

THEORY

OF

DECISION TREE

This page intentionally left blank

Decision Tree

Function Approximation

credits: Tom Mitchell

(3)

I. Problem Setting:

- Set of possible instances X [feature vector]
- Unknown target function f . [discrete-valued for decision tree]
- Set of function hypotheses: $H = \{h \mid h: X \rightarrow Y\}$ [set of all possible decision trees]

II. Input:

- Training examples $\{\langle x^{(i)}, y^{(i)} \rangle\}$ of unknown target function f $\rightarrow i^{th}$ example

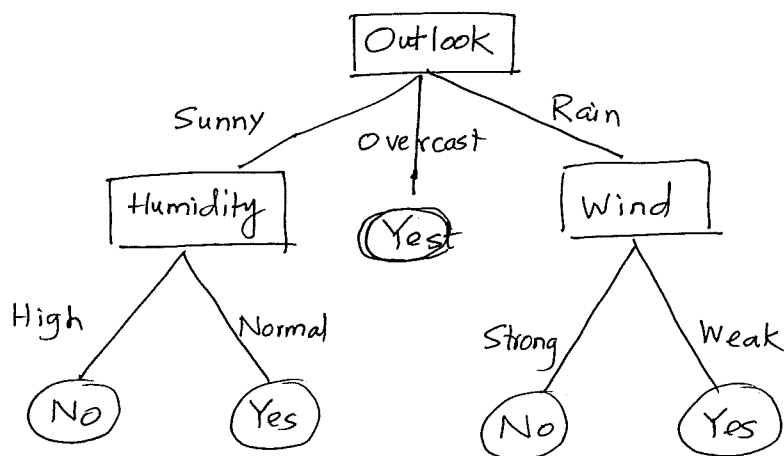
III. Output:

- Hypothesis $h \in H$ that best approximates f .

Example:

	X				Y
Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
	Sunny - 0	Hot - 0	High - 0	Strong - 0	Yes - 1
	Overcast - 1	Mild - 1	Normal - 1	Weak - 1	No - 0
	Rain - 2	Cold - 2			
Sample (feature-label combination)					
1	0	1	0	1	0
2	1	0	1	1	1
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

A Decision Tree for $f: \langle \text{Outlook, Temperature, Humidity, Wind} \rangle \rightarrow \text{PlayTennis}$



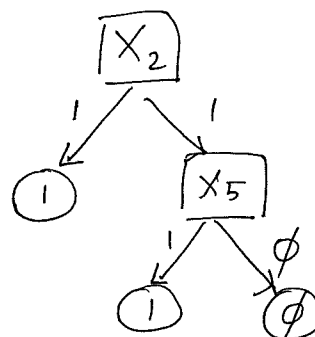
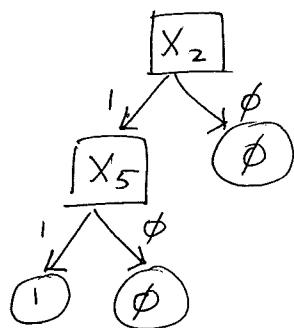
0 = Sunny
 T = ~~Hot~~ Hot
 H = Normal
 Wind = Weak / Strong

- Internal node: test one attribute X_i
- Branch from node: select one value for X_i
- Leaf node: predicts Y

Example Decision Trees

Suppose $X = \langle X_1, \dots, X_n \rangle$ $X_i \in \{0, 1\}$ boolean

Q.1. How would you represent $Y = X_2 \cdot X_5$ and $Y = X_2 \vee X_5$

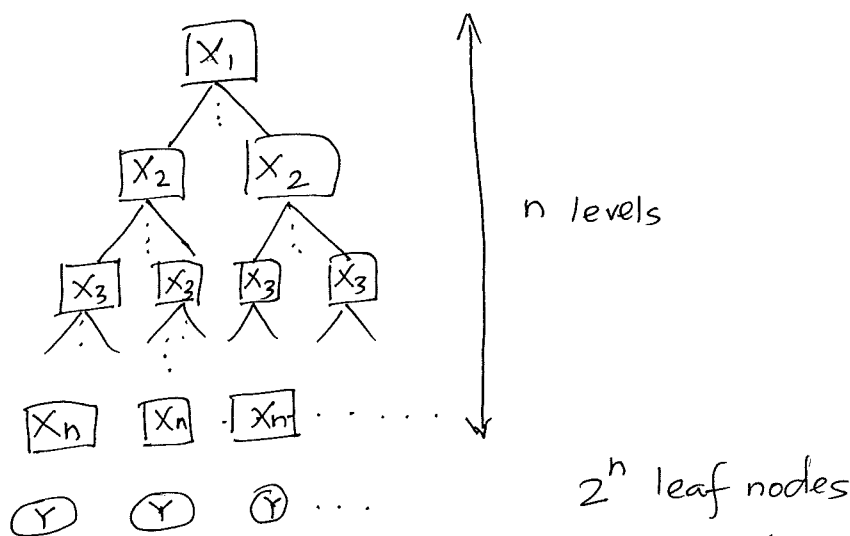


• \equiv AND
 $\vee \equiv$ OR
 $\neg \equiv$ NOT

or a more complicated one $Y = X_2 X_5 \vee X_3 X_4 (\neg X_1)$
 (or discrete-valued?)

Q.2. Can we represent arbitrary boolean functions using decision trees?

Yes!



(discrete-valued)

— so decision trees are expressive (can ~~be~~ represent any boolean function)

— This makes decision trees universal function approximator.

Top-Down Decision Tree Learning : ID3

node = Root

[ID3, C4.5, Quinlan]

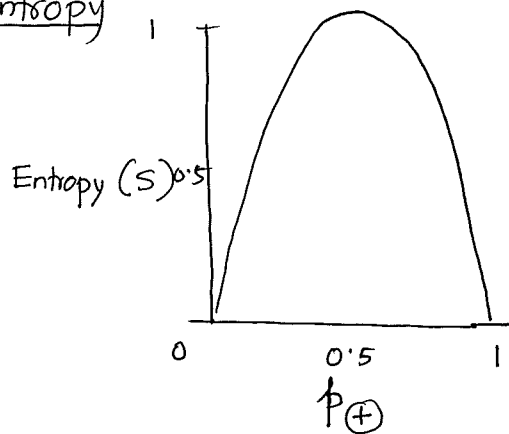
Main Loop:

1. $A \leftarrow$ the "best" decision attribute for next node
2. Assign A as decision attribute for node
3. For each value of A , create new descendant of node
4. Sort training examples to leaf nodes
5. If training examples are perfectly classified, STOP.
else iterate over new leaf nodes

- Top-down greedy growth of decision tree using "best" attribute until all examples are perfectly classified

- Which attribute is best?

Sample Entropy



- S is a sample of training examples
- p_+ is proportion of +ve examples
- p_- is proportion of -ve examples
- Entropy measures impurity of S

$$H(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

all +ve examples \rightarrow entropy is 0

all -ve examples \rightarrow entropy is 0

$p_+ = p_- \rightarrow$ entropy is 1

EntropyEntropy $H(X)$ of r.v. X

$$H(X) = - \sum_{i=1}^n P(X=i) \log_2 P(X=i)$$

 $H(X)$ is based on information theory:

- Most efficient possible code assigns $-\log_2 P(X=i)$ bits to encode the message $X=i$
- So, expected number of bits to code one random X is

$$\sum_{i=1}^n P(X=i) (-\log_2 P(X=i))$$

Some Notations

Specific conditional entropy $H(X|Y=v)$ of X given $Y=v$:

$\left[\begin{array}{l} \text{- How much } H \text{ is reduced if we sort data by attributes} \\ \text{- choose the attribute which reduces entropy the most} \end{array} \right]$

$$H(X|Y=v) = - \sum_{i=1}^n P(X=i|Y=v) \log_2 P(X=i|Y=v)$$

Conditional entropy $H(X|Y)$ of X given Y :

$$H(X|Y) = \sum_{v \in \text{values}(Y)} P(Y=v) H(X|Y=v)$$

Mutual Information (a.k.a. information gain) of X and Y :

$$I(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

Information gain is the expected reduction in entropy of target variable Y for data sample S , due to sorting on variable A

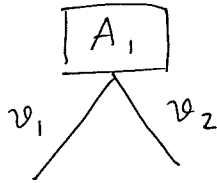
$$\text{Gain}(S, A) = I_S(A, Y) = H_S(Y) - H_S(Y|A)$$

choose the attribute with highest Gain

Selecting the "best" attribute

$$S : [x_1+, x_2-]$$

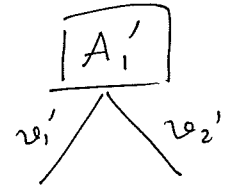
Sample : $E_\phi = ?$
Entropy



$$S_1 : [x_3+, x_4-] \quad S_2 : [x_5+, x_6-]$$

Conditional Entropy : $E_1 = ?$ $E_2 = ?$

$$\begin{aligned} \text{Gain}(S, A_1) &= G \\ &= E_\phi - \frac{|S_1|}{|S|} \cdot E_1 - \frac{|S_2|}{|S|} \cdot E_2 \end{aligned}$$



$$S'_1 : [x'_3+, x'_4-] \quad S'_2 : [x'_5+, x'_6-]$$

$$E'_1 = ? \quad E'_2 = ?$$

$$\begin{aligned} \text{Gain}(S', A'_1) &= G' \\ &= E_\phi - \frac{|S'_1|}{|S'|} E'_1 - \frac{|S'_2|}{|S'|} E'_2 \end{aligned}$$

- Find whether G is greater than G'
- If so, choose A_1 over A'_1 , else : choose A'_1

Q. Is there more than one DT that will perfectly sort the data? If yes, which one do you choose?

A. Function Approximation : Big Picture

$$f : X \rightarrow Y \quad X = \langle x_1, x_2, \dots, x_n \rangle \quad x_i \in \{0, 1\} \quad Y \in \{0, 1\}$$

$$H = \{h \mid h : X \rightarrow Y\}$$



Hypothesis Space

Instance space

$$|X| = 2^n$$

of DTs that can represent all possible functions = 2^{2^n}

$$\# \text{ possible functions} = 2^{|X|} = 2^{2^n}$$

sad fact about inductive inference

No Free Lunch!

training examples we need to label so that there is just one DT (unique) in the hypothesis space = All

PAC learning

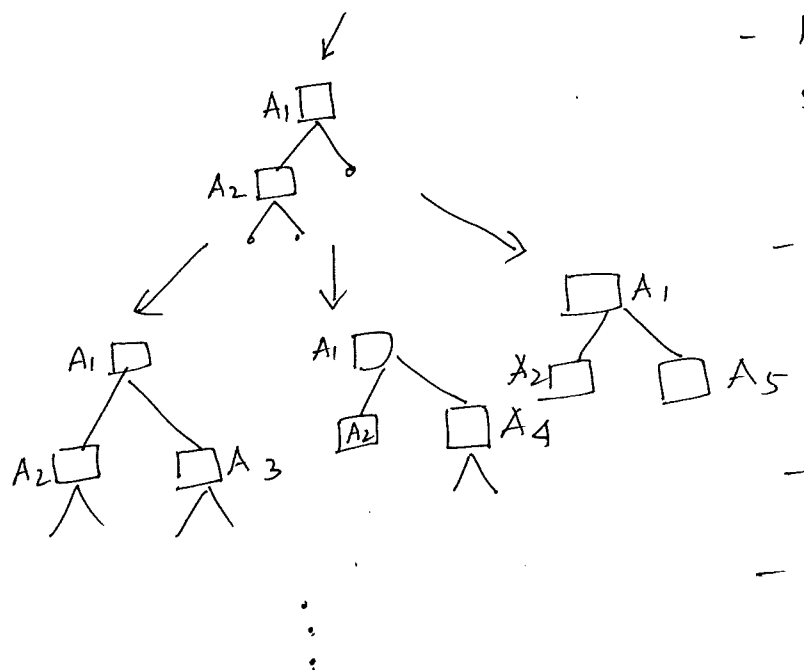
"All models are wrong" - George Box 1976. J. American Stat. Ass

In practical ML applications, we don't have all possible labeled training data. So, we make some assumptions during learning. (8)

Q. In decision tree learning algorithm (ID3), what's the assumption?

A. It stops at the smallest acceptable tree.

Function Approximation as Search for the best hypothesis



- ID3 performs heuristic search through the space of decision trees

- It stops at smallest acceptable tree. Why?

- William of Ockham ~1300

- Occam's razor: prefer the simplest hypothesis that fits the data

All models are wrong; some are useful

Measuring Accuracy of a trained model using a test data

- Accuracy/Success
- Error rate/failure

		True Condition	
		+	-
Predicted Condition	+	TP	FP
	-	FN	TN

$$\text{Precision} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \frac{TP}{TP + FN}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$F_1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

$$\text{MCC} = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Overfitting in Decision Tree

- Fits the training data very well, but does not perform very well on test data
- Overfitting can happen due to:
 - Noisy training example too.

Generally,

consider a hypothesis h and

- Error rate over training data : $\text{error}_{\text{train}}(h)$
- True error rate over all data : $\text{error}_{\text{true}}(h)$

h overfits the training data if

$$\text{error}_{\text{true}}(h) > \text{error}_{\text{train}}(h)$$

$$\text{Amount of overfitting} = \text{error}_{\text{true}}(h) - \text{error}_{\text{train}}(h)$$

Avoiding overfitting in Decision Tree : Reduced Error Pruning

Split data into training and validation set

Create tree that classifies training set correctly

Do until further pruning is harmful:

- 1) Evaluate impact on validation set of pruning each node (and below)
- 2) Greedily remove the one that most improves validation accuracy

Guaranteed to produce smallest version of most accurate subtree.