# 10-301/601: Introduction to Machine Learning Lecture 3 — Decision Trees

Henry Chai & Matt Gormley 9/7/22

#### **Front Matter**

- Announcements:
  - HW1 released 8/29, due 9/7 (today!) at 11:59 PM
    - Keep an eye out on Piazza for the HW1 exit poll,
       which will count towards your participation grade
      - You must complete this poll within one week of its release to receive full credit
  - HW2 released 9/7 (today!), due 9/19 at 11:59 PM
    - Unlike HW1, you will only have:
      - 1 submission for the written portion
      - 10 submissions of the programming portion to our autograder

#### Q & A:

# Do I have to use LaTex to complete the HWs?

- Technically no... but we do strongly encourage it! So much so that we are offering one bonus point for each assignment that you complete in LaTex.
- Follow-up: but what if I don't know LaTex?
  - Don't worry, we have you covered! Over on Piazza, we've pinned some resources on how to use LaTex and our TAs have recorded an awesome tutorial which you can find in our course Panopto folder.
- Regardless of whether or not you use LaTex, you are responsible for making sure that your submission is aligned with our PDF template; do not omit any pages!

#### Q & A:

### How do these in-class polls work?

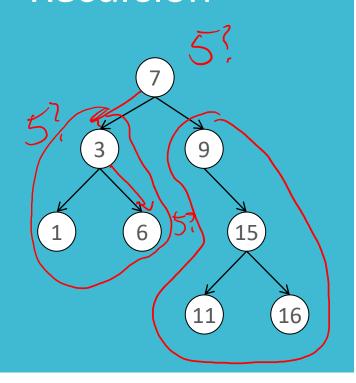
- Open the poll, either by clicking the [Poll] link on the schedule page of our course website or going to <a href="http://poll.mlcourse.org">http://poll.mlcourse.org</a>
- Sign into Google Forms using your Andrew email
- Answer all poll questions during lecture for full credit or within 24 hours for half credit
- Avoid the toxic option (will be clearly specified in lecture) which gives negative poll points
- You have 8 free "poll points" for the semester that will excuse you from all polls from a single lecture; you cannot use more than 3 poll points consecutively.

#### Poll Question 1:

Which of the following did you bring to class today? Select all that apply

- A. A smartphone
- B. A flip phone
- C. A payphone
- D. No phone

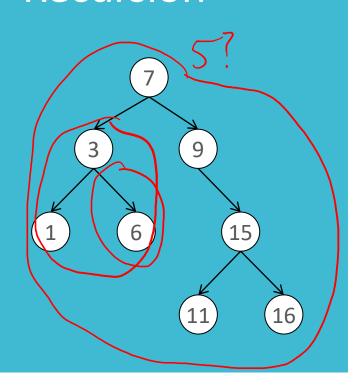
### Background: Recursion



- A binary search tree (BST) consists of nodes, where each node:
  - has a value, v
  - up to 2 children, a left descendant and a right descendant
  - all its left descendants have values less than v and its right descendants have values greater than v
- We like BSTs because they permit search in O(log(n)) time, assuming n nodes in the tree

```
def contains_iterative(node, key):
    cur = node
    while true:
        if key < cur.value & cur.left != null:
            cur = cur.left
        else if cur.value < key & cur.right != null:
            cur = cur.right
        else:
            break
    return key == cur.value</pre>
```

### Background: Recursion



- A binary search tree (BST) consists of nodes, where each node:
  - has a value, v
  - up to 2 children, a left descendant and a right descendant
  - all its left descendants have values less than v and its right descendants have values greater than v
- We like BSTs because they permit search in O(log(n)) time, assuming n nodes in the tree

  def contains\_recursive(node, key):

  if key < node.value & node.left != null:

  return contains(node.left, key)

else if node.value < key & node.right != null:
 return contains(node.right, key)</pre>

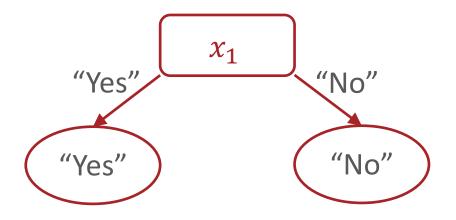
else:

return key == node.value

### Recall: Decision Stumps

• Alright, let's actually (try to) extract a pattern from the data

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?	$\hat{y}$ Predictions
Yes	Low	Normal	No	Yes
No	Medium	Normal	No	No
No	Low	Abnormal	Yes	No
Yes	Medium	Normal	Yes	Yes
Yes	High	Abnormal	Yes	Yes



# Decision Stumps: Questions

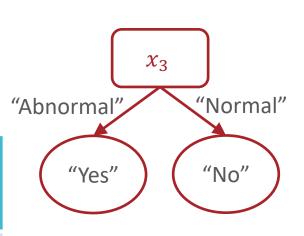
1. Why stop at just one feature?

2. How can we pick which feature to split on?

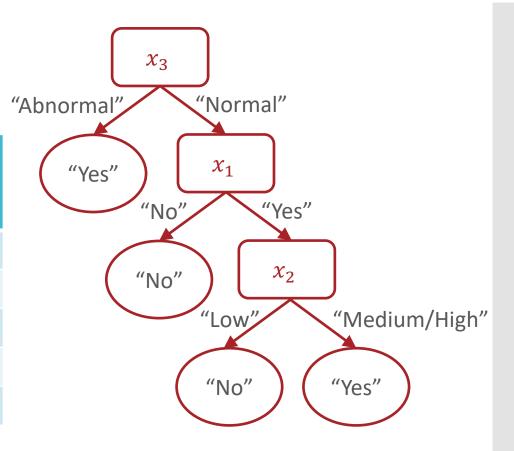
# From Decision Stump

• • •

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes



$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes



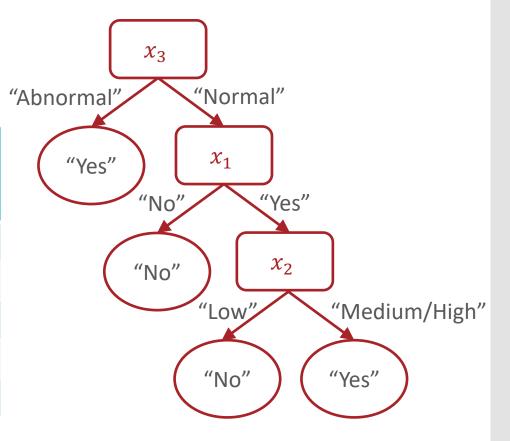
9/7/22 **11** 

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

Normal

High

No

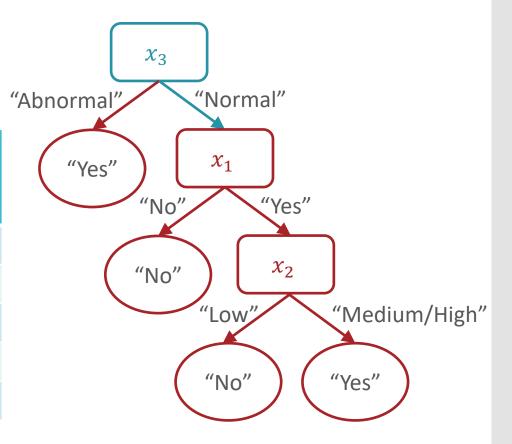


9/7/22

No

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

No	High	Normal	No
	9		



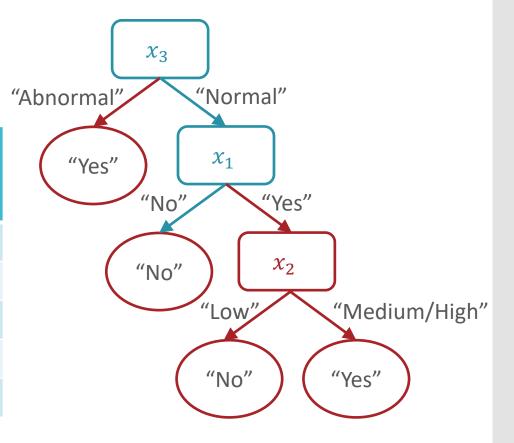
13

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

Normal

High

No



9/7/22

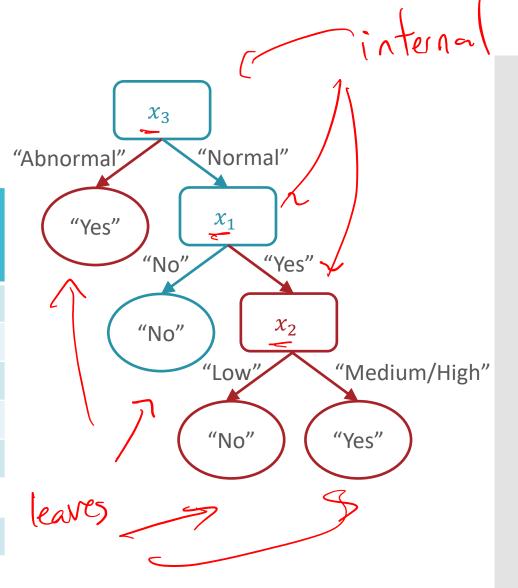
No

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

Normal

High

No



9/7/22

No

# Decision Tree: Pseudocode

```
det predict_thrative(x'):
     - Walk from the root node to a leaf
     while (five):
           if the current node is internal (non-kat):
              chick the associated attribute X
              move down the branch according to X
              return the label stored at that made
```

### Decision Tree: Example

Learned from medical records of 1000 women Negative examples are C-sections

```
[833+,167-] .83+ .17-
-> Fetal_Presentation = 1: [822+,116-] .88+ .12-
\rightarrow | Previous_Csection = 0: [767+,81-] .90+ .10-
-> | | Primiparous = 0: [399+,13-] .97+ .03-
 ->| | Primiparous = 1: [368+,68-] .84+ .16-
   | \ | \ | Fetal_Distress = 0: [334+,47-] .88+ .12-
     | | Fetal_Distress = 1: [34+,21-] .62+ .38-
 \rightarrow | Previous_Csection = 1: [55+,35-] .61+ .39-
 >Fetal_Presentation = 2: [3+,29-] .11+ .89-
   Fetal_Presentation = 3: [8+,22-] .27+ .73-
```

9/7/22 Figure courtesy of Tom Mitchell

# Decision Stumps: Questions

1. Why stop at just one feature?

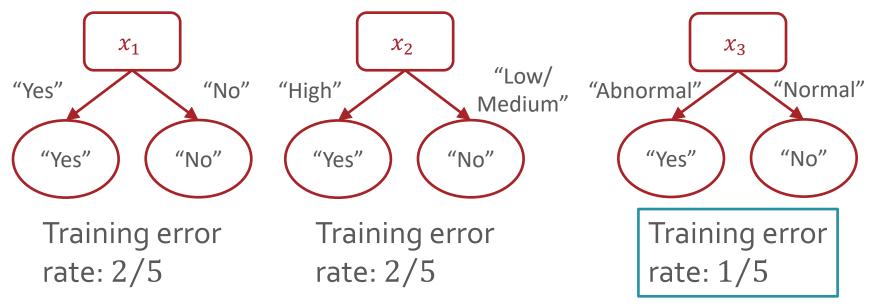
2. How can we pick which feature to split on as well as the order of the splits?

### Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion

# Training Error Rate as a Splitting Criterion

$x_1$ Family History	$x_2$ Resting Blood Pressure	$x_3$ Cholesterol	<i>y</i> Heart Disease?
Yes X	Low	Normal	No
No	Medium	Normal	No
No X	Low X	Abnormal	Yes
Yes	Medium X	Normal X	Yes
Yes	High	Abnormal	Yes



9/7/22 **20** 

#### Poll Question 2:

Which feature would you split on using training error rate as the splitting criterion?

$x_1$	$x_2$	у
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1
1	1	1

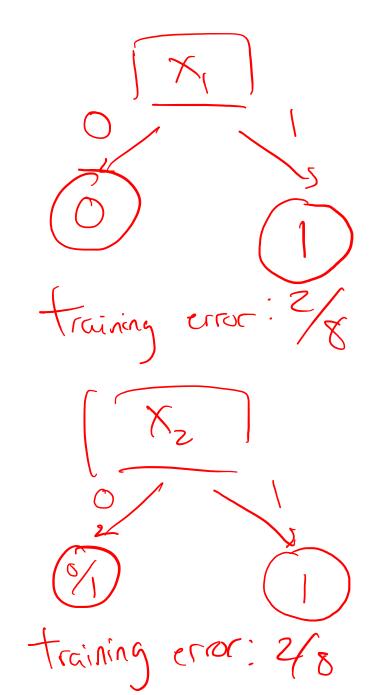
- A.  $x_1$
- B.  $x_2$
- C. Either  $x_1$  or  $x_2$
- D. Neither  $x_1$  nor  $x_2$

9/7/22 **21** 

#### Poll Question 2:

Which feature would you split on using training error rate as the splitting criterion?

$x_1$	$x_2$	У
1 ×	7 0	0
1 ×	0 ×	0
1	<b>y</b> 0	1
1	<b>✓</b> 0	1
1	1	1
1	1	1
1	1	1
1	1	1



### Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion
- Potential splitting criteria:
  - Training error rate (minimize)
  - Gini impurity (minimize) → CART algorithm
  - Mutual information (maximize) → ID3 algorithm

9/7/22 **23** 

### Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Idea: when deciding which feature to split on, use the one that optimizes the splitting criterion
- Potential splitting criteria:
  - Training error rate (minimize)
  - Gini impurity (minimize) → CART algorithm
  - Mutual information (maximize) → ID3 algorithm

24

#### Entropy

• Entropy describes the purity or uniformity of a collection of values: the lower the entropy, the more pure

$$H(s) = -\frac{\sum_{v \in V(S)} |S_v|}{|S_v|} \log_2 \frac{|S_v|}{|S_v|}$$
when  $S$  is a set of values
$$V(S) \approx \text{all unique values in } S$$

$$S_v \approx \text{all elements in } S_v \neq \text{value } v$$

If all the elements in S are the same, then

$$H(S) = -\left(\frac{|S|}{|S|} \log_2 \frac{|S|}{|S|}\right) = -\left(\frac{|S|}{|S|}\right) = 0$$

#### Entropy

• Entropy describes the purity or uniformity of a collection of values: the lower the entropy, the more pure

$$H(S) = -\sum_{v \in V(S)} \frac{|S_v|}{|S|} \log_2 \left(\frac{|S_v|}{|S|}\right)$$

where *S* is a collection of values,

V(S) is the set of unique values in S

 $S_v$  is the collection of elements in S with value v

• If S is split fifty-fifty between two values, then

$$H(S) = -\left(\frac{1}{2}\log_z \frac{1}{2} + \frac{1}{2}\log_z \frac{1}{2}\right)$$

$$= -\left(\frac{1}{2}(-1) + \frac{1}{2}(-1)\right) =$$

### Mutual Information

 Mutual information describes how much information or clarity a particular feature provides about the label

### Mutual Information: Example

	+, =	74 0	- + = /4
	$x_d$	y	
	1	1	7
	1	1	
	<b>Q</b>	0	fraction
~/-	0	0	) Jof points where $X_i = V$
$I(\mathcal{H}_{a}, X) =$		$\sum_{V(x_d)} (f_v) \Big($	$H(Y_{x_d=v})$
		7 H(7/2=1	$\left(\frac{1}{2} + \frac{1}{2} + $
		$\mathcal{T}$	

# Mutual Information: Example

$x_d$	y	
1	1	1
0	1	
1	0	\ \(  \)
0	0	2

$$I(\mathbf{x}_{d}(\mathbf{x})) = H(Y) - \sum_{v \in V(x_d)} (f_v) \left( H(Y_{x_d=v}) \right)$$

$$= \left( \frac{1}{Z} \left( \left( Y_{x_d=v} \right) + \frac{1}{Z} H(Y_{x_d=v}) \right) \right)$$

$$= \left( \frac{1}{Z} \left( \left( Y_{x_d=v} \right) + \frac{1}{Z} H(Y_{x_d=v}) \right) \right)$$

#### **Poll Question 3:**

Which feature would you split on using mutual information as the splitting criterion?

$x_1$	$x_2$	У
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1
1	1	1

A.  $x_1$ 

B.  $x_2$ 

C. Either  $x_1$  or  $x_2$ 

D. Neither  $x_1$  nor  $x_2$