

**Non-Parametric Methods:
Nearest Neighbor
and
Decision Tree Classification
and
Bagging for Constructing Ensembles**

Parametric vs. Non-Parametric Learning

- **Parametric learning**

- ▶ define a space of models parameterized by a fixed number of parameters
- ▶ find model that best fits the data (by searching over the parameters)
- ▶ Example: Perceptron, PA, and SVMs

- **Non-Parametric learning**

- ▶ define a space of models that can grow in size with data
- ▶ find model that best fits the data
- ▶ “Non-parametric” means “Not-fixed”
- ▶ Example: K-nearest neighbor and decision trees

Nearest-Neighbor Classifier

Nearest-neighbor Classifier

- Given a set of training examples $(x_i, y_i): i = 1, 2, \dots, N$, and a test example x
- **Prediction Rule:**
 - ▲ Find the training point x_j such that the distance between x and x_j is minimum (If tied, break ties uniformly at random)
 - ▲ Output y_j

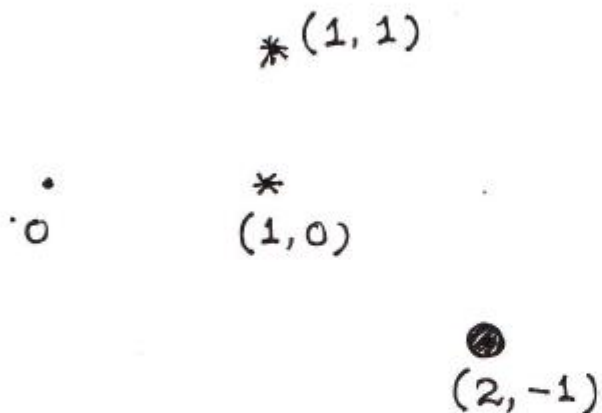
Example 1

Training data:

$((1, 0), 0)$, $((1, 1), 0)$, $((2, -1), 1)$

Test points:

$(0, 0)$, $(2, 1)$, $(1.5, -0.5)$



- $\text{dist}((0, 0), (1, 0)) = 1$

$\text{dist}((0, 0), (1, 1)) = \sqrt{2}$

$\text{dist}((0, 0), (2, -1)) = \sqrt{5}$

so: closest point to $(0, 0)$ is $(1, 0)$, Label = 0

Example 1 (contd.)

- $\text{dist}((2, 1), (1, 0)) = \sqrt{2}$

$\text{dist}((2, 1), (1, 1)) = 1$

$\text{dist}((2, 1), (2, -1)) = 2.$

closest: $(1, \overset{1}{\bullet})$

output: $y = \overset{1}{\bullet} 0$

- $\text{dist}((1.5, -0.5), (1, 0)) = \frac{1}{\sqrt{2}}$

$\text{dist}((1.5, -0.5), (1, 1)) = \sqrt{\frac{5}{2}}$

$\text{dist}((1.5, -0.5), (2, -1)) = \frac{1}{\sqrt{2}}$

closest: $(1, 0), (2, -1)$

Break ties at random

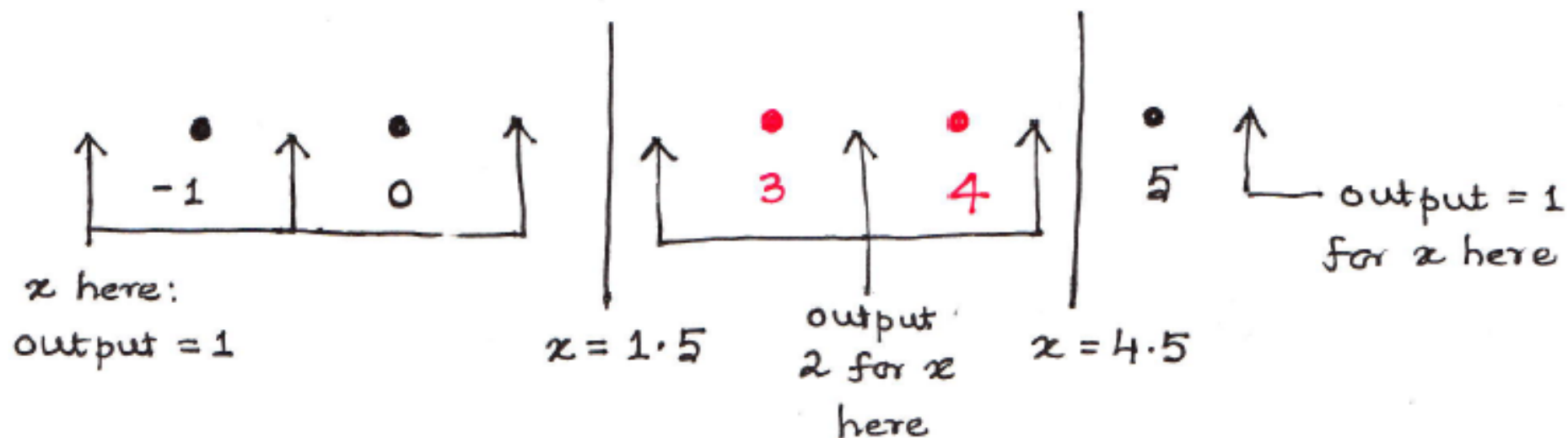
(report $y = 0$ w.p. $1/2$

$y = 1$ w.p. $1/2$)

Example 2

Training data:
 x is a scalar.

$(-1, 1)$, $(0, 1)$, $(3, 2)$, $(4, 2)$, $(5, 1)$

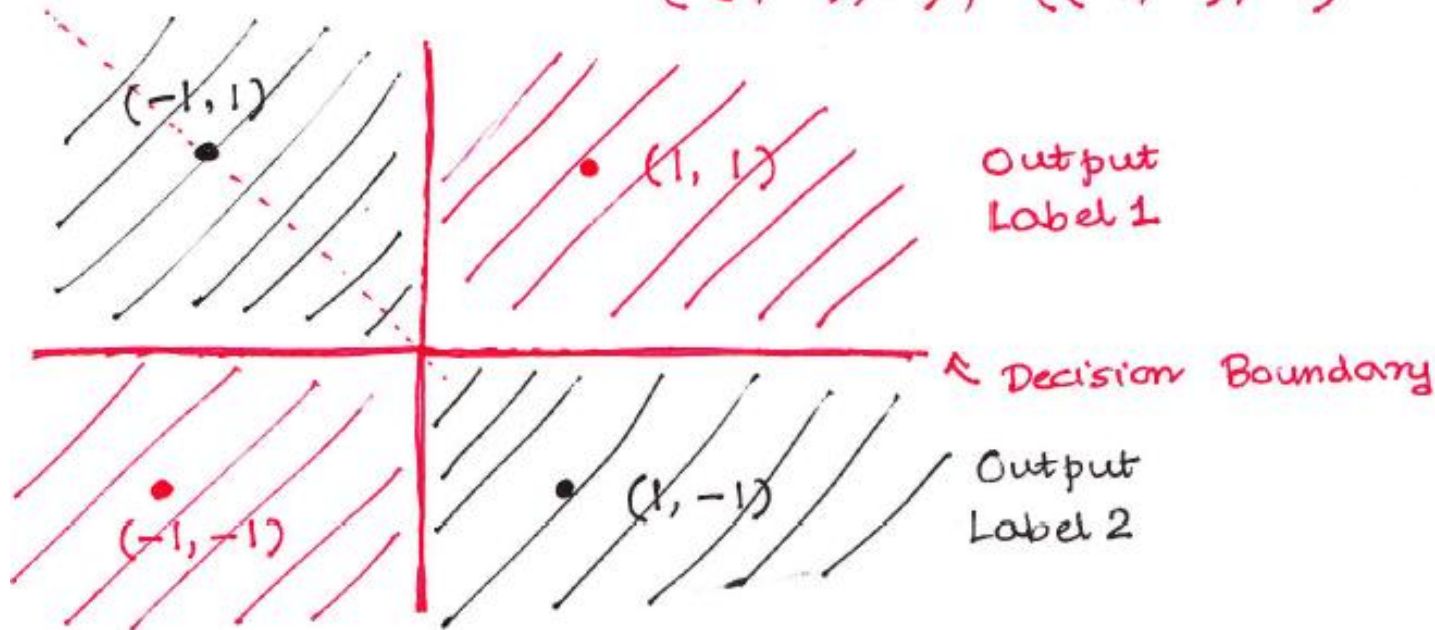


- **Decision boundary:**

- ▶ boundary between regions of different classes
- ▶ the output changes at the boundary

Example 3

Training data: $((1, 1), 1)$, $((-1, -1), 1)$,
 $((1, -1), 2)$, $((-1, 1), 2)$

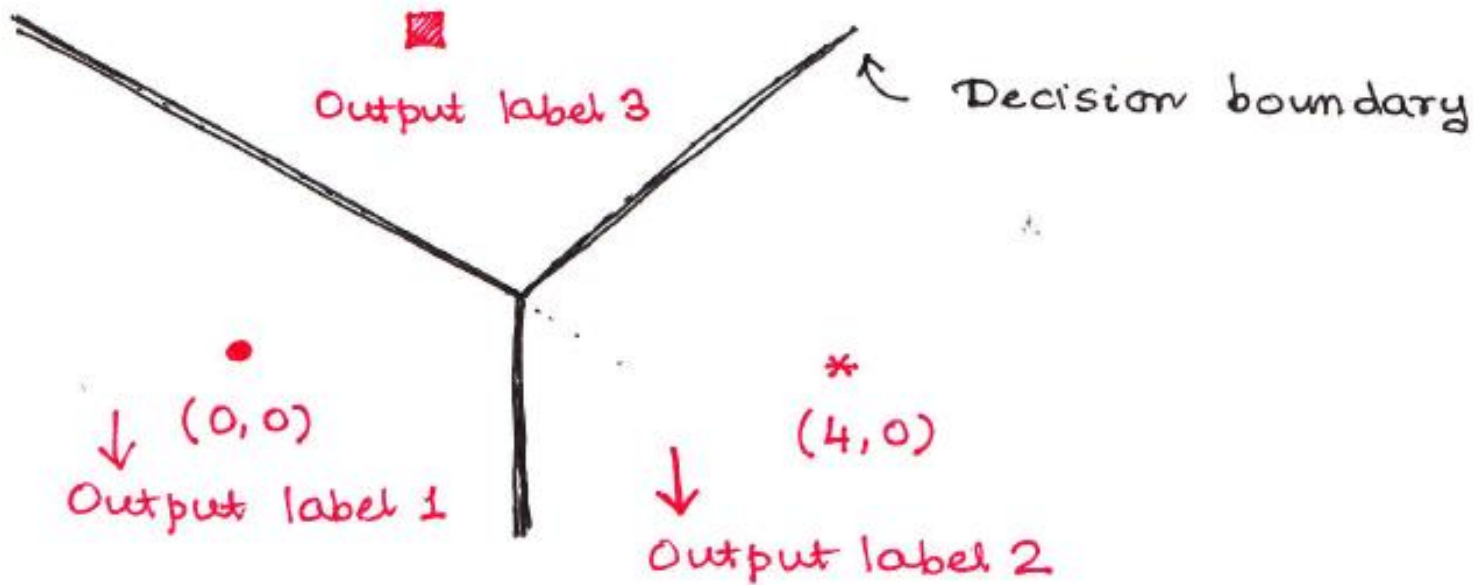


$$\|x - (1, 1)\| \leq \min [\|x - (1, -1)\|, \|x - (-1, -1)\|, \|x - (-1, 1)\|]$$

(Equation represents all vectors x which are closer to $(1, 1)$ than any other data point.)

Example 4

Training data: $((0,0), 1)$, $((4,0), 2)$, $((1,3), 3)$



When does NN work?

- Works well away from the decision boundary
- Not so well at the boundary
- Also does not work well when data is noisy

eg:



↑
suppose ~~this~~ noisy point
NN classifier does badly
around this point.

- How can we make NN more robust to noise?

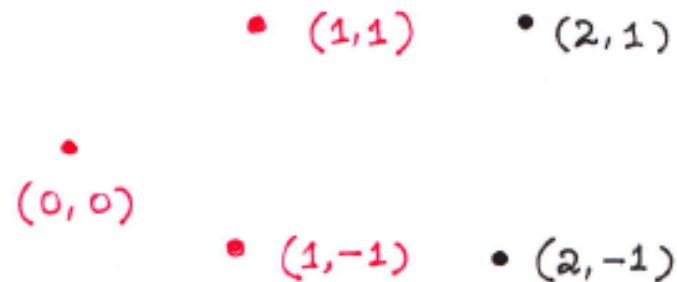
K-Nearest Neighbor Classifier

- Given a set of training examples $(x_i, y_i): i = 1, 2, \dots, N$, and a test example x
- **Prediction Rule:**
 - ▶ Find j_1, j_2, \dots, j_k , the indices of k points closest to x in the training data
 - ▶ Output the majority label among $y_{j_1}, y_{j_2}, \dots, y_{j_k}$
 - ▶ Majority label == one that occurs most often
 - ▶ If there is tie, resolve uniformly at random

Example 1: 3-NN

Training data:

$((0, 0), 0)$ $((1, 1), 0)$ $((1, -1), 0)$
 $((2, 1), 1)$ $((2, -1), 1)$



Test points: $(1, 0)$.

$$\text{dist}((1, 0), (0, 0)) = 1$$

$$\text{dist}((1, 0), (2, 1)) = \sqrt{2}$$

$$\text{dist}((1, 0), (1, 1)) = 1$$

$$\text{dist}((1, 0), (2, -1)) = \sqrt{2}$$

$$\text{dist}((1, 0), (1, -1)) = 1$$

closest 3 points: $(0, 0)$, $(1, 1)$, $(1, -1)$

Their labels : 0 0 0

So output = 0.

Example 1: 3-NN

Test point: $(2, 0.5)$

$$\begin{aligned} \text{dist}((2, 0.5), (0, 0)) &= \frac{\sqrt{17}}{2} & \text{dist}((2, 0.5), (2, 1)) &= \frac{1}{2} \\ \text{dist}((2, 0.5), (1, 1)) &= \frac{\sqrt{5}}{2} & \text{dist}((2, 0.5), (2, -1)) &= \frac{5}{2} \\ \text{dist}((2, 0.5), (1, -1)) &= \frac{\sqrt{13}}{2} \end{aligned}$$

Closest 3 points: $(2, 1), (1, 1), (2, -1)$

Labels: $1 \quad 0 \quad 1$

Majority: 1 = output label.

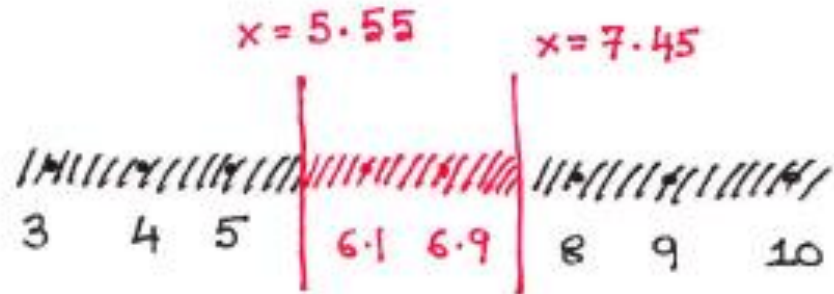
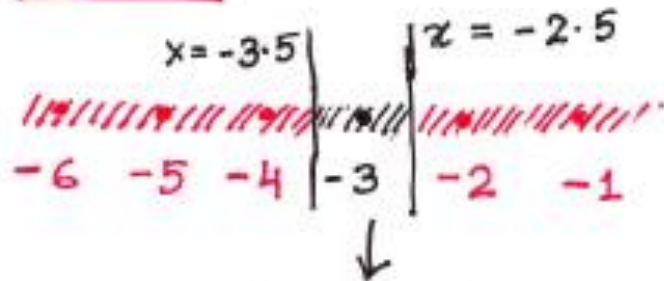
Example 2

• • • • •
-6 -5 -4 -3 -2 -1

• • • • •
3 4 5 6.1 6.9 8 9 10

Suppose the points at -3 and 6.1 and 6.9 are noisy.

1-NN:



Example 2

3-NN:



Output label is red
on this entire region



5-NN



Output label red
in this region

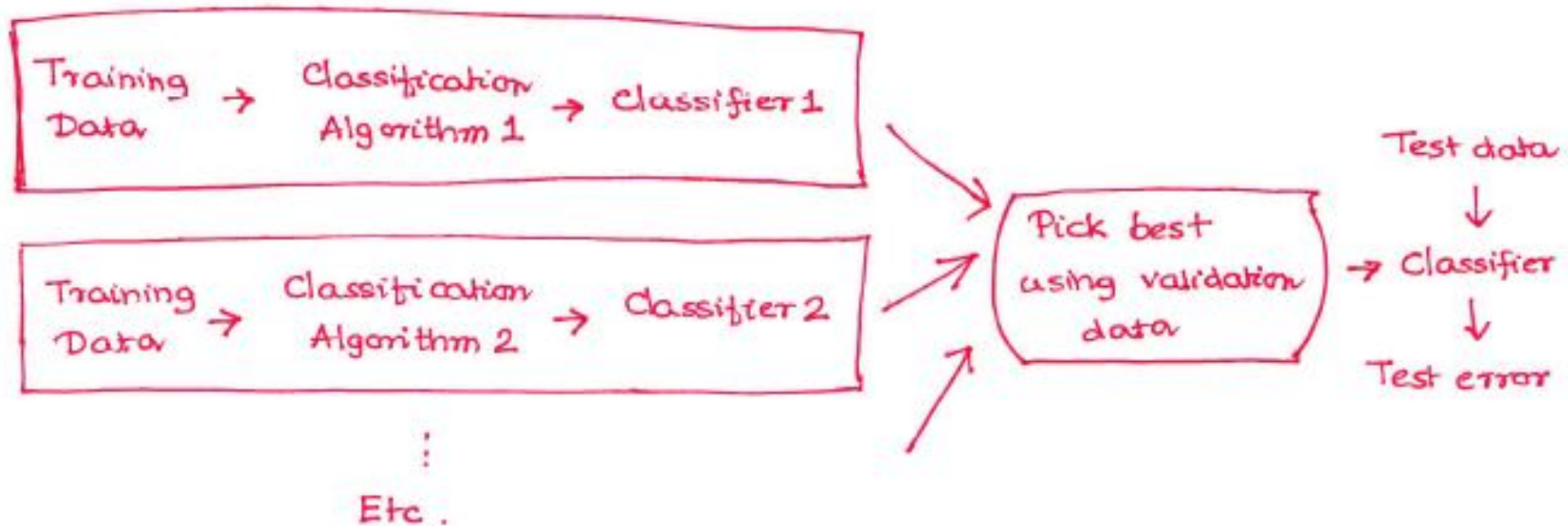


Output label black in this
entire region.

- How to choose “K” for K-NN classifier?

Choosing “K” via Cross-Validation

1. Split data into training set and validation set.
2. Train classifier on training set for $k = 1, 3, 5, \dots$
3. Evaluate the error of each classifier trained on validation set and pick the one with the lowest error.



Other design choices

- Distance measure?
 - ▲ Most common is Euclidean distance
 - ▲ Others used too
- How to find NNs?
 - ▲ In 1D, binary search will do -- $O(\log n)$
 - ▲ In higher-dimensional space, advanced data structures such as Locality Sensitive Hashing (LSH) is used
 - ▲ LSH: Hash the data points such that closer points are mapped to the same buckets with very high probability

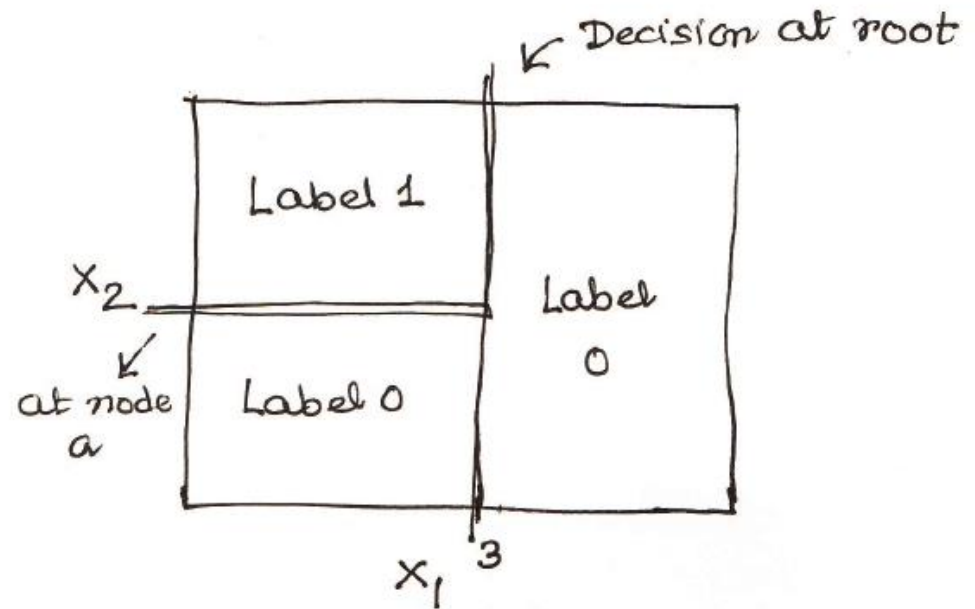
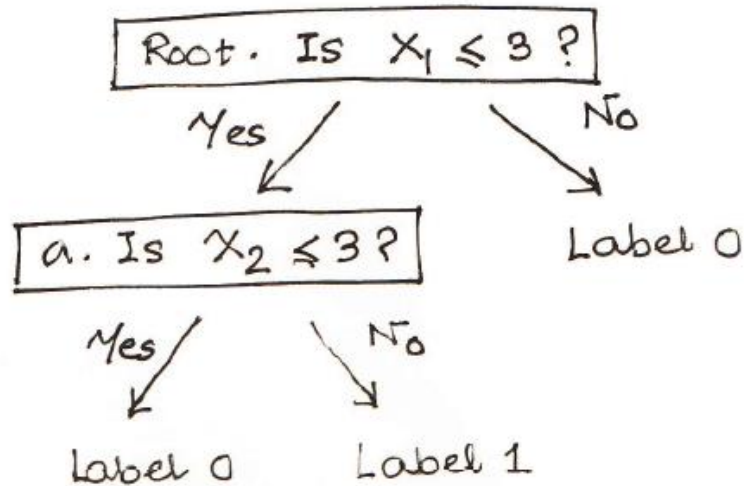
Summary

- K-NN is a non-parametric classifier
 - ▲ no training needed! Just store all the training data 😊
 - ▲ simple and easy to implement
 - ▲ flexible hypothesis space: infinitely complex
- Disadvantages
 - ▲ Classification time is high (NN computation!)
 - ▲ Space requirement is high
 - ▲ Doesn't work very well in high dimensions

Decision Tree Classifier

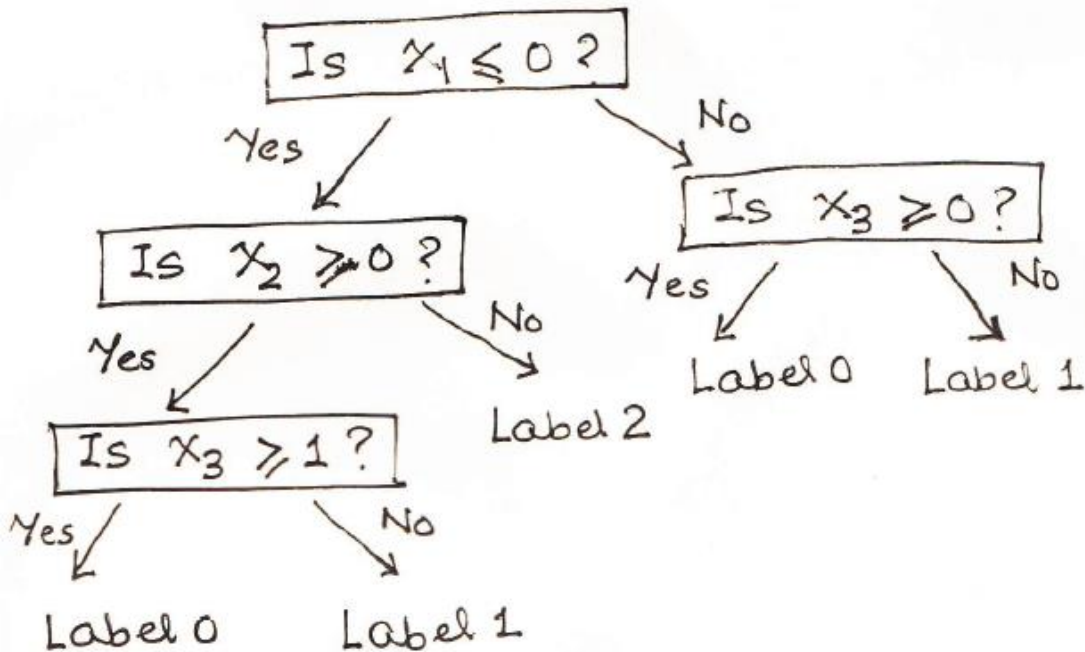
Decision Tree: Examples

Example 1: Features x_1, x_2



Decision Tree: Examples

Example 2: Features x_1, x_2, x_3



- Each node is based on a feature
- Each node: a decision is made based on the value of this feature
- Each leaf corresponds to a label (same label may appear on many leaves)

Decision Tree: Examples

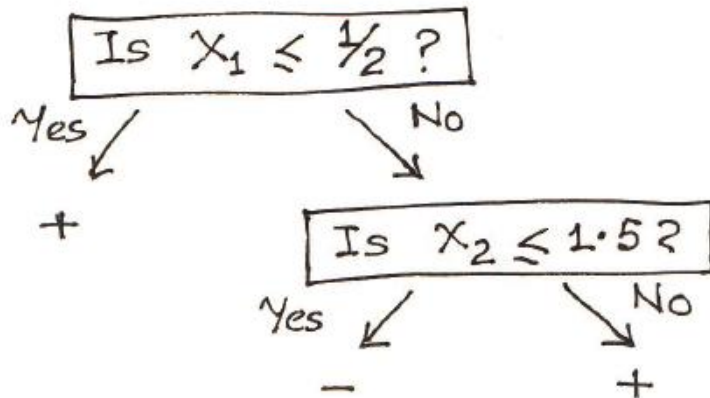
$(0, 2) +$ $+$ $(1, 2)$

$(0, 1) +$ $-$ $(1, 1)$

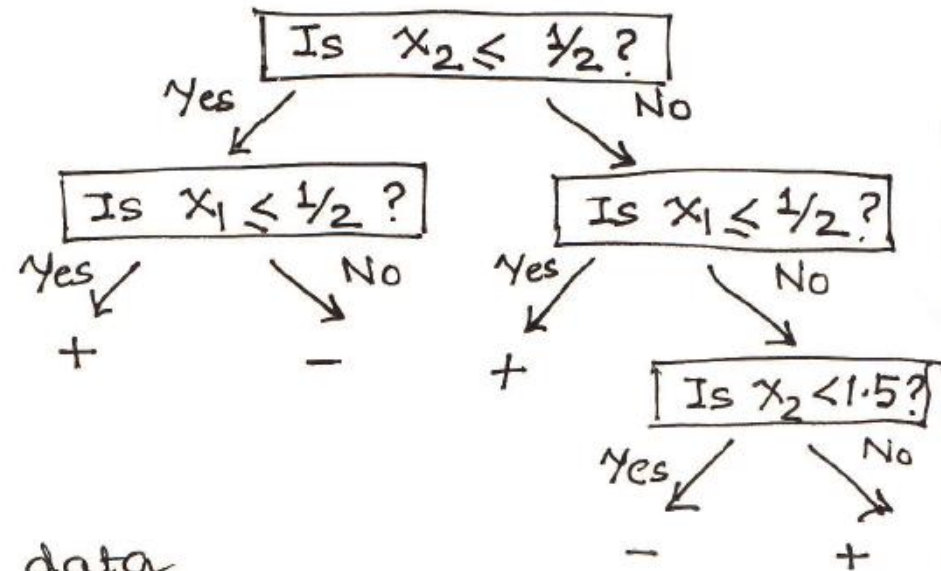
$+$ $-$

$(0, 0)$ ~~$(1, 0)$~~

Tree 1:



Tree 2:



- Both trees classify the training data correctly; but which one is better?

Decision Tree: Examples

- Answer: The simpler one. (here, tree 1)
- In general, finding the smallest decision tree to fit or training data set is computationally difficult.
- A ~~good~~ commonly used algorithm that proceeds greedily is the ID3 Decision Tree algorithm.
- Convention: Rule at each node of the form:
" Is $X_f \leq t$? "

- **Occam's Razor:**

- ▶ Pick the simplest hypothesis that explains the data
- ▶ Simpler solutions generalize well

ID3 Decision Tree Learning Algorithm

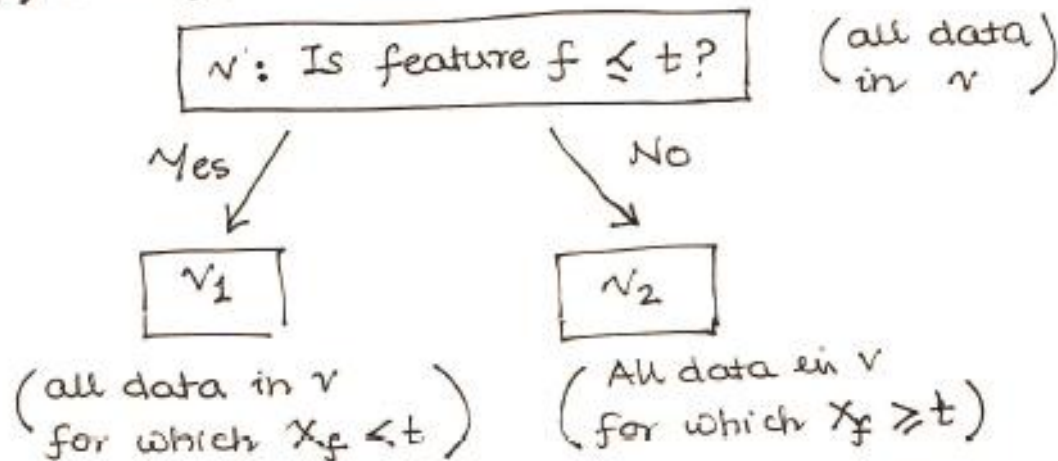
1. Initially, Whole training data is at root.

2. While there is an impure node:

(a) Pick any impure node v

(b) Pick a feature f and threshold t along which to "split" the data at v . Done according to "Splitting Rule" (to be described later)

(c) Modify tree as:



If any of v_1 or v_2 is pure (i.e. has data of only one label), then ~~predict~~ make it a leaf that predicts this label.

Impure Node: Node with data points with multiple labels

Each node in the algorithm corresponds to a subset of the training data.

ID3 Algorithm: Running Example

Notes:

1. Convention: Use a threshold "in the middle" of two values
eg. ~~✖~~ Node has data points: $(0,0)$, $(0,1)$, $(1,1)$, and if we split along feature 1, use threshold $\frac{1}{2}$.
2. More discussion of splitting rule coming up.

Example: Training data: $\begin{matrix} a & b & c \\ ((0,0), 1), & ((0,1), 0), & ((1,0), 0), \\ & d \\ & ((1,1), 0) \end{matrix}$

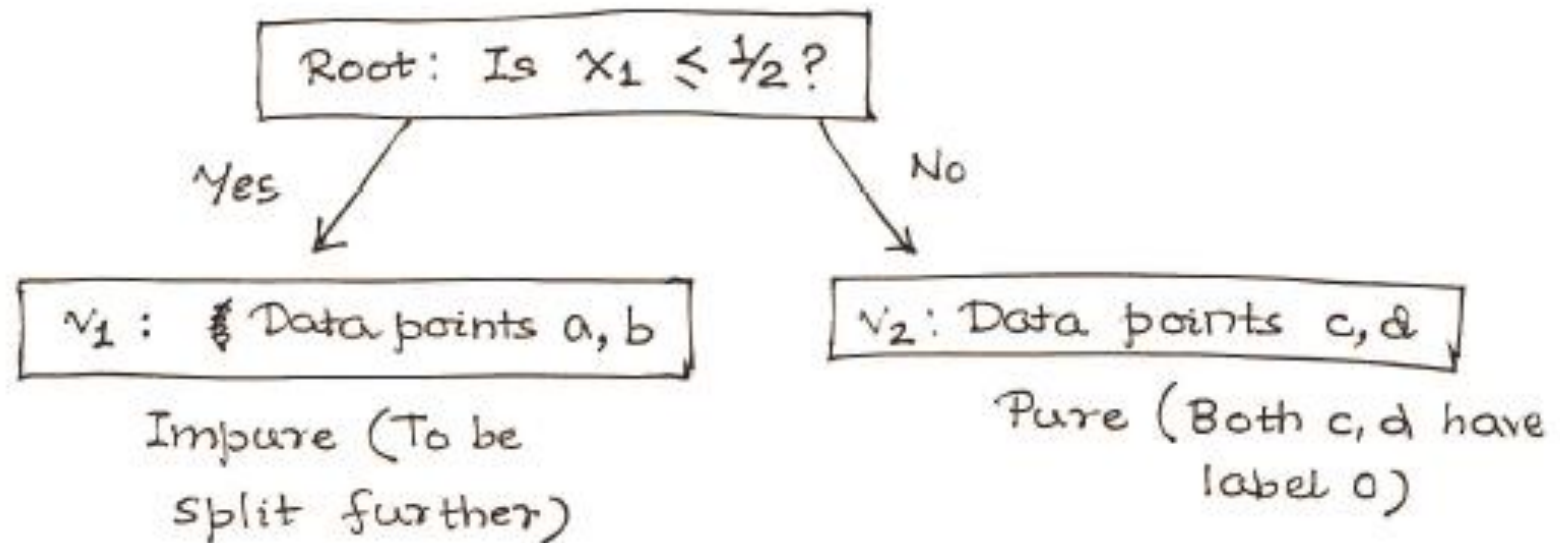
1. Initially: Root has (a,b,c,d)

Root: (a,b,c,d)

Impure.

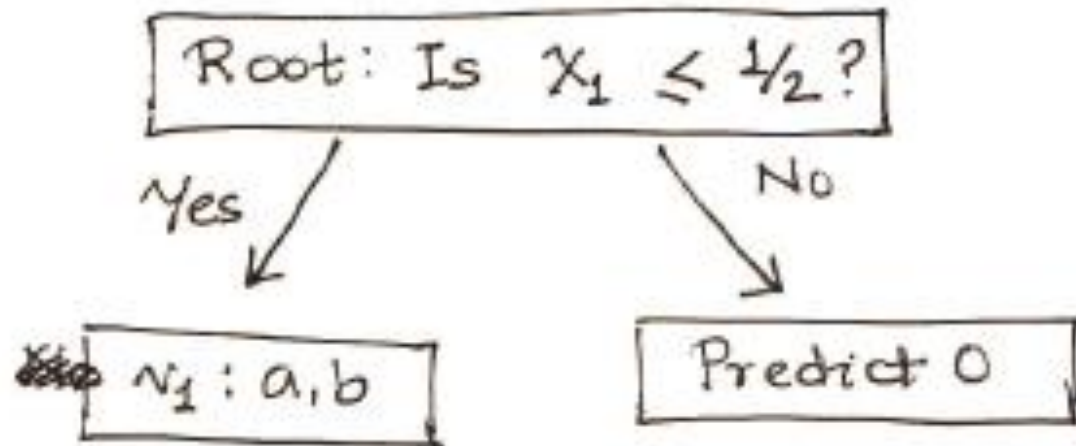
ID3 Algorithm: Running Example

2. Suppose we pick X_1 , threshold $\frac{1}{2}$ (thru Selection Rule).



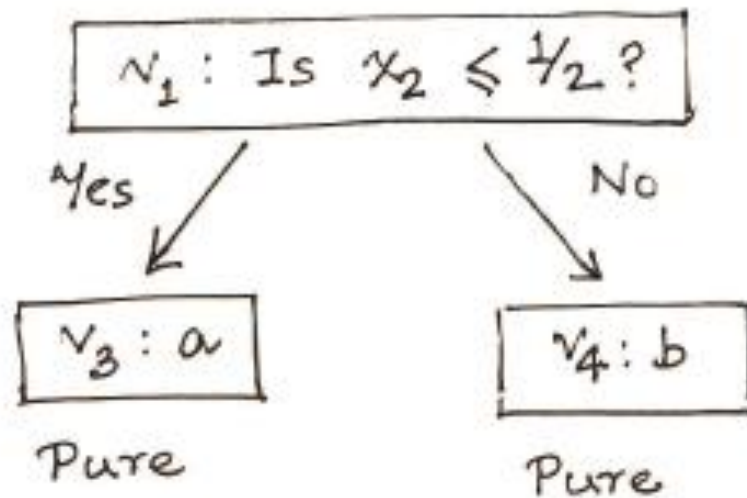
ID3 Algorithm: Running Example

3.



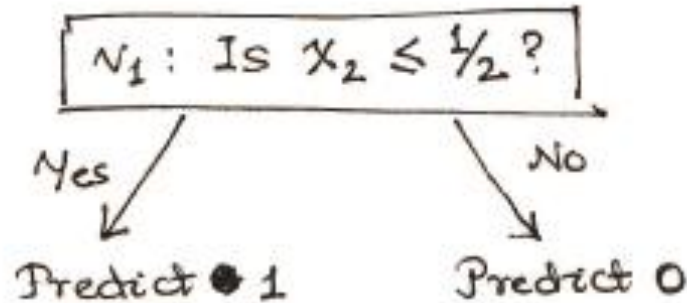
ID3 Algorithm: Running Example

4. Suppose we pick feature x_2 , threshold $\frac{1}{2}$ (to split v_1)



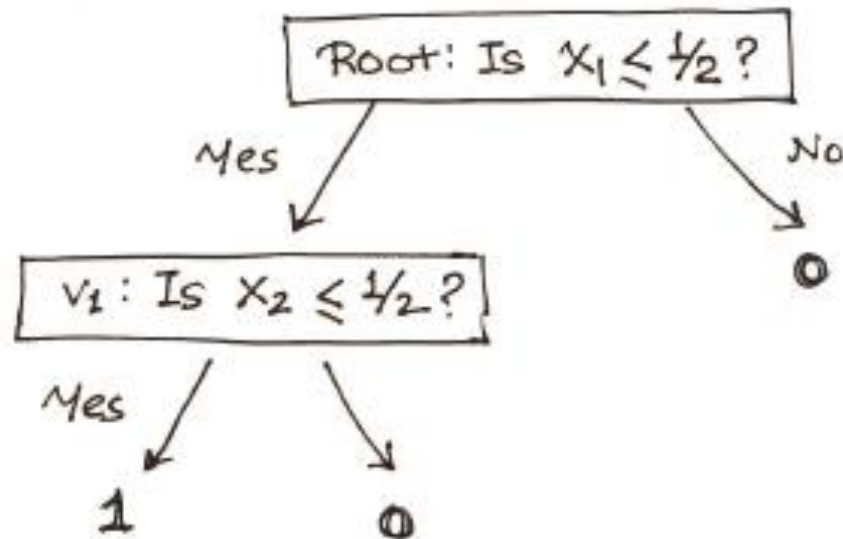
ID3 Algorithm: Running Example

5.



All nodes are pure at this point, so we are done.

Complete Decision Tree:



Splitting Rule

Splitting Rule: How to choose a feature and threshold along which to split a node?

In ID3 Decision Trees, we choose one that reduces uncertainty the most.

How to measure uncertainty? Through a notion in information theory, called entropy

Entropy: Let X be a random variable that takes values v_1, \dots, v_k with probabilities p_1, \dots, p_k . Then, entropy of X , denoted by $H(X)$ is defined as:

$$H(X) = -p_1 \log p_1 - p_2 \log p_2 - \dots - p_k \log p_k$$

(Note: Use the convention that $0 \cdot \log 0 = 0$)

(Optional) Entropy: Examples

Some Entropy Values:

1. X : 0/1 random variable (r.v). $\Pr(X=1) = 1/2$.

$$p_0 = \Pr(X=0) = 1/2$$

$$p_1 = \Pr(X=1) = 1/2$$

$$\begin{aligned} H(X) &= -p_0 \log p_0 - p_1 \log p_1 = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = \log 2 \\ &= 0.69 \quad (\text{Usually, we use natural logarithm for these calculations.}) \end{aligned}$$

2. X : 0/1 r.v. $\Pr(X=1) = 0$.

$$p_0 = 1, \quad p_1 = 0.$$

$$H(X) = -1 \log 1 - 0 \log 0 = 0$$

3. X is a r.v. which takes values $1, 2, \dots, k$ w.p. $\frac{1}{k}$ each.

$$\begin{aligned} H(X) &= -\frac{1}{k} \log \frac{1}{k} - \frac{1}{k} \log \frac{1}{k} - \dots - \frac{1}{k} \log \frac{1}{k} \\ &= -k \cdot \frac{1}{k} \log \frac{1}{k} = \log k. \end{aligned}$$

For 0/1 r.v, closer $\Pr(X=0)$ is to $1/2$, higher is the entropy.

Entropy: Properties

Properties of Entropy:

1. $H(X)$ does not depend on the exact values taken by X , only on the probabilities of distinct values.

2. $H(X) \geq 0$ (always). Why?

$$H(X) = - \sum_{i=1}^k p_i \log p_i, \text{ where } \log p_i \leq 0, \text{ as } p_i \leq 1.$$

3. If X takes k values, $H(X) \leq \log k$. This maximum is achieved only when each value occurs w.p. $1/k$.

Conditional Entropy

Conditional Entropy:

Let X, Z be two r.v.s. The conditional entropy of X given Z is defined as:

$$H(X|Z) = \sum_z \Pr(Z=z) H(X|Z=z)$$

of X

(Intuitively, average entropy ^{of X} given that we know Z).

Conditional Entropy: Examples

Joint Distribution of X, Z :

$Z \backslash X$	0	1
0	$\frac{1}{3}$	$\frac{1}{3}$
1	$\frac{1}{3}$	0

$$\Pr(Z=0) = \frac{2}{3} \quad \Pr(Z=1) = \frac{1}{3}$$

~~For any x ,~~ For any x ,

$$\Pr(X=x | Z=0) = \frac{\Pr(X=x, Z=0)}{2/3}$$

$$\Pr(X=x | Z=1) = \frac{\Pr(X=x, Z=1)}{1/3}$$

So: Conditional distributions are:

$$X | Z=0 : \quad \Pr(X=0 | Z=0) = 1/2$$

$$\Pr(X=1 | Z=0) = 1/2$$

$$X | Z=1 : \quad \Pr(X=0 | Z=1) = 1$$

$$\Pr(X=1 | Z=1) = 0$$

$$H(X | Z=0) = -1/2 \log 1/2 - 1/2 \log 1/2 = \log 2$$

$$H(X | Z=1) = -1 \log 1 - 0 \log 0 = 0$$

$$\text{So: } H(X | Z) = \frac{2}{3} \cdot \log 2 + \frac{1}{3} \cdot 0 = \frac{2}{3} \log 2.$$

Conditional Entropy: Properties

1. Suppose $X = Z$. What is $H(X|Z)$?

$$H(X|Z) = \sum_z \Pr(Z=z) H(X|Z=z).$$

Now, given that $Z=z$, we know that $X=z$ w.p. 1.

$$\text{So, } H(X|Z=z) = 0.$$

$$\text{Thus, } H(X|Z) = \sum_z \Pr(Z=z) \cdot 0 = 0.$$

(This may hold even if $H(X)$ is very large!)

2. Suppose X, Z are independent. What is $H(X|Z)$?

Since X, Z are independent, $\Pr(X=x|Z=z)$ for any z and x ,
is equal to $\Pr(X=x)$.

i.e. $X|Z=z$ has exactly the same distribution as X .

$$\text{So, } H(X|Z=z) = H(X)$$

$$H(X|Z) = \sum_z \Pr(Z=z) H(X|Z=z) = H(X) \sum_z \Pr(Z=z) = H(X)$$

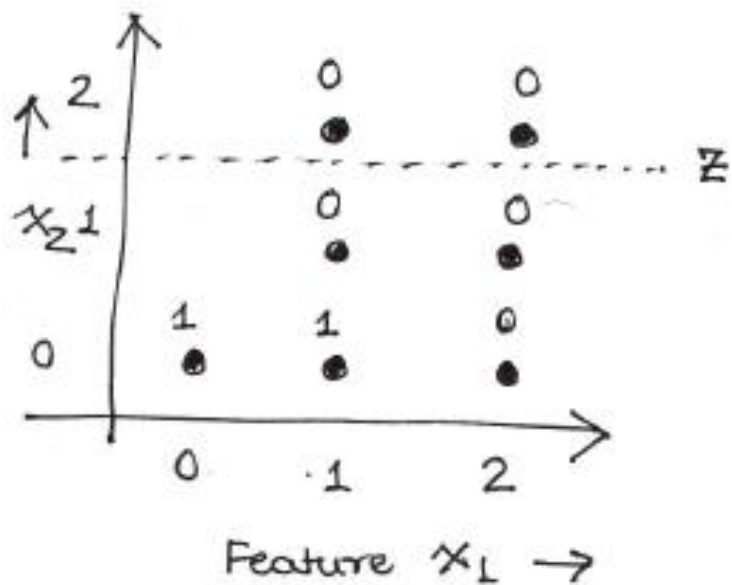
Information Gain

$$\text{Information Gain}(Z) = H(X) - H(X|Z)$$

(essentially, how much entropy of X is reduced because we know Z).

It can be shown that information gain (IG) is ≥ 0 (always)

Information Gain: Example



Suppose Figure shows the training data at a node v .

So, v has 5 label 0, 2 label 1 data points, and

$$H(X) = -\frac{5}{7} \log \frac{5}{7} - \frac{2}{7} \log \frac{2}{7}$$

Suppose we would like to split along Z . ($X_1 < 1.5$ or not).

Suppose $Z=0$ corresponds to $X_1 \geq 1.5$ (1)

$Z=1$ " " $X_1 < 1.5$ (2)

Then, $X|Z=0$ has 2 label 0s, 0 label 1s.

$$\Pr(X=0|Z=0) = 1, \quad \Pr(X=1|Z=0) = 0.$$

$$\text{So, } H(X|Z=0) = -1 \log 1 - 0 \log 0 = 0.$$

Information Gain: Example

$x|z=1$ has 3 label 0's and 2 label 1's.

$$\Pr(X=0|Z=1) = 3/5, \quad \Pr(X=1|Z=1) = 2/5.$$

$$H(X|Z=1) = -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{5}$$

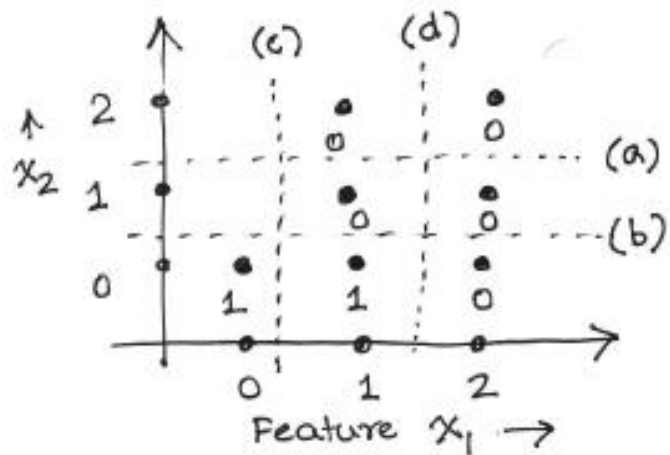
Also :

Also: $\Pr(Z=0) = 2/7$ (as 2 out of the total 7 points in X will lie on this branch)

and: $\Pr(Z=1) = 5/7$ (as 5 is not -1 and 7 is not -1)

$$\text{so: } H(X|Z) = P_r(Z=0) H(X|Z=0) + P_r(Z=1) H(X|Z=1) \\ = 0.48$$

How to Split a node: Example 1



4 possible (feature, threshold) pairs to split along.

Call them (a), (b), (c), (d).

Initially: (5 0, 2 1)

(5 label 0's, 2 label 1's)

(a): Results: (2 0, 0 1); (3 0, 2 1)

$$H(X|Z) = 0.48$$

(b): Results: (4 0, 0 1); (1 0, 2 1)

$$H(X|Z) = 0.27$$

(c): Results: (5 0, 1 1); (0 0, 1 1)

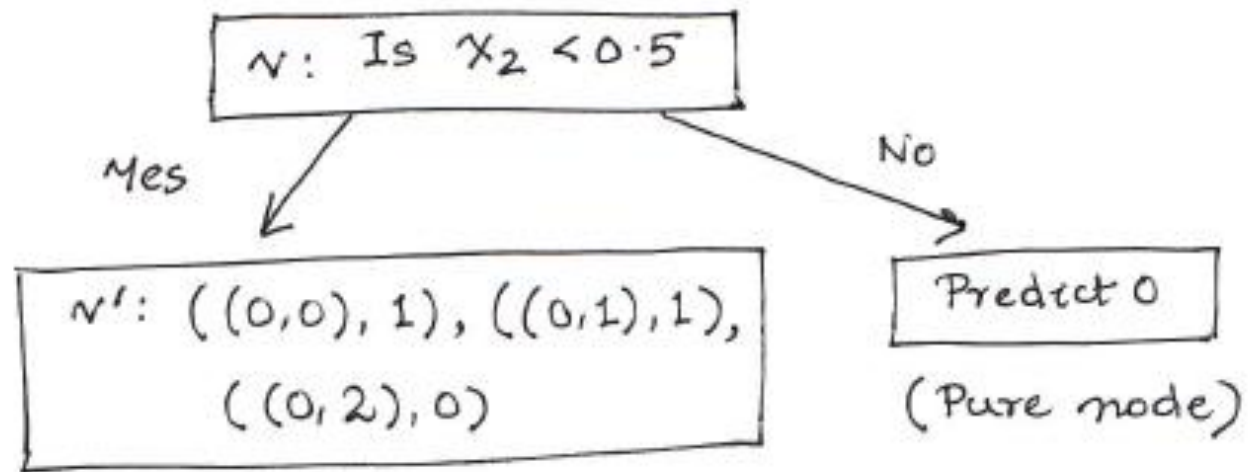
$$H(X|Z) = 0.38$$

(d): Results: (2 0, 2 1); (3 0, 0 1)

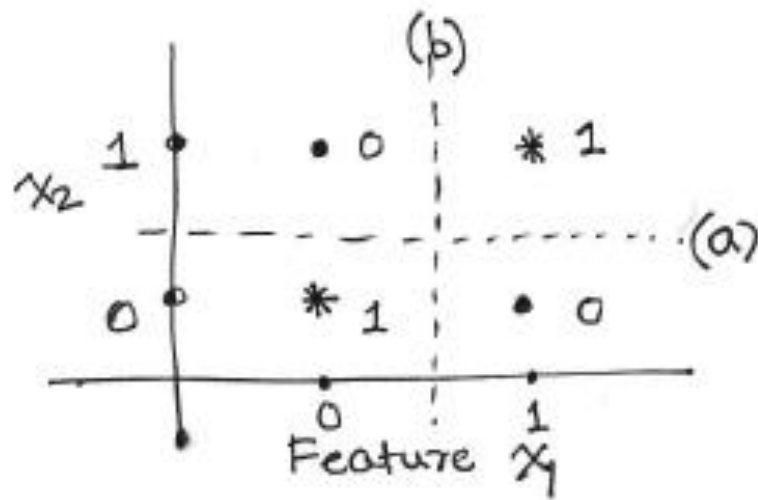
$$H(X|Z) = 0.39$$

How to Split a node: Example 1

So, ID3 will select (b) to get:



How to Split a node: Example 2



Initially; $\bullet(2 \ 0, \ 2 \ 1)$

$$H(X) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = \log 2 = 0.69.$$

2 possible ways to split: (a) and (b).

For (a); Result: $(1 \ 0, \ 1 \ 1), \ (1 \ 0, \ 1 \ 1)$

$$H(X|Z) = 0.69$$

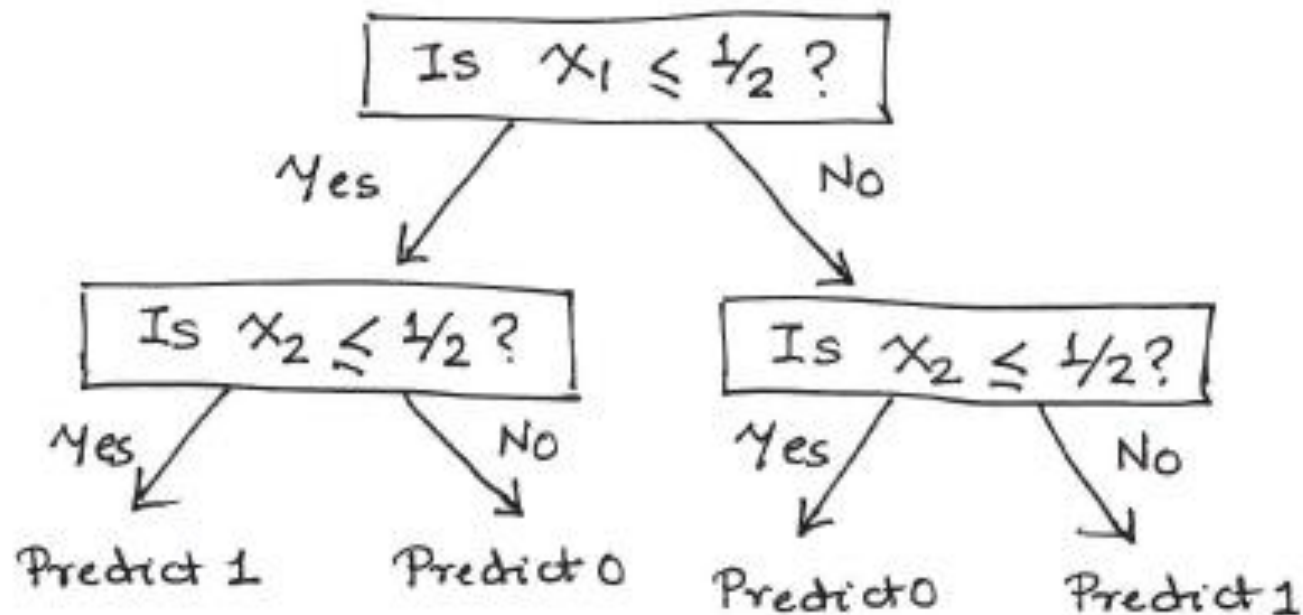
For (b): Result: $(1 \ 0, \ 1 \ 1) ; \ (1 \ 0, \ 1 \ 1)$

$$H(X|Z) = 0.69$$

So, $IG(Z) = 0.$

How to Split a node: Example 2

However, we shouldn't stop constructing the tree even if $IG = 0$. In this case, correct tree is:



Over-fitting: Example

Stopping Rule : Stop when every node is pure.

But this rule may overfit the data

Example:

x	x	x	x
x	x	o	x
x	x	x	x

The o may be an outlier on some noisy data.

If we build a decision tree that stops when all nodes are pure, then we may overfit the training data.

o	o	o
o	o	o

What is Over-fitting?

Data comes from some underlying true distribution \mathcal{D} .

Training data, test data \rightarrow all samples from \mathcal{D} .

Training error of a classifier:

$$\hat{err}(h) = \frac{1}{n} \sum_{i=1}^n 1(h(x_i) \neq y_i)$$

$1(\cdot)$ means an indicator variable which is 1 when the condition is true, 0 o/w.

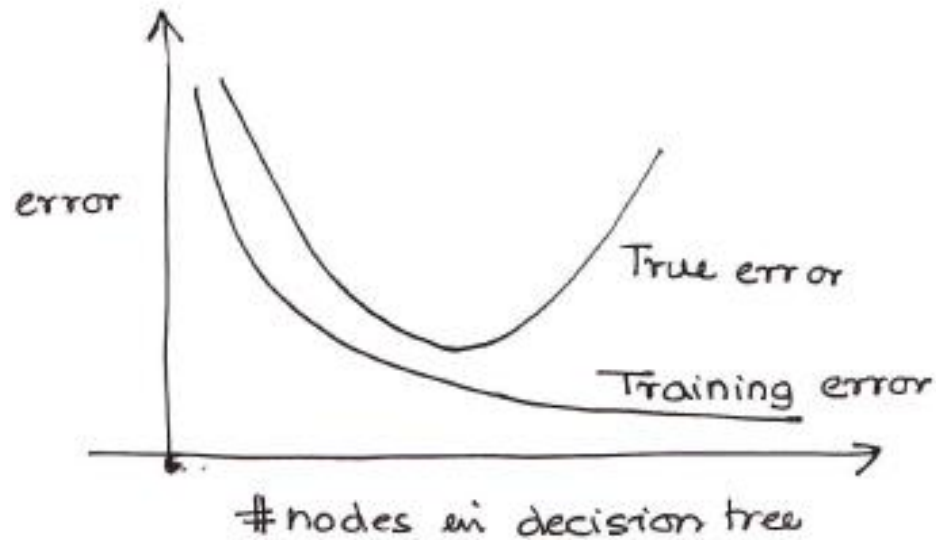
True error of a classifier:

$$err(h) = \Pr_{(x,y) \sim \mathcal{D}}(h(x) \neq y)$$

For a fixed classifier, with more data, training error should approach true error.

Over-fitting in Decision Trees

Overfitting happens when:



As we make the tree more and more complicated, training error decreases. At some point, true error stops improving and may get worse.

True data distribution has some structure, which we would like to capture. Training data captures some of this structure, but may have some "chance" structure of its own, which we should not model into a classifier.

Pruning to avoid Over-fitting in Decision Trees

1. Split the training ~~set~~ data into training set S and validation set V .
2. Build ID3 tree T using training set S
3. Prune using V :

Repeat:

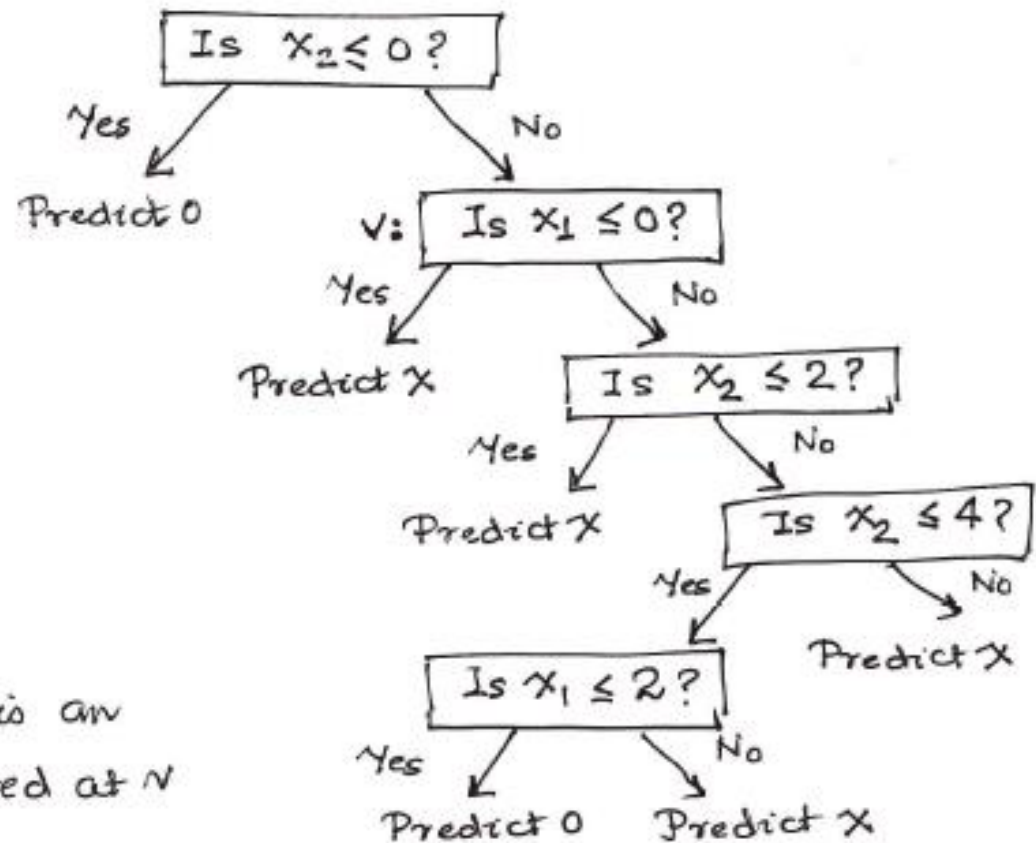
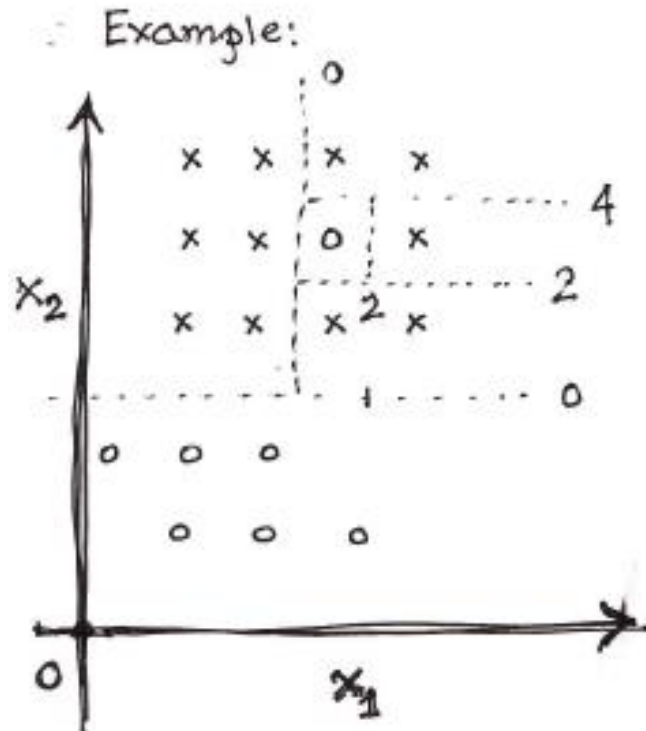
For each node v in T :

Replace subtree rooted at v by single node that predicts majority label in v to get tree T'

If $\text{error}(T') \text{ on } V \leq \text{error}(T) \text{ on } V$, then $T = T'$

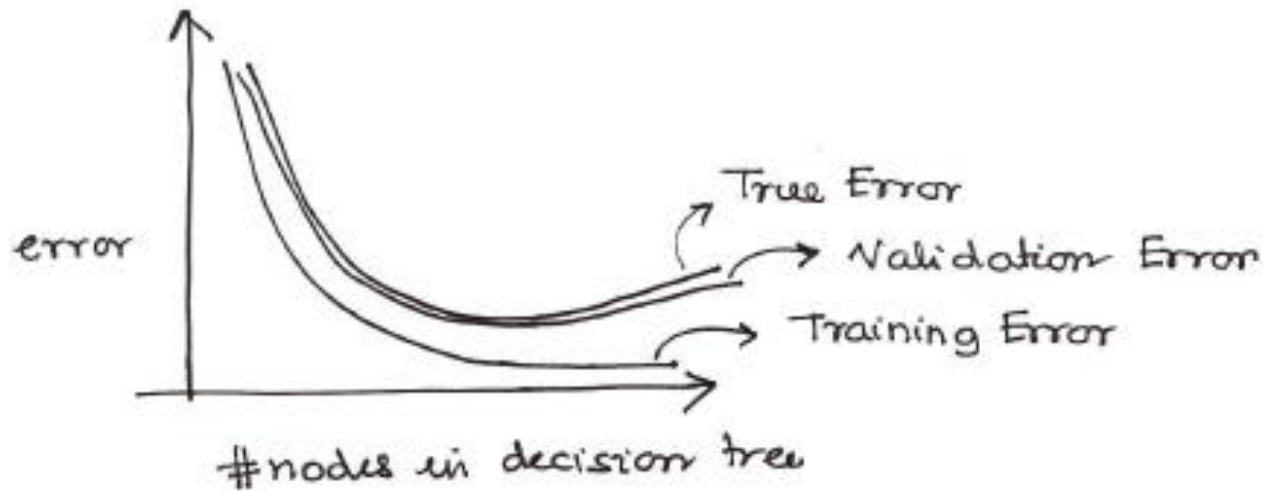
Continue till there is no such node v .

Pruning: Example



If the 0 point in between is an error, then the subtree rooted at v will have higher validation error than predicting x at v .

Pruning



As tree grows more complicated, training error decreases.
At some point, true error stops improving and may even get worse. This will be reflected in the validation error, provided we have enough validation data.

Bagging for Constructing Ensembles

What is Ensemble Learning?

In ensemble learning, the idea is to combine multiple classifiers into a single one. Ensemble learning usually works very well in practice.

Two methods (for this class):

1. Bagging
2. Boosting

Bagging

BAGGING: (Bootstrap AGGREGATING)

1. Input: n labelled training examples $(x_1, y_1), \dots, (x_n, y_n)$

2. Algorithm:

Repeat k times:

(a) Select m samples out of n with replacement from the training set to get training set S_i

(b) Train classifier h_i on S_i (usually, h_i 's are the same type of classifier)

3. Output: Classifiers h_1, \dots, h_k

Testing: Given test example x , output the majority of $h_1(x), h_2(x), \dots, h_k(x)$ (break ties at random as usual)

Choice Points in Bagging

1. How to pick k ?

Higher k is better, but also increases training time, storage requirement and classification time. So pick a k which is feasible.

Choice Points in Bagging

2. How to pick m ?

Popular choice for $m = n$.

But this is still very different from working with the entire training set!

$$\Pr(S_i = S) = \frac{n!}{n^n} \quad \left(\begin{array}{l} \text{\# ways of choosing } n \text{ samples} \\ \text{with replacement} = n^n \end{array} \right)$$

\Rightarrow Only $n!$ of these ways give you the entire training set!

$\Rightarrow \frac{n!}{n^n} \approx$ a very tiny number $\ll 2^{-n/2}$

For any (x_j, y_j) , $\Pr((x_j, y_j) \text{ is not in } S_i) = \left(1 - \frac{1}{n}\right)^n \approx \frac{1}{e}$ (for large n)
 $1/e \approx 0.37$; so about 37% of the data is left out of S_i .