### Lecture 5: Decision Trees and Their Representation COMP 343, Spring 2022 Victoria Manfredi





**Acknowledgements:** These slides are based primarily on content from the book "Machine Learning" by Tom Mitchell, and on slides created by Vivek Srikumar (Utah), Dan Roth (Penn), and Jessica Wu (Harvey Mudd College)

# **Today's Topics**

### Homework 2 out

- Due Wed., February 16 by 5p

### Decision trees (non-linear classifiers)

- Our first learning algorithm!
- Representation: what are decision trees?

**Looking Forward** 

# Using supervised learning

### ✓ What is our instance space?

• What are the inputs to the problem? What are the features?

### ✓ What is our label space?

• What kind of learning task are we dealing with?

### What is our hypothesis space?

- What functions should the learning algorithm search over?
- 4. What is our learning algorithm?
  - How do we learn the model from the labeled data?
- 5. What is our loss function or evaluation metric?
  - How do we measure success? What drives learning?

Much of the rest of the semester

# Coming up ...

### Different hypothesis spaces and learning algorithms

#### Decision trees and ID3 algorithm

 (Non-parametric) regression and classification

#### Linear regression

Least mean squares regression

#### **Linear classifiers**

Perceptron

#### **Non-linear classifiers**

Multi-layer perceptron (neural networks)

Provides a "gentle" introduction to machine learning concepts Rule-based decision algorithm

# Coming up ...

**1**. Representation: What are decision trees?

### 2. Algorithm: Learning decision trees

- The ID3 algorithm: a greedy heuristic
- 3. Some extensions

# Decision Trees REPRESENTATION

Data can be represented as a big table, with columns denoting different attributes

Name	Label
Norman Danner	+
Karen Collins	-
Dan Licata	+
Danny Krizanc	+
Saray Shai	+
Wai Kiu Chan	-

Data can be represented as a big table, with columns denoting different attributes

Name	Label
Norman Danner	+
Karen Collins	-
Dan Licata	+
Danny Krizanc	+
Saray Shai	+
Wai Kiu Chan	-

# Problem: this representation does not generalize.

If name is not in table, we don't know what to do. So we need to extract attributes from the data

#### Extract attributes from the data

		Attri	<i>y<sub>i</sub></i>		
	Name	Faculty department area	Name contains Dan	Second character of first name	Label
ce	Norman Danner	CS	Yes	0	+
instan	Karen Collins	Math	No	а	-
r is an	Dan Licata	CS	Yes	а	+
ch row	Danny Krizanc	CS	Yes	а	+
Ea(	Saray Shai	CS	No	а	+
	Wai Kiu Chan	Math	No	а	-

10

Once we add other (informative) columns, maybe we don't need original data. Instead work with featurized version of data

Name	Faculty department area	Name contains Dan	Second character of first name	Label
Norman Danner	CS	Yes	0	+
Karen Collins	Math	No	а	-
Dan Licata	CS	Yes	а	+
Danny Krizanc	CS	Yes	а	+
Saray Shai	CS	No	а	+
Wai Kiu Chan	Math	No	а	-

#### How big is this table?

Name	Faculty department area	Name contains Dan	Second character of first name	Label
Norman Danner	CS	Yes	0	+
Karen Collins	Math	No	а	-
Dan Licata	CS	Yes	а	+
Danny Krizanc	CS	Yes	а	+
Saray Shai	CS	No	а	+
Wai Kiu Chan	Math	No	а	-

#### How big is this table?

Name	Faculty department area	Name contains Dan	Second character of first name	Label
Norm With thes	e three attribute $2 \times$	es, how many ur $2 \times 26 = 104$	nique rows are p	ossible?
Karen	y he an infinite r	umber of name	s but we can be	ave at most
Dan L 104 possi	ble buckets for r	names!	s, but we can ne	
Danny Krizanc	CS	Yes	а	+
Saray Shai	CS	No	а	+
Wai Kiu Chan	Math	No	а	-

Name	Faculty department area	Name contains Dan	Second character of first name	Label	
Norm	With these three attribut $2 \times 2$	es, how many units $2 \times 26 = 104$	nique rows are p	ossible?	
Karen					
Dan L	There may be an infinite number of names, but we can have at most 104 possible buckets for names!				
Danny	If there are 100 attributes, all binary, how many unique rows are possible?				
Saray					
Wai K					

Name	Faculty department area	Name contains Dan	Second character of first name	Label	
Norm	With these three attribut $\frac{2}{2}$ ×	tes, how many uncertainty $2 \times 26 = 104$	nique rows are p	ossible?	
Karen					
Dan L	There may be an infinite number of names, but we can have at most 104 possible buckets for names!				
Danny	If there are 100 attributes, all binary, how many unique rows are possible?				
Saray	(100 times) $2 \times 2 \times$	$2 \times 2 \times \cdots \times 2$	$= 2^{100}$		
Wai K					

Name	Faculty department area	Name contains Dan	Second character of first name	Label	
Norm	With these three attribut	es, how many units $2 \times 26 - 104$	nique rows are p	ossible?	
Karen		2 × 20 = 10+			
Danl	There may be an infinite	number of name	es, but we can ha	ave at most	
Dan L	104 possible buckets for	names!			ļ
Danny	If there are 100 attributes, all binary, how many unique rows are possible?				
Saray	(100 times) $2 \times 2 \times 10^{-10}$	$2 \times 2 \times \cdots \times 2$	$= 2^{100}$		
Wai K	If we wanted to store all to represent data in bette	possible rows, th er, more efficient	is number is <mark>toc</mark> way?	large. How	

Name	Faculty department area	Name contains Dan	Second character of first name	Label	
Norm One v	vay of thinkin	g about decis	ion trees is as	a more	
Karen	efficient way of organizing data				
Dan L					
Danny Krizanc	CS	Yes	а	+	
Saray Shai	CS	No	а	+	
Wai Kiu Chan	Math	No	а	-	

A hierarchical data structure that represents data using a divide-and-conquer strategy

A hierarchical data structure that represents data using a divide-and-conquer strategy

Decision trees are also a hypothesis class. Can be used as a hypothesis class for non-parametric classification or regression

A hierarchical data structure that represents data using a divide-and-conquer strategy

Decision trees are also a hypothesis class. Can be used as a hypothesis class for non-parametric classification or regression

General idea: given a collection of examples, learn a decision tree that represents it. Use this representation to classify new examples

Decision trees are a family of classifiers (our focus) for instance that are represented by collections of attributes (i.e., feature vectors: color=; shape=; label=)

Decision trees are a family of classifiers (our focus) for instance that are represented by collections of attributes (i.e., feature vectors: color=; shape=; label=)

Color?

How do decision trees work?

Nodes are tests for feature values

Decision trees are a family of classifiers (our focus) for instance that are represented by collections of attributes (i.e., feature vectors: color=; shape=; label=)

How do decision trees work?



- Nodes are tests for feature values
- There is one branch for every value that the feature can take

Decision trees are a family of classifiers (our focus) for instance that are represented by collections of attributes (i.e., feature vectors: color=; shape=; label=)

How do decision trees work?

Nodes are tests for feature values



- There is one branch for every value that the feature can take
- Leaves of the tree specify the category (class labels)

Decision trees are a family of classifiers (our focus) for instance that are represented by collections of attributes (i.e., feature vectors: color=; shape=; label=)

How do decision trees work?

Nodes are tests for feature values



- There is one branch for every value that the feature can take
- Leaves of the tree specify the category (class labels)

### Like playing twenty questions game



#### Our dataset comprises shapes



What are some attributes of the shapes?



What are some attributes of the shapes? *Color, shape, size* 



What are some attributes of the shapes? *Color, shape, size* 

Let's build a decision tree



What are some attributes of the shapes? *Color, shape, size* 

First, what is the color of the shape?



What are some attributes of the shapes? Color, shape, size Blue



3 options for color, so 3 edges



What are some attributes of the shapes? Color, shape, size



If color is red, you know label immediately. What is it?



What are some attributes of the shapes? Color, shape, size



*If color is red, you know label immediately. What is it?* 



What are some attributes of the shapes? Color, shape, size If color is blue, If color is blue,

what should we ask about next?














# 3 min: What is the label for a red triangle? And why?

Even though no red triangle in data, can still label it. Start at root and go down tree. If color is red, go down that edge and label is B.







- 1. How do we learn a decision tree? Coming up soon
- 2. How to use a decision tree for prediction?
  - What is the label for a red triangle?

 $\implies$  Just follow a path from the root to a leaf





If we knew shapes could be circles, squares or triangle, then add branch





3 min: Think about how to label branch?











One option: add another edge in tree for "unseen"









#### Set of possible instances, $\boldsymbol{X}$

- Each instance  $\mathbf{x} \in X$  is a feature vector
- E.g., <shape=square, color=red>

#### Set of possible instances, $\boldsymbol{X}$

- Each instance  $\mathbf{x} \in X$  is a feature vector
- E.g., <shape=square, color=red>

#### Set of possible labels, $\boldsymbol{Y}$

Y is discrete-valued

#### Set of possible instances, $\boldsymbol{X}$

- Each instance  $\mathbf{x} \in X$  is a feature vector
- E.g., <shape=square, color=red>

#### Set of possible labels, $\boldsymbol{Y}$

Y is discrete-valued

Unknown target function,  $f: X \to Y$ 

#### Set of possible instances, $\boldsymbol{X}$

- Each instance  $\mathbf{x} \in X$  is a feature vector
- E.g., <shape=square, color=red>

#### Set of possible labels, $\boldsymbol{Y}$

Y is discrete-valued

Unknown target function,  $f: X \to Y$ 

Set of function hypotheses,  $H = \{h \mid h : X \rightarrow Y\}$ 

- Each hypothesis h is a decision tree
- Trees sort  $\mathbf{x}$  to leaf which assigns  $y \in Y$

What is the hypothesis class that decision trees represent?

Remember, a hypothesis is just a function that maps instances to labels

So what class of functions can decision trees represent?

3min: what do you think this class of functions is?

What functions can decision trees represent?

What functions can decision trees represent?



What functions can decision trees represent?



What functions can decision trees represent?



What functions can decision trees represent?

• Any Boolean function (assuming all features are Boolean)



(Color=Blue AND Shape=Triangle => Label=B) AND (Color=Blue AND Shape=Square => Label=A) AND (Color=Blue AND Shape=Circle => Label=C) AND ....

What functions can decision trees represent?

• Any Boolean function (assuming all features are Boolean)

3 min: how would you represent this function as a decision tree?

Α	В	A xor B
F	F	F
F	т	т
Т	F	т
т	т	F

What functions can decision trees represent?



Α	В	A xor B
F	F	F
F	Т	т
Т	F	т
т	т	F

What functions can decision trees represent?

• Any Boolean function (assuming all features are Boolean)



In worst case, how many nodes needed?

Α	В	A xor B
F	F	F
F	Т	т
т	F	т
т	т	F

What functions can decision trees represent?

• Any Boolean function (assuming all features are Boolean)



In worst case, how many nodes needed? Exponentially many!

Decision trees have variable-sized hypothesis space: as # of nodes (or depth) increases, hypothesis space grows (can express more complex functions)

Α	В	A xor B
F	F	F
F	Т	т
т	F	т
т	т	F

#### **Decision Trees**

Outputs are discrete categories (classification)

#### **Decision Trees**

Outputs are discrete categories (classification)

But real valued outputs are possible (regression trees)

### **Decision Trees**

Outputs are discrete categories (classification)

But real valued outputs are possible (regression trees)

There are methods for handling noisy data (noise in the label or in the features) and for handling missing attribute values. Pruning trees helps with noise. Smaller trees also tend to generalize better. More on this later ...

We have seen instances represented as attribute-value pairs (color=Blue, shape=Square, second letter=a)

– Values have been categorical

We have seen instances represented as attribute-value pairs (color=Blue, shape=Square, second letter=a)

– Values have been categorical

How do we deal with numeric feature values (e.g., length = ?)

What can we do?

We have seen instances represented as attribute-value pairs (color=Blue, shape=Square, second letter=a)

– Values have been categorical

How do we deal with numeric feature values (e.g., length = ?)

Discretize values

– Use thresholds on the values for splitting nodes












# Decision Trees TAKEAWAYS

## Summary: decision trees

Decision trees can represent any Boolean function A way to represent a lot of data A natural representation (think playing 20 questions) Predicting with a decision tree classifier is easy

## Summary: decision trees

Decision trees can represent any Boolean function A way to represent a lot of data A natural representation (think playing 20 questions) Predicting with a decision tree classifier is easy

Clearly, given a dataset, there are many decision trees that can represent it

Learning a good representation from data is the challenge

#### Problem setting for decision tree learning

#### Set of possible instances, X

- Each instance  $\mathbf{x} \in X$  is a feature vector
- E.g., <shape=square, color=red>

#### Set of possible labels, $\boldsymbol{Y}$

Y is discrete-valued

Unknown target function,  $f: X \to Y$ 

Set of function hypotheses,  $H = \{h \mid h : X \rightarrow Y\}$ 

- Each hypothesis h is a decision tree
- Trees sort  $\mathbf{x}$  to leaf which assigns  $y \in Y$