

Furong Huang / furongh@cs.umd.edu



Last week: introducing machine learning

What does "learning by example" mean?

Classification tasks

Learning requires examples + inductive bias Generalization vs. memorization

Formalizing the learning problem Function approximation Learning as minimizing expected loss0poo





Consider systems that apply a function f() to input items x and return outputs y = f(x).





In (supervised) machine learning, we deal with systems whose $f(\mathbf{x})$ is **learned from** examples.





We typically use machine learning when the function $f(\mathbf{x})$ we want the system to apply is unknown to us, and we cannot "think" about it.







Supervised Learning: Training



Given the training examples in $\mathcal{D}^{\text{train}}$ The learner returns a model $g(\mathbf{x})$



Supervised Learning: Testing





Supervised Learning: Testing





Supervised Learning: Testing



Apply the model to the raw test data?
Evaluate by the failing predicted labels against the test labels



Learning formally

- Given: examples $(\mathbf{x}, f(\mathbf{x}))$ of unknown function f
- Find: A good approximation of f
- x provides some representation of the input
 - Feature extraction: the process of mapping a domain element into a representation.
 - > $\mathbf{x} \in \{0,1\}^n, \, \mathbf{x} \in \mathbb{R}^n$
- The target function f() (label)
 - ► $f(\mathbf{x}) \in \{-1, +1\}$ Binary classification
 - ► $f(\mathbf{x}) \in \{1, 2, 3, ..., k 1\}$ Multi-class classification
 - ► $f(\mathbf{x}) \in \mathbb{R}$

Regression



Learning formally - continued

Hypothesis space

- > Set of possible instances $X = \{x \in \mathcal{X}\}$
- Set of possible outputs $Y = \{y \in \boldsymbol{\mathcal{Y}}\}$
- The set of function hypotheses that provide some mapping from the input to output space $H = \{h \mid h: \mathbf{x} \to y, \mathbf{x} \in \mathcal{X}, y \in \mathcal{Y}\}$
- Size of hypothesis space



Machine Learning as Function Approximation

Problem setting

- Set of possible instances $X = \{x \in \mathcal{X}\}$
- ► Unknown target function $f: Y = \{y \in \boldsymbol{\mathcal{Y}}\}$
- ▶ Set of function hypotheses $H = \{h \mid h: \mathbf{x} \to \mathbf{y}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}\}$
- Input
 - Training examples {(x⁽¹⁾, y⁽¹⁾), ... (x^(N), y^(N))} of unknown target function f
- Output
 - Hypothesis $h \in H$ that best approximates target function f



Formalizing induction: Loss Function

l(y, f(x)) where y is the truth and f(x) is the system's prediction

e.g.
$$l(y, f(x)) = \begin{cases} 0 & if \ y = f(x) \\ 1 & otherwise \end{cases}$$

Captures our notion of what is important to learn



Formalizing induction: Data generating distribution

Where does the data come from?

- Data generating distribution
 A probability distribution *D* over (*x*, *y*) pairs
- We don't know what D is!

We only get a random sample from it: our training data



Formalizing induction: Expected loss

- *f* should make good predictions
 - as measured by loss l
 - on **future** examples that are also drawn from D
- Formally
 - ε , the expected loss of f over D with respect to l should be small

$$\varepsilon \triangleq \mathbb{E}_{(x,y)\sim D}\{l(y,f(x))\} = \sum_{(x,y)} D(x,y)l(y,f(x))$$



Formalizing induction: Training error

- We can't compute expected loss because we don't know what D is
- We only have a sample of D
 - training examples $\{(x^{(1)}, y^{(1)}), ..., (x^{(N)}, y^{(N)})\}$
- All we can compute is the training error

$$\hat{\varepsilon} \triangleq \sum_{n=1}^{N} \frac{1}{N} l(y^{(n)}, f(x^{(n)}))$$



Formalizing Induction

Given

- a loss function *l*
- a sample from some unknown data distribution *D*
- Our task is to compute a function f that has low expected error over *D* with respect to *l*.

$$\mathbb{E}_{(x,y)\sim D}\left\{l(y,f(x))\right\} = \sum_{(x,y)} D(x,y)l(y,f(x))$$



Today: Decision Trees

- What is a decision tree?
- How to learn a decision tree from data?
- What is the inductive bias?
- Generalization?



An example training set

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



A decision tree to decide whether to play tennis





Decision Trees

- Representation
 - Each internal node tests a feature



- Each branch corresponds to a feature value
- Each leaf node assigns a classification
 - or a probability distribution over classifications
- Decision trees represent functions that map examples in X to classes in Y
- f: <Outlook, Temperature, Humidity, Wind> \rightarrow PlayTennis?





How would you represent the following Boolean functions with decision trees?

$(A \cap B)$

 $(A \cup B)$

Take home exercise: $(A \cap B) \cup (C \cap \neg D)$



Today: Decision Trees

- What is a decision tree?
- How to learn a decision tree from data?
- What is the inductive bias?
- Generalization?



Function Approximation with Decision Trees

- Problem setting
 - Set of possible instances X
 - ♦ Each instance $x \in X$ is a feature vector $x = [x_1, ..., x_D]$
 - > Unknown target function $f: X \rightarrow Y$
 - ✤ Y is discrete valued
 - ▶ Set of function hypotheses $H = \{h \mid h: X \rightarrow Y\}$
 - Each hypothesis h is a decision tree
- Input
 - Training examples {(x⁽¹⁾, y⁽¹⁾), ... (x^(N), y^(N))} of unknown target function f
- Output
 - > Hypothesis $h \in H$ that best approximates target function f

Decision Trees Learning

- Finding the hypothesis $h \in H$
 - That minimizes training error
 - Or maximizes training accuracy
- How?
 - H is too large for exhaustive search!
 - We will use a heuristic search algorithm which
 - Picks questions to ask, in order
 - Such that classification accuracy is maximized

Top-down Induction of Decision Trees

CurrentNode = Root

DTtrain (examples for CurrentNode, features at CurrentNode):

- Find F, the "best" decision feature for next node
- For each value of F, create new descendant of node
- 3. Sort training examples to leaf nodes
- 4. If training examples perfectly classified Stop

Else

Recursively apply DTtrain over new leaf nodes



How to select the "best" feature?

 A good feature is a feature that lets us make correct classification decision

- One way to do this:
 - select features based on their classification accuracy

Let's try it on the PlayTennis dataset



Let's build a decision tree using features O, T, H, W

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



Partitioning examples according to Humidity feature

Day	Outlook	Temperature	Humidity	Wind	PlayTennis?	
D1	Sunny	Hot	High	Weak	No	
D2	Sunny	Hot	High	Strong	No	
D3	Overcast	Hot	High	Weak	Yes	
D4	Rain	Mild	High	Weak	Yes	
D5	Rain	Cool	Normal	Weak	Yes	
D6	Rain	Cool	Normal	Strong	No	
D7	Overcast	Cool	Normal	Strong	Yes	J
D8	Sunny	Mild	High	Weak	No	
D9	Sunny	Cool	Normal	Weak	Yes	
D10	Rain	Mild	Normal	Weak	Yes	
D11	Sunny	Mild	Normal	Strong	Yes	
D12	Overcast	Mild	High	Strong	Yes	
D13	Overcast	Hot	Normal	Weak	Yes	
D14	Rain	Mild	High	Strong	No	



Partitioning examples: H = Normal

	Day	Outlook	Temperature	Humidity	Wind	PlayTennis?	
	D1	Sunny	Hot	High	Weak	No	
	D2	Sunny	Hot	High	Strong	No	
	D3	Overcast	Hot	High	Weak	Yes	
	D4	Rain	Mild	High	Weak	Yes	
	D5	Rain	Cool	Normal	Weak	Yes	
	D6	Rain	Cool	Normal	Strong	No	
	D7	Overcast	Cool	Normal	Strong	Yes	ノ
	D8	Sunny	Mild	High	Weak	No	
	D9	Sunny	Cool	Normal	Weak	Yes	
	D10	Rain	Mild	Normal	Weak	Yes	
	D11	Sunny	Mild	Normal	Strong	Ves	J
	D12	Overcast	Mild	High	Strong	Yes	
C	D13	Overcast	Hot	Normal	Weak	Yes	
	D14	Rain	Mild	High	Strong	No	



Partitioning examples: H = Normal and W = Strong

	Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
	D1	Sunny	Hot	High	Weak	No
	D2	Sunny	Hot	High	Strong	No
	D3	Overcast	Hot	High	Weak	Yes
	D4	Rain	Mild	High	Weak	Yes
	D5	Rain	Cool	Normal	Weak	Yes
	D6	Rain	Cool	Normal	Strong	No
	D7	Overcast	Cool	Normal	Strong	Yes
	D8	Sunny	Mild	High	Weak	No
	D9	Sunny	Cool	Normal	Weak	Yes
	D10	Rain	Mild	Normal	Weak	Yes
	D11	Sunny	Mild	Normal	Strong	Ves
	D12	Overcast	Mild	High	Strong	Yes
C	D13	Overcast	Hot	Normal	Weak	Yes
	D14	Rain	Mild	High	Strong	No



Decision Trees



- Can represent any Boolean Function
- Can be viewed as a way to compactly represent a lot of data.
- The evaluation of the Classifier is easy
- Clearly, given data, there are many ways to represent it as a decision tree.
- Learning a good representation from data is the challenge.



Features

- Outlook:
- > Temperature:
- Humidity:
- ► Wind:

{Sun, Overcast, Rain} {Hot, Mild, Cool} {High, Normal, Low} {Strong, Weak}

Labels

Binary classification task: Y = {+, -}



	0	Т	н	W	Play?
1	S	Н	Н	W	-
2	S	Н	Н	S	-
3	0	Н	Н	W	+
4	R	Μ	Н	W	+
5	R	С	Ν	W	+
6	R	С	Ν	S	-
7	0	С	Ν	S	+
8	S	М	Н	W	-
	S	С	Ν	W	+
)	R	М	Ν	W	+
1	S	Μ	Ν	S	+
2	0	М	Н	S	+
3	0	Н	Ν	W	+
4	R	М	Н	S	-







Top-down Induction of Decision Trees

CurrentNode = Root

DTtrain (examples for CurrentNode, features at CurrentNode):

- Find F, the "best" decision feature for next node
- For each value of F, create new descendant of node
- 3. Sort training examples to leaf nodes
- 4. If training examples perfectly classified Stop

Else

Recursively apply DTtrain over new leaf nodes



Picking the Root Attribute

- The goal is to have the resulting decision tree as small as possible (Occam's Razor)
 - However, finding the minimal decision tree consistent with the data is NP-hard
- The recursive algorithm is a greedy heuristic search for a simple tree, but cannot guarantee optimality.
- The main decision in the algorithm is the selection of the next attribute to condition on.



Picking the Root Attribute

Training Data with 2 Boolean attributes (A,B)

((A=0,B=0), -): 50 examples ((A=0,B=1), -): 50 examples ((A=1,B=0), -): 0 examples ((A=1,B=1), +): 100 examples

What should be the first attribute we select?



Picking the Root Attribute

Training Data with 2 Boolean attributes (A,B) ((A=0,B=0), -): 50 examples ((A=0,B=1), -): 50 examples ((A=1,B=0), -): 3 examples ((A=1,B=1), +): 100 examples

Trees looks structurally similar; which attribute should we choose?



Another feature selection criterion: Entropy

- Used in the ID3 algorithm [Quinlan, 1963]
 - pick feature with smallest entropy to split the examples at current iteration
- Entropy measures impurity of a sample of examples







H(*X*) is the expected number of bits needed to encode a randomly drawn value of *X* (under most efficient code)

Why? Information theory:

- Most efficient possible code assigns -log₂ P(X=i) bits to encode the message X=i
- So, expected number of bits to code one random *X* is:

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



Sample Entropy



High Entropy – High level of uncertainty

Low Entropy – Low level of uncertainty

- $\bullet~S$ is a sample of training examples
- p_{\oplus} is the proportion of positive examples in S
- p_{\ominus} is the proportion of negative examples in S
- \bullet Entropy measures the impurity of S

 $H(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$

MARYLAND

Information Gain

 The information gain of an attribute *a* is the expected reduction in entropy caused by partitioning on this attribute

$$Gain(S, a) = Entropy(S) - \sum_{v \in values(a)} \frac{|S_v|}{|S|} Entropy(S_v$$

- Original set is S.
- > S_{v} is the subset of S for which attribute a has value v.
- The entropy of partitioning the data is calculated by weighing the entropy of each partition by its size relative to the original set.
 - Partitions of low entropy (imbalanced splits) lead to high gain
- Take Home Exercise: go back and check which of the A, B splits is better.



	0	Т	н	W	Play?
1	S	Н	Н	W	-
2	S	Н	Н	S	-
3	0	Н	Н	W	+
4	R	Μ	Н	W	+
5	R	С	Ν	W	+
6	R	С	Ν	S	-
7	0	С	Ν	S	+
8	S	М	Н	W	-
	S	С	Ν	W	+
)	R	М	Ν	W	+
1	S	Μ	Ν	S	+
2	0	М	Н	S	+
3	0	Н	Ν	W	+
4	R	М	Н	S	-



	0	Т	Н	W	Play?
1	S	Н	Н	W	-
2	S	Н	Н	S	-
3	0	Н	Н	W	+
4	R	Μ	Н	W	+
5	R	С	Ν	W	+
6	R	С	Ν	S	-
7	0	С	Ν	S	+
8	S	М	Н	W	-
9	S	С	Ν	W	+
10	R	М	Ν	W	+
11	S	Μ	Ν	S	+
12	0	М	Н	S	+
13	0	Н	Ν	W	+
14	R	М	Н	S	-



		-	<u> </u>	147	Dlaw2	
	U		п	VV	Play?	Outlook = sunny:
1	S	Н	Н	W	-	p=2/5 n =3/5 H _s =0.971
2	S	Н	Н	S	-	
3	0	Н	Н	W	+	Outlook =overcast:
4	R	Μ	Н	W	+	p=4/4 n =0 H _o =0
5	R	С	Ν	W	+	
6	R	С	Ν	S	-	Outlook = rainy:
7	0	С	Ν	S	+	p=3/5 n =2/5 H _R =0.971
8	S	Μ	Н	W	-	
9	S	С	Ν	W	+	Expected entropy
10	R	Μ	Ν	W	+	$\left(\frac{5}{11}\right) \times 0.971 + \left(\frac{4}{11}\right) \times 0$
11	S	Μ	Ν	S	+	(14) (14)
12	0	Μ	Н	S	+	$+\left(\frac{3}{14}\right) \times 0.971 = 0.694$
13	0	Н	Ν	W	+	Information gain:
14	R	Μ	Н	S	-	0.940-0.694 = 0.246
14	R	М	Н	S	-	0.940-0.694 = 0.246



	0	Т	Н	W	Play?	Humidity = sunny:
1	S	Н	Н	W	-	p=3/7 n =4/7 H _H =0.985
2	S	Н	н	S	-	
3	0	Н	н	W	+	Humidity =overcast:
4	R	М	н	W	+	p=6/7 n =1/7 H _o =0.592
5	R	С	Ν	W	+	
6	R	С	Ν	S	-	
7	0	С	Ν	S	+	Expected entropy
8	S	М	н	W	-	$\binom{7}{-} \times 0.985 + \binom{7}{-} \times 0.592$
9	S	С	Ν	W	+	
10	R	М	Ν	W	+	= 0.7/85
11	S	Μ	Ν	S	+	0 0/0_0 7785 – 0 1515
12	0	М	н	S	+	$0.340^{-}0.7703^{-}0.1313$
13	0	Н	Ν	W	+	
14	R	М	н	S	_	



Which feature to split on?

			<u> </u>		
	0	T	Н	W	Play?
1	S	Н	Н	W	-
2	S	Н	Н	S	-
3	0	Н	Н	W	+
4	R	М	Н	W	+
5	R	С	Ν	W	+
6	R	С	Ν	S	-
7	0	С	Ν	S	+
8	S	М	Н	W	-
9	S	С	Ν	W	+
10	R	М	Ν	W	+
11	S	М	Ν	S	+
12	0	М	н	S	+
13	-			14/	
	0	н	N	VV	+



An Illustrative Example (III)

	0	Т	Н	W	Play?	
1	S	Н	Н	W	-	Outlook
2	S	Н	Н	S	-	
3	0	Н	Н	W	+	
4	R	Μ	Н	W	+	Sunny Overcast Rain
5	R	С	Ν	W	+	1.2.8.9.11 3.7.12.13 4.5.6.10.14
6	R	С	Ν	S	-	2+,3- 4+,0- 3+,2-
7	0	С	Ν	S	+	? <u>Yes</u> ?
8	S	Μ	Н	W	-	
9	S	С	Ν	W	+	
10	R	Μ	Ν	W	+	Continue until:
11	S	Μ	Ν	S	+	• Every attribute is included in path, or,
12	0	Μ	Н	S	+	 All examples in the leaf have same label
13	0	Н	Ν	W	+	
14	R	М	Н	S	_	



An Illustrative Example (IV)



Day	Outlook	Temperature	Humidit	y Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	Νο
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes



An Illustrative Example (V)





induceDecisionTree(S)

- 1. Does *S* uniquely define a class? if all $s \in S$ have the same label *y*: return *S*;
- 2. Find the feature with the most information gain:
 - $i = \operatorname{argmax}_i Gain(S, X_i)$
- 3. Add children to S:

for k in Values(X_i): $S_k = \{s \in S | X_i = k\}$ addChild(S, S_k) induceDecisionTree(S_k) return S;



An Illustrative Example (VI)





Hypothesis Space in Decision Tree Induction

- Conduct a search of the space of decision trees which can represent all possible discrete functions. (pros and cons)
- Goal: to find the best decision tree
- Finding a minimal decision tree consistent with a set of data is NP-hard.
- Performs a greedy heuristic search: hill climbing without backtracking.
- Makes statistically based decisions using all data.



A decision tree to distinguish homes in New York from homes in San Francisco

Take a look at home:

http://www.r2d3.us/visual-intro-to-machine-learning-part-1/



Check out course webpage, Canvas, Piazza

Submit HW01 I due Thursday 10:30am

Do the readings!



Furong Huang 3251 A.V. Williams, College Park, MD 20740 301.405.8010 / furongh@cs.umd.edu