

Decision Trees and Boosting

CSci 5525: Machine Learning

Instructor: Nicholas Johnson

November 10, 2020

Announcements

- HW3 due tonight (11:59 PM CST)
- Project progress report due in 1 week (Nov 17)
- Exam 2 coming up (**Monday** Nov 23, due 48 hours later)
 - Covers lectures 11 (Deep Learning I) - 21 (tentatively PCA)

Attribute-based representations

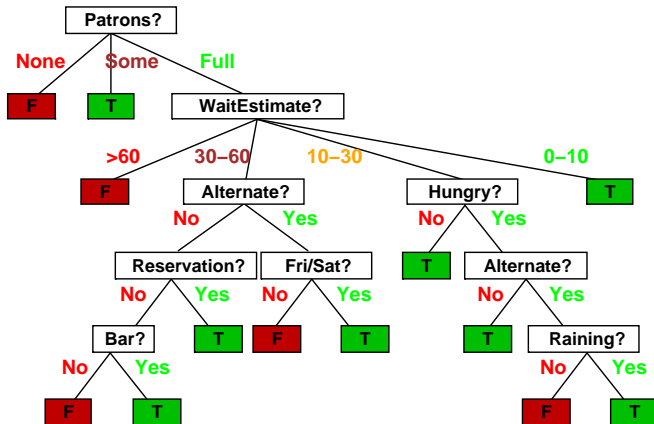
Examples described by attribute values (Boolean, discrete, continuous)

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

Classification of examples is positive (T) or negative (F)

Decision Tree Example (Restaurant)

One possible representation for hypotheses

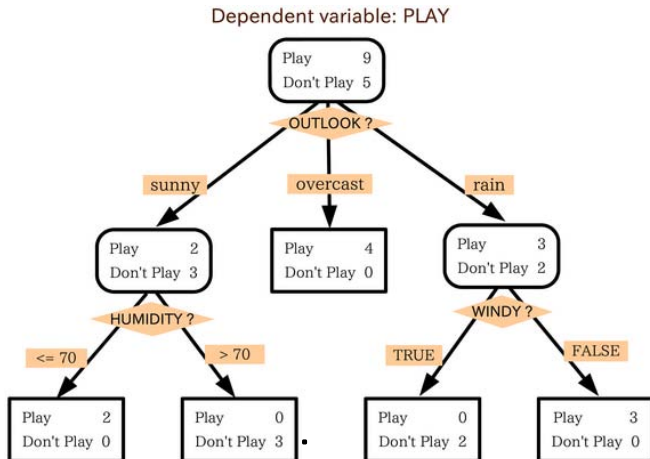


Example: Playing Golf

Play golf dataset

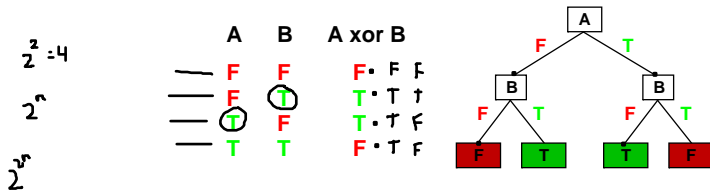
Independent variables					Dep. var
OUTLOOK	TEMPERATURE	HUMIDITY	WINDY	PLAY	
sunny	85	85	FALSE	Don't Play	•
sunny	80	90	TRUE	Don't Play	•
overcast	83	78	FALSE	Play	
rain	70	96	FALSE	Play	
rain	68	80	FALSE	Play	
rain	65	70	TRUE	Don't Play	
overcast	64	65	TRUE	Play	
sunny	72	95	FALSE	Don't Play	•
sunny	69	70	FALSE	Play	•
rain	75	80	FALSE	Play	
sunny	75	70	TRUE	Play	•
overcast	72	90	TRUE	Play	
overcast	81	75	FALSE	Play	
rain	71	80	TRUE	Don't Play	

Decision Tree Example (Golf)



Expressiveness

- Decision trees can express any function of the input attributes
 - For Boolean functions, truth table row \rightarrow path to leaf



- The basic trade-off
 - There is a consistent decision tree for any training set
 - Unless f is nondeterministic in x
 - One path to leaf for each example
 - But it probably won't generalize to new examples
- Prefer to find more compact decision trees

Hypothesis Spaces

- How many distinct decision trees with n Boolean attributes
 - Number of boolean functions
 - Number of distinct truth tables with 2^n rows $= 2^{2^n}$
 - 6 Boolean attributes give 18,446,744,073,709,551,616 trees
- How many purely conjunctive hypotheses ($Hungry \cap \neg Rain$)
 - Each attribute can be in (positive), in (negative), or out
 $\implies 3^n$ distinct conjunctive hypotheses
- More expressive hypothesis space
 - Increases chance that target function can be expressed
 - Increases number of hypotheses consistent w/ training set
 \implies may get worse predictions

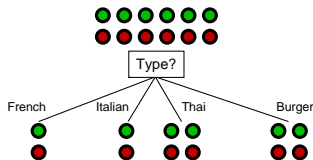
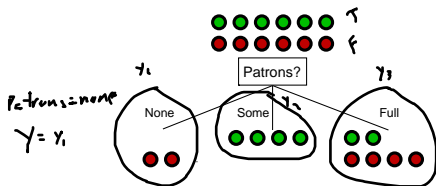
3^n

Decision Tree Learning

Aim: Find a small tree consistent with the training examples

Recursively choose “most significant” attribute as root of (sub)tree

Good attribute splits the examples (ideally) into “pure” subsets



Patrons? is a better choice

—gives *information* about the classification

Information Gain

- $$X = \begin{matrix} x_1 \\ x_2 \end{matrix} \quad P(X_i) = \begin{matrix} \frac{1}{2} \\ \frac{1}{2} \end{matrix}$$

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

- $$X = \begin{pmatrix} x_1 \\ -x_2 \\ \vdots \\ x_n \end{pmatrix}$$

Return

- now
same
full

$$IG(X|Y) = H(X) - H(X|Y)$$

- Navigation icons: back, forward, search, and other controls.

Information Gain (Contd.)

- An attribute splits the examples E into subsets E_i
 - Each E_i should need less information to classify
- Let E_i have p_i positive and n_i negative examples
 - $H(\langle p_i/(p_i + n_i), n_i/(p_i + n_i) \rangle)$ bits needed to classify
 - Expected number of bits per example over all branches is

$$H(X|Y) = \sum_i \frac{p_i + n_i}{p + n} H(\langle p_i/(p_i + n_i), n_i/(p_i + n_i) \rangle)$$

C.E. $\xrightarrow{\hspace{2cm}}$

- For 'Patrons', this is 0.459 bits, for 'Type' this is (still) 1 bit
- Choose the attribute that maximizes information gain

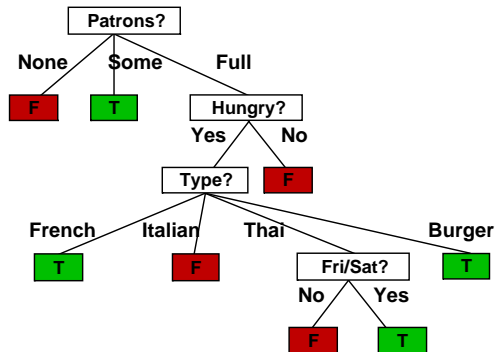
I 6

Patrons: $1 - 0.459 > 0$

Type: $1 - 1 = 0$

A Simple Decision Tree

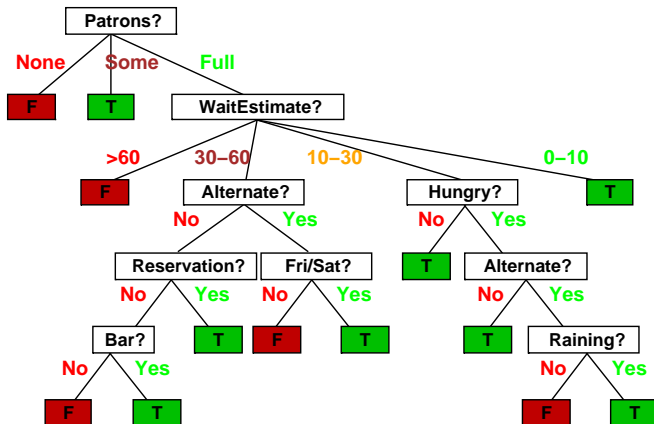
- Decision tree learned from the 12 examples:



- Substantially simpler than “true” tree
 - More complex hypothesis is not justified by a small dataset

The Original Tree

The original (complex) tree with same performance on the training set



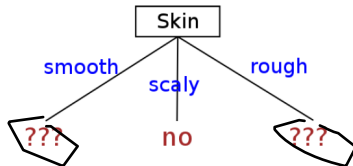
Another Example: Edibility

target

Skin	Color	Thorny?	Flowering?	Edible?
smooth	pink	no	yes	yes
smooth	pink	no	no	yes
scaly	pink	no	yes	no
rough	purple	no	yes	no
rough	orange	yes	yes	no
scaly	orange	yes	no	no
smooth	purple	no	yes	yes
smooth	orange	yes	yes	no
rough	purple	yes	yes	no
smooth	purple	yes	no	no
scaly	purple	no	no	no
scaly	pink	yes	yes	no
rough	purple	no	no	yes
rough	orange	yes	no	yes

Edibility Example (Contd.)

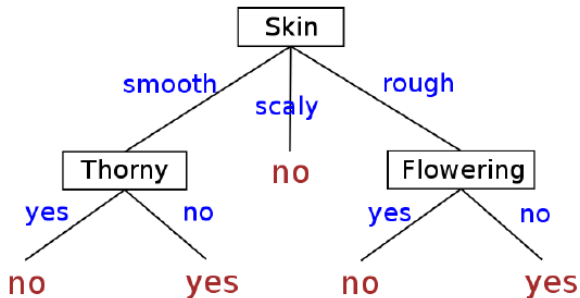
- Entropy $H(\text{Edible}) = 0.9403$
- Information Gain for each attribute
 - $IG(\text{Edible}|\text{Skin}) = 0.2467$ —
 - $IG(\text{Edible}|\text{Color}) = 0.0292$
 - $IG(\text{Edible}|\text{Thorny}) = 0.1519$
 - $IG(\text{Edible}|\text{Flowering}) = 0.0481$



Edibility Example (Contd.)

- Consider $Skin = Smooth$
 - Entropy $H(Edible|Skin = smooth) = 0.9710$
 - $IG(Edible|Color, Skin = smooth) = 0.4000$
 - $IG(Edible|Thorny, Skin = smooth) = 0.9710$
 - $IG(Edible|Flowering, Skin = smooth) = 0.0200$
- Consider $Skin = rough$
 - Choose $Flowering$

Edibility Decision Tree



Other Methods for Feature Selection

- Issues with Information Gain
 - Attributes that can take many values (e.g. customer ID)
 - High information gain but does not generalize well
- Gain Ratio

$$\text{GainRatio}(X|Y) = \frac{IG(X|Y)}{H(Y)}$$

↑

- Gini Index

$$\text{Gini}(X) = \sum_{i \neq j} p(i)p(j)$$

$$\text{Gini}(X|Y) = \sum_j p(y_j) \text{Gini}(X|y_j)$$

$$\text{GiniGain}(X|Y) = \text{Gini}(X) - \text{Gini}(X|Y)$$

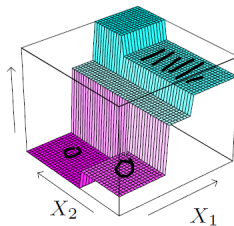
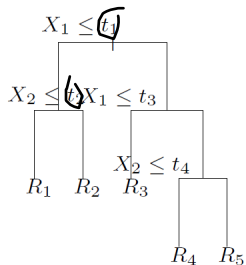
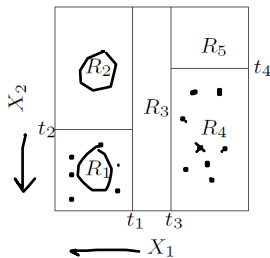
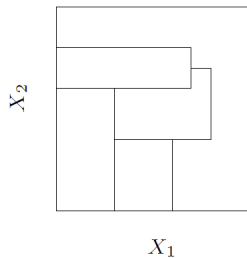
Challenges: Overfitting, Continuous features, Stability

- Overfitting: May generate a large tree
 - Good training set performance, poor test set performance
 - Decision tree pruning: Remove irrelevant attributes
 - Information gain is small after splitting
 - Significance test with null hypothesis of no underlying pattern
- Continuous valued attributes
 - Split into regions, convert to categorical variable
 - Find the split(s) so as to maximize information gain

Temperature	40	48	60	72	80	90
Play Tennis	No	No	Yes	Yes	Yes	No

- Stability
 - Mild change in attributes can change the tree
 - High variance, unstable, but interpretable
 - Ensembles approaches to reduce variance

Regression Trees



Regression Trees

- Assume regions R_m with 'constant' values
- The regression model is given by

$$f(\mathbf{x}) = \sum_{m=1}^M \hat{c}_m \mathbb{I}[\mathbf{x} \in R_m]$$

- Using criterion: $\min \sum (y_i - f(\mathbf{x}_i))^2$ we have

$$\hat{c}_m = \text{average}(y_i | \mathbf{x}_i \in R_m) = \frac{1}{N_m} \sum_{i: \mathbf{x}_i \in R_m} y_i$$

- Difficult to find general regions (in high-d)
- Difficult to even find best binary partition
- Greedy approach: For each feature j and split point s

$$R_1(j, s) = \{\mathbf{X} | X_j \leq s\} \quad R_2(j, s) = \{\mathbf{X} | X_j > s\}$$

- The splitting variable and split point is chosen by solving

$$\rightarrow \min_{j, s} \left[\min_{c_1} \sum_{\mathbf{x}_i \in R_1(j, s)} (y_i - c_1)^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j, s)} (y_i - c_2)^2 \right]$$

Weak Learning

- Often we stop decision tree learning early
 - After 1 level - “decision stump”
- Short decision trees are weak learners
- A weak learner predicts “slightly better” than random
- The PAC setting
 - Let \mathcal{X} be an instance space, $c : \mathcal{X} \mapsto \{0, 1\}$ be a target concept, \mathcal{H} be a hypothesis space ($h : \mathcal{X} \mapsto \{0, 1\}$)
 - Let D be a fixed (but unknown) distribution on \mathcal{X}
- An algorithm, after training on $(x_i, c(x_i)), [i]_1^m$, selects $(h) \in \mathcal{H}$ such that

$$P_{x \sim D}[h(x) \neq c(x)] \leq \frac{1}{2} - \gamma$$

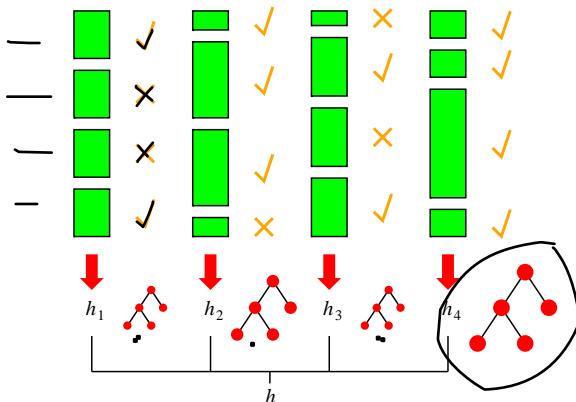
- The algorithm is called a γ -weak learner
- We assume the existence of such a learner



The Boosting Model

- Boosting converts a weak learner to a strong learner
- Boosting proceeds in rounds
 - Booster constructs D_t on X , the train set
 - Weak learner produces a hypothesis $h_t \in \mathcal{H}$ so that
$$P_{x \sim D_t}[h_t(x) \neq c(x)] \leq \frac{1}{2} - \gamma_t$$
 - After T rounds, the weak hypotheses $h_t, [t]_1^T$ are combined into a final hypothesis h_{final}
- We need procedures
 - for obtaining D_t at each step
 - for combining the weak hypotheses

Boosting Algorithms



- Weight decreased on correct samples
- Weight increased on incorrect samples