?

Intuitively, you would prefer the one that *separates* the training examples as much as possible.



Decision tree: divide and conquer

Supervised learning: classification and regression

Internal decision nodes implements a test function

↙ Univariate: Uses a single attribute, x_i – **This is used most frequently**

- Numeric x_i : Binary split : $x_i > w_m$ ↩
- Discrete x_i : n-way split for n possible values or binary

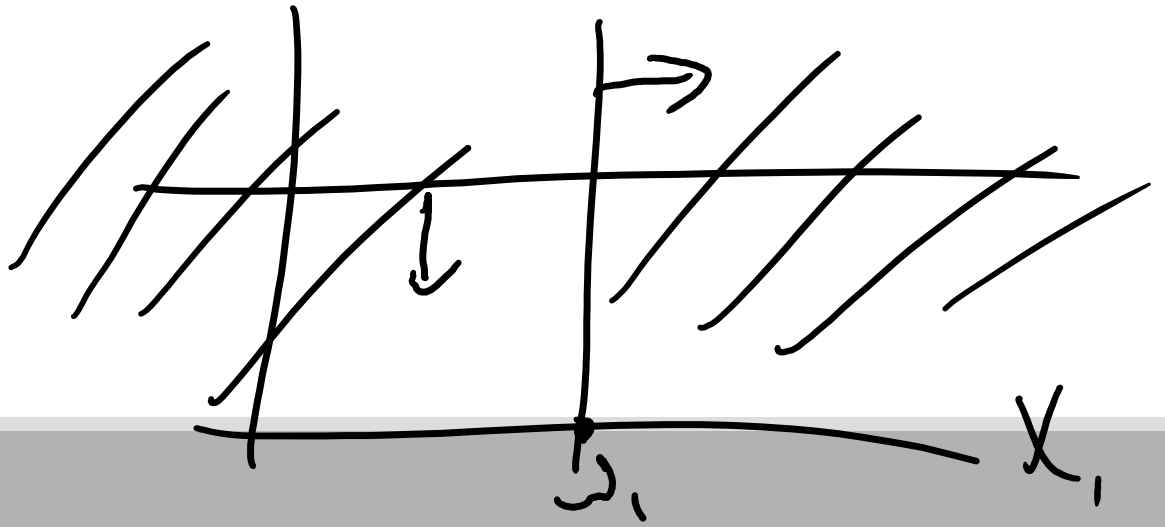
↓ Multivariate: Uses multiple/all attributes, x

Leaves

- ↘
- Classification: Class labels, or proportions
 - Regression: Numeric; r average, or local fit

Highly interpretable

only ask a question
about 1 feature
at a time



Classification Trees

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
- $f_m(x) : x_j > w_{m0}$ (threshold)

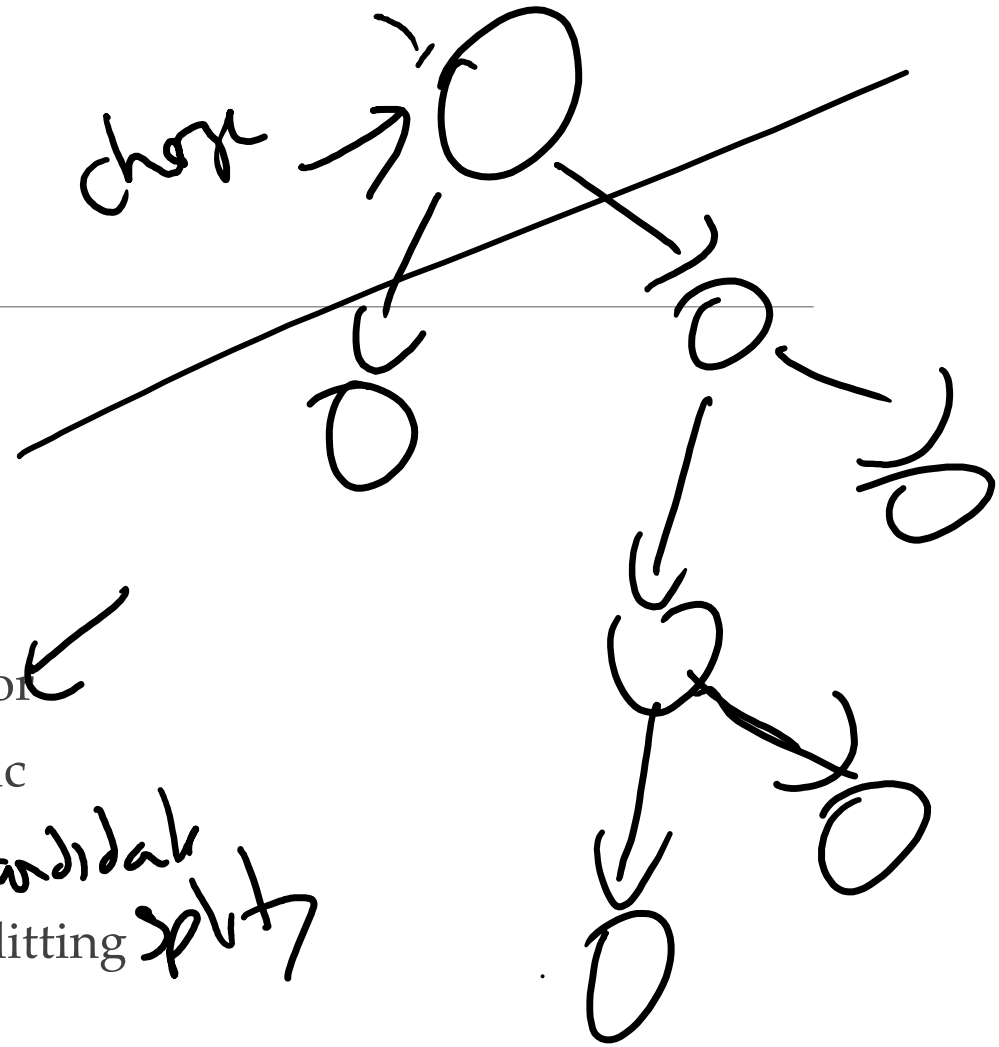
Goal: find the smallest tree that has low/zero training error

→ **NP-complete**, forced to use local search based on heuristic

Greedy: each step we look for the best split among candidates
 Score(D_1, D_2, \dots, D_k) measuring "goodness" of splitting the data into k subsets

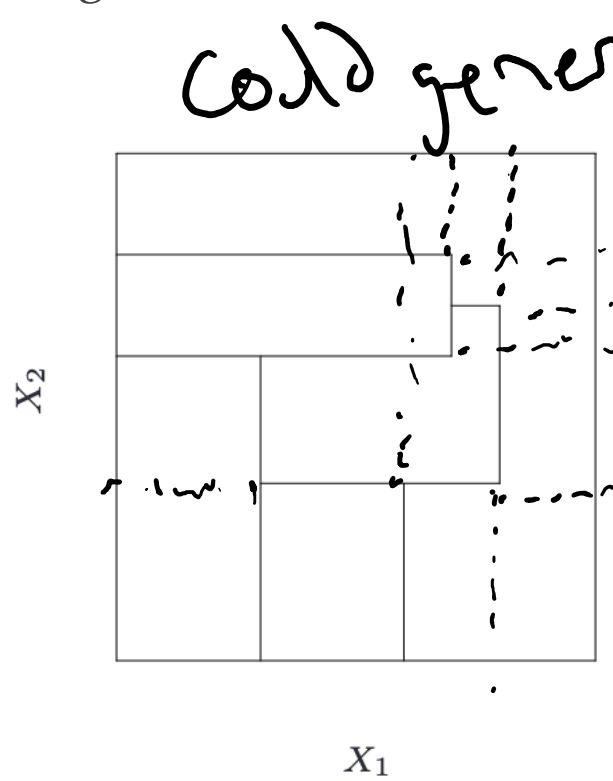
Continue recursively until no more split (leaf node)

↑
purity threshold is

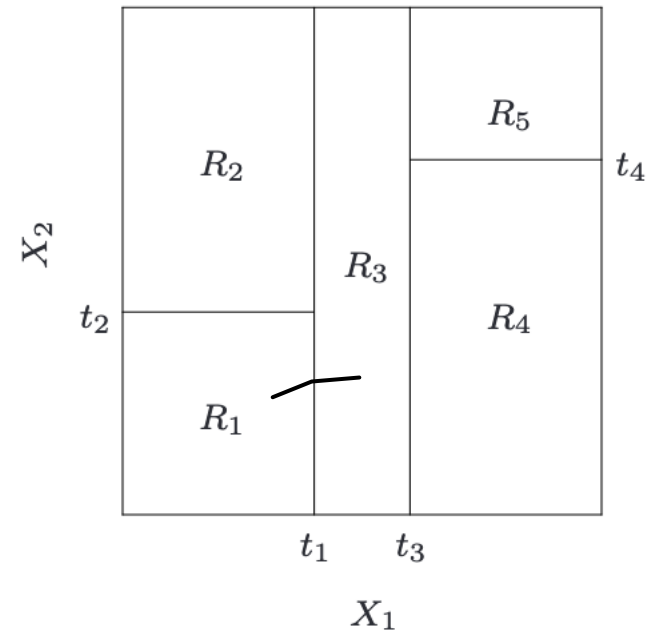


Region Types

What sorts of regions can the decision tree make?



a

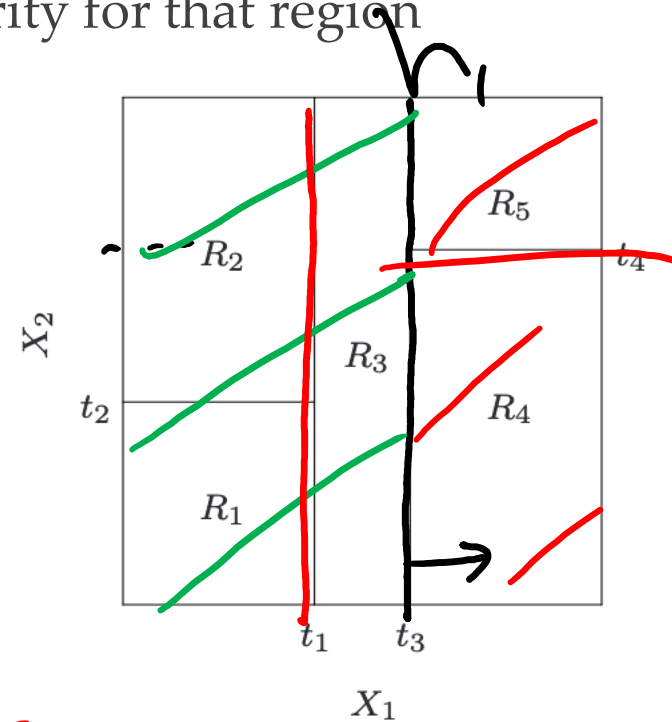
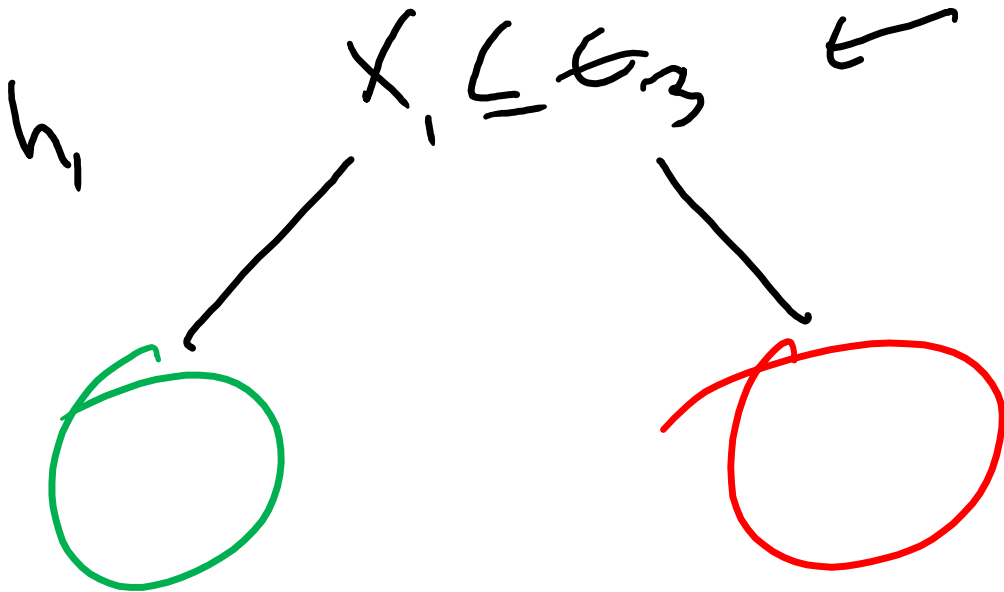


b

Region Types

What sorts of regions can the decision tree make?

Each “bin” is decided by whatever is the majority for that region

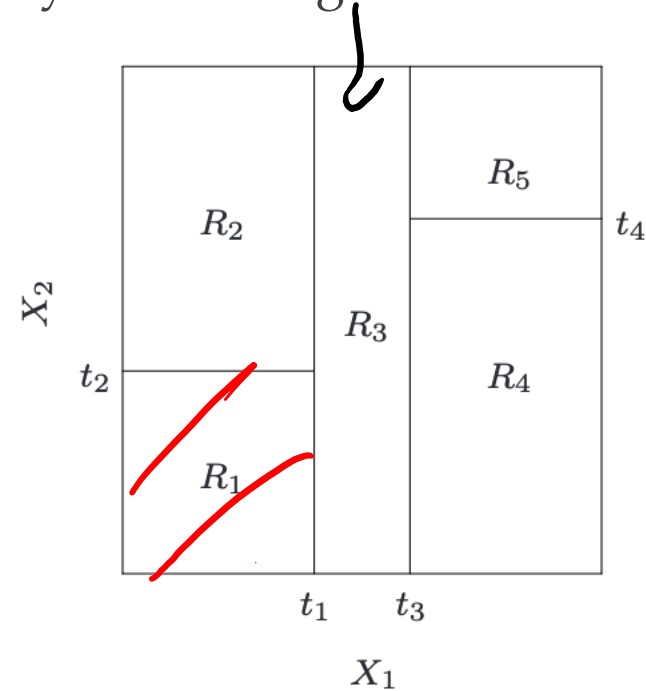
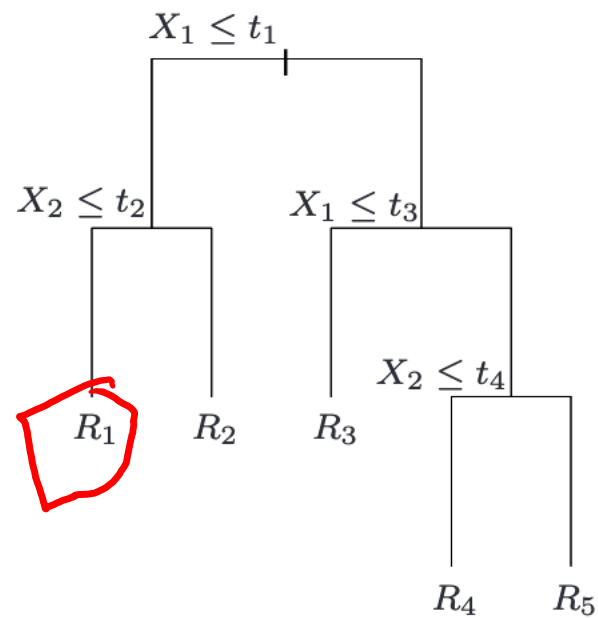


equivalent to DT
on next slide

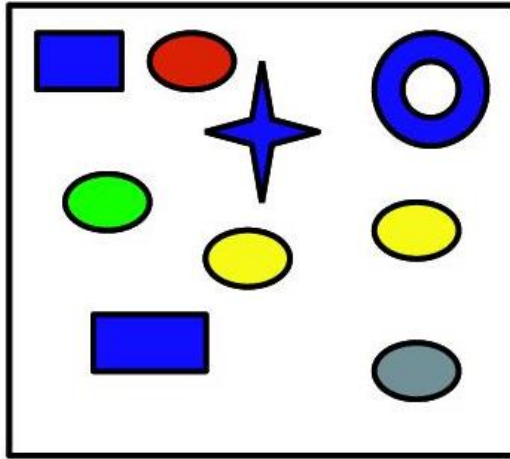
Region Types

What sorts of regions can the decision tree make?

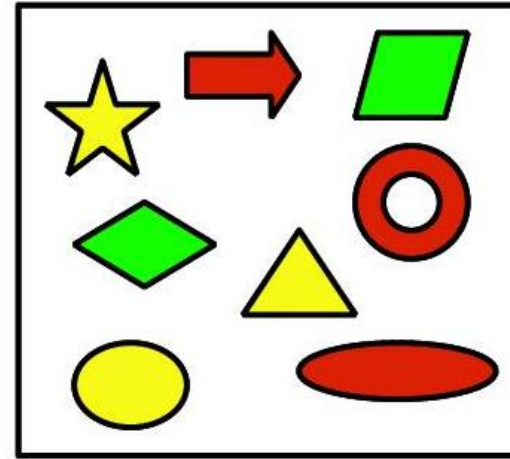
Each “bin” is decided by whatever is the majority for that region



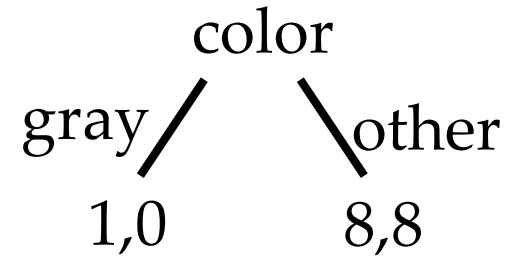
yes



no



Error rate: 5/17



Error rate: 8/17

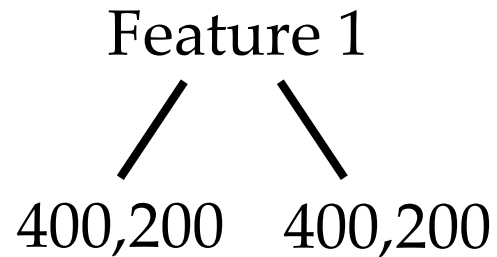
Use majority label at the leaf, then compute error rate

Accuracy score pitfall

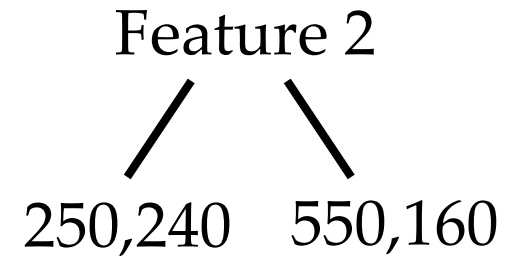
Feature 1
/ \
400,200 400,200

Feature 2
/ \
250,240 550,160

Accuracy score pitfall

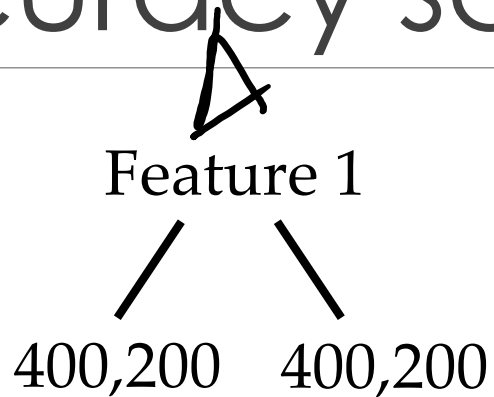


Error rate: $(200+200)/1200$

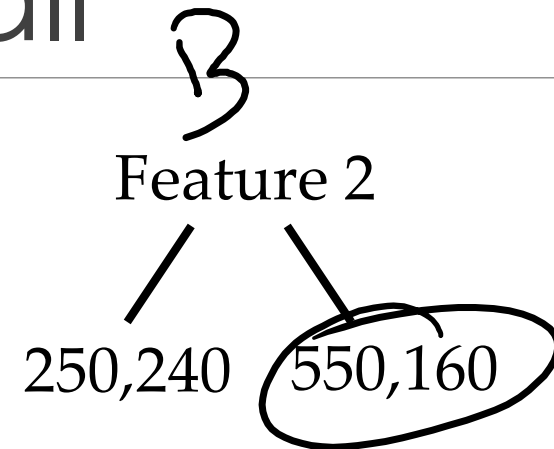


Error rate: $(240+160)/1200$

Accuracy score pitfall



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



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

more push

Both have the same error rate!!!

Which is “progressing more” towards a lower error?