

10-301/601: Introduction to Machine Learning

Lecture 2 – Decision Trees

Henry Chai

5/18/22

Front Matter

- Announcements:
 - HW1 released 5/17 due 5/24 at 1 PM
 - Recitation 1 on 5/19: review of prerequisite material
 - General advice for the summer:
 - Start HWs early!
 - Go to office hours! Starting today, 5/18
- Recommended Readings:
 - Daumé III, Chapter 1: Decision Trees

Our second Machine Learning Classifier

- A **classifier** is a function that takes feature values as input and outputs a label
- Memorizer: if a set of features exists in the **training** dataset, predict its corresponding label; otherwise, predict the majority vote

Family History	Resting Blood Pressure	Cholesterol	Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes

Notation

- Feature space, \mathcal{X}
- Label space, $\mathcal{Y} = \{\text{"No"}, \text{"Yes"}\}$
- (Unknown) Target function, $c^*: \mathcal{X} \rightarrow \mathcal{Y}$
- Training dataset:

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, c^*(\mathbf{x}^{(1)}) = y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}) \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$$

- Data point:

$$(\mathbf{x}^{(n)}, y^{(n)}) = (\underbrace{[x_1^{(n)}, x_2^{(n)}, \dots, x_D^{(n)}]}_{\text{feature vector}}, y^{(n)})$$

- Classifier, $h : \mathcal{X} \rightarrow \mathcal{Y}$
- Goal: find a classifier, h , that best approximates c^*

Evaluation

- Loss function, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
 - Defines how “bad” predictions, $\hat{y} = h(\mathbf{x})$, are compared to the true labels, $y = c^*(\mathbf{x})$
 - Common choices
 1. Squared loss (for regression): $\ell(y, \hat{y}) = (y - \hat{y})^2$
 2. Binary or 0-1 loss (for classification):
$$\ell(y, \hat{y}) = \begin{cases} 1 & \text{if } \underline{y \neq \hat{y}} \\ 0 & \text{otherwise} \end{cases}$$

Evaluation

- Loss function, $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
 - Defines how “bad” predictions, $\hat{y} = h(\mathbf{x})$, are compared to the true labels, $y = c^*(\mathbf{x})$
 - Common choices
 1. Squared loss (for regression): $\ell(y, \hat{y}) = (y - \hat{y})^2$
 2. Binary or 0-1 loss (for classification):

$$\ell(y, \hat{y}) = \mathbb{1}(y \neq \hat{y})$$

indicator function

- Error rate:

$$err(h, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^N \mathbb{1}(y^{(n)} \neq \hat{y}^{(n)})$$

Notation: Example

- Memorizer: if a set of features exists in the **training dataset**, predict its corresponding label; otherwise, predict the majority vote

	x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?	\hat{y} Predictions
	Yes	Low	Normal	No	No
$x^{(2)}$	No	Medium	Normal	No	No
	No	Low	Abnormal	Yes	Yes
	Yes	Medium	Normal	Yes	Yes
	Yes	High	Abnormal	Yes	Yes

- $N = 5$ and $D = 3$
- $x^{(2)} = (x_1^{(2)} = \text{"No"}, x_2^{(2)} = \text{"Medium"}, x_3^{(2)} = \text{"Normal"})$

Our second Machine Learning Classifier

- Memorizer:

```
def train( $\mathcal{D}$ ):
```

store \mathcal{D}

```
def majority_vote( $\mathcal{D}$ ):
```

return $\text{mode}(y^{(1)}, y^{(2)}, \dots, y^{(N)})$

```
def predict( $x'$ ):
```

{ if $\exists x^{(n)} \in \mathcal{D}$ s.t. $x^{(n)} == x'$:
return $y^{(n)}$

else:
return majority_vote(\mathcal{D})

Our third Machine Learning Classifier

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

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

- Decision stump: based on a single feature, x_d , predict the most common label in the training dataset among all data points that have the same value for x_d

Our third Machine Learning Classifier: Example

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

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

- Decision stump on x_1 :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_d) = \begin{cases} ??? & \text{if } x'_1 = \text{"Yes"} \\ ??? & \text{otherwise} \end{cases}$$

Our third Machine Learning Classifier: Example

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

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

- Decision stump on x_1 :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_d) = \begin{cases} \text{"Yes"} & \text{if } x'_1 = \text{"Yes"} \\ \text{???} & \text{otherwise} \end{cases}$$

Our third Machine Learning Classifier: Example

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

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

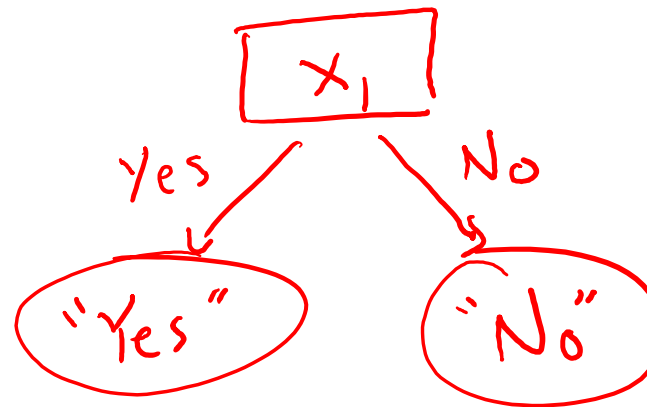
- Decision stump on x_1 :

$$h(\mathbf{x}') = h(x'_1, \dots, x'_n) = \begin{cases} \text{"Yes"} & \text{if } x'_1 = \text{"Yes"} \\ \text{"No"} & \text{otherwise} \end{cases}$$

Our third Machine Learning Classifier: Example

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

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y 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: Pseudocode

```
def train( $\mathcal{D}$ ):
```

1. pick a feature, x_d
2. Split \mathcal{D} according to x_d
for v in $V(x_d)$ = all possible values of x_d
 $\mathcal{D}_v = \{ (x^{(n)}, y^{(n)}) \in \mathcal{D} \mid x_d^{(n)} = v \}$
3. Compute the majority vote
for v in $V(x_d)$
 $\hat{y}_v = \text{majority_vote}(\mathcal{D}_v)$

```
def predict( $x'$ ):
```

- for v in $V(x_d)$
if $x'_d = v$, return \hat{y}_v

Decision Stumps: Questions

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

Lecture 2 Polls

0 done

 **0 underway**

Which feature do you think we should split on for this data set?

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

x_1

x_2

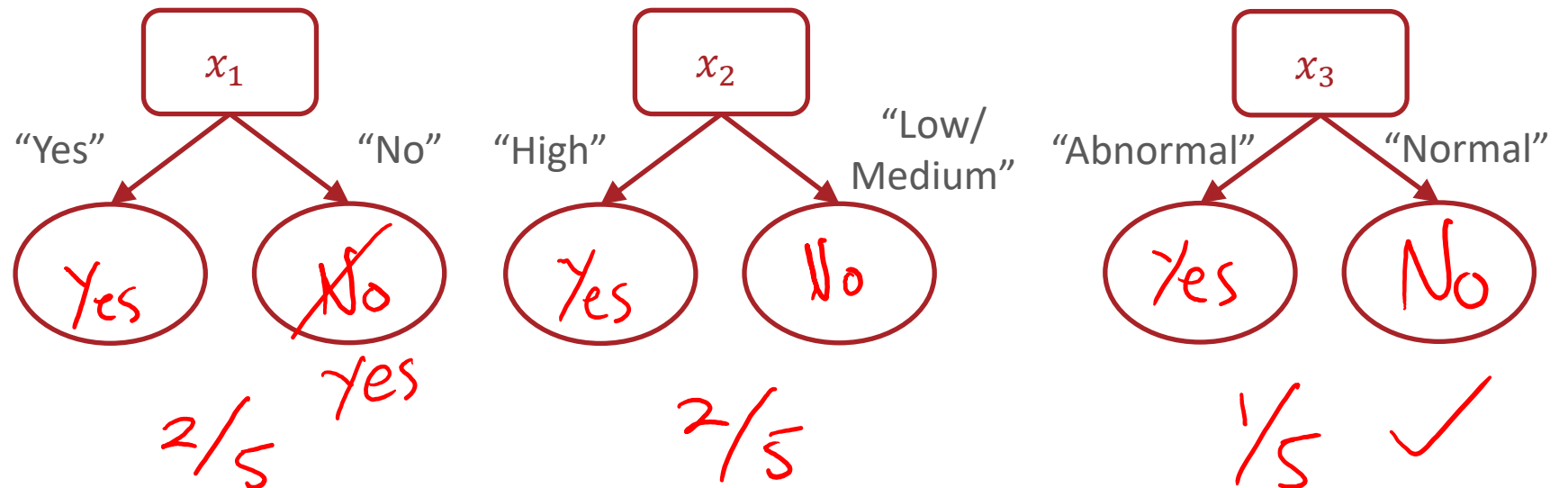
x_3

Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Insight: use the feature that optimizes the splitting criterion for our decision stump.

Training error rate as a Splitting Criterion

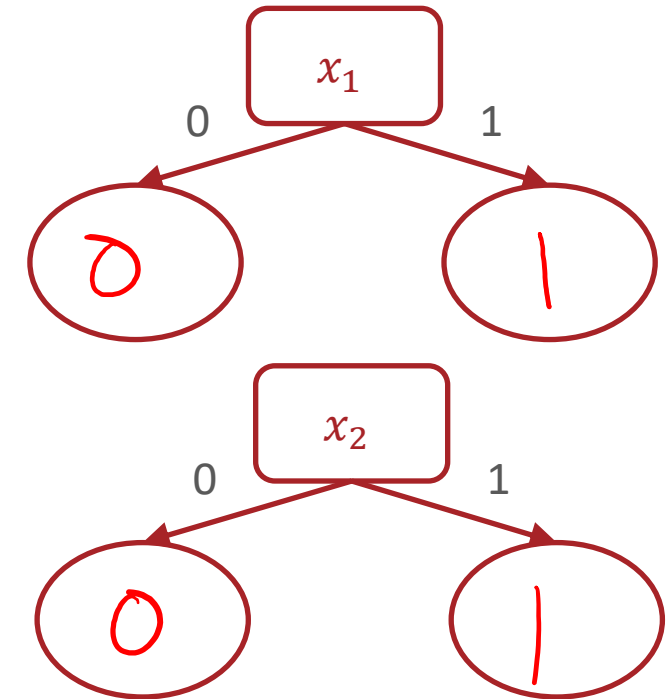
x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y Heart Disease?
Yes	Low	Normal	No
No	Medium	Normal	No
No	Low	Abnormal	Yes
Yes	Medium	Normal	Yes
Yes	High	Abnormal	Yes



Training error rate as a Splitting Criterion?

x_1	x_2	y
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1
1	1	1

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



2/8 for both

Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Insight: use the feature that optimizes the splitting criterion for our decision stump.
- Potential splitting criteria:
 - Training error rate (minimize)
 - Gini impurity (minimize) → CART algorithm
 - Mutual information (maximize) → ID3 algorithm

Splitting Criterion

- A **splitting criterion** is a function that measures how good or useful splitting on a particular feature is *for a specified dataset*
- Insight: use the feature that optimizes the splitting criterion for our decision stump.
- Potential splitting criteria:
 - Training error rate (minimize)
 - Gini impurity (minimize) → CART algorithm
 - Mutual information (maximize) → ID3 algorithm

Entropy

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

$$H(S) = - \sum_{v \in V(S)} \underbrace{\frac{|S_v|}{|S|}}_{1 = \text{size of}} \log_2 \left(\frac{|S_v|}{|S|} \right)$$

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 all the elements in S are the same, then

$$\frac{|S_v|}{|S|} = 1 \Rightarrow H(S) = -1 \log_2(1) = 0$$

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,

$\{Yes, No\} = V(S)$ is the set of unique values in S

$\{Yes, Yes, Yes\}$
 $\{No, No\}$ S_v is the collection of elements in S with value v

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

$$\frac{|S_0|}{|S|} = \frac{|S_1|}{|S|} = \frac{1}{2} \Rightarrow H(S) = - \left(\frac{1}{2} \log_2 \left(\frac{1}{2} \right) + \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) = 1$$

Mutual Information

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

$$I(x_d, Y) = H(Y) - \sum_{v \in V(x_d)} \underbrace{(f_v)} (H(Y_{x_d=v}))$$

where x_d is a feature

Y is the collection of all labels

$V(x_d)$ is the set of unique values of x_d

f_v is the fraction of inputs where $x_d = v$

$Y_{x_d=v}$ is the collection of labels where $x_d = v$

Mutual Information: Example

x_d	y
1	1
1	1
0	0
0	0



$$I(x_d, Y) = H(Y) - \sum_{v \in V(x_d)} (f_v) (H(Y_{x_d=v}))$$

$$= 1 - \left(\frac{1}{2} (0) + \frac{1}{2} (0) \right)$$
$$= 1$$

Mutual Information: Example

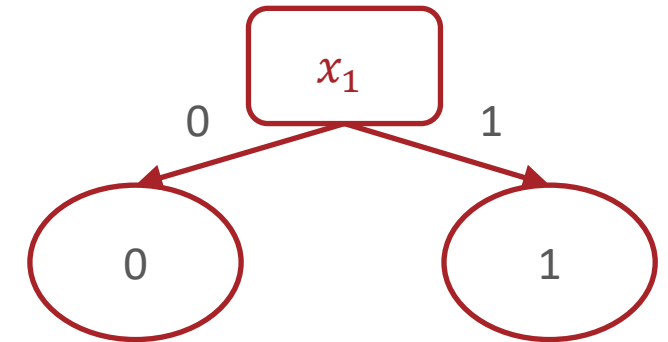
x_d	y
1	1
0	1
1	0
0	0

$$\begin{aligned} I(x_d, Y) &= H(Y) - \sum_{v \in V(x_d)} (f_v) (H(Y_{x_d=v})) \\ &= 1 - \left(\frac{1}{2} (1) + \frac{1}{2} (1) \right) \\ &= 0 \end{aligned}$$

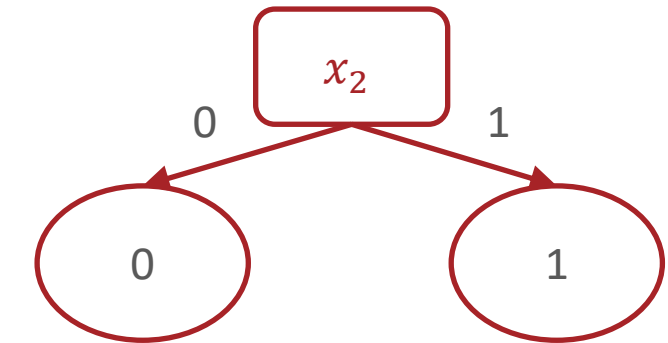
Mutual Information as a Splitting Criterion

x_1	x_2	y
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1
1	1	1

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



Mutual Information: 0

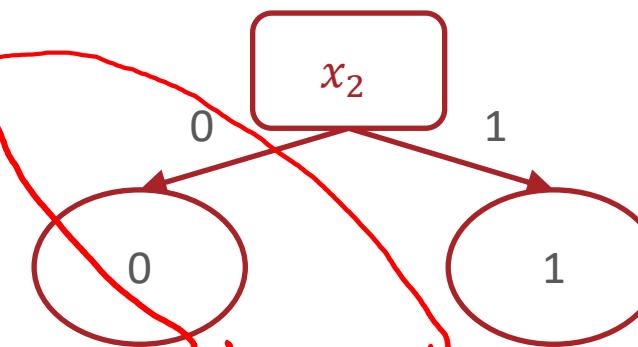
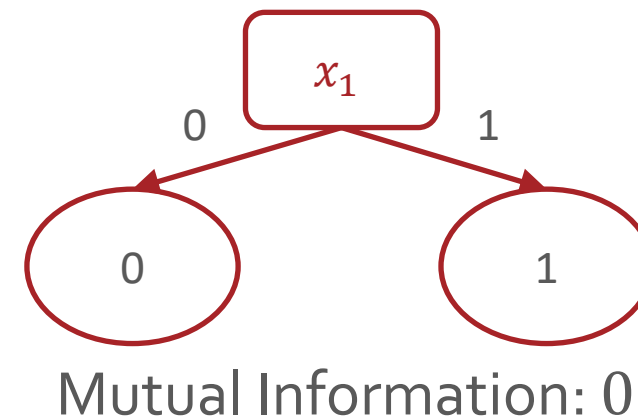


Mutual Information: $H(Y) - \frac{1}{2}H(Y_{x_2=0}) - \frac{1}{2}H(Y_{x_2=1})$

Mutual Information as a Splitting Criterion

x_1	x_2	y
1	0	0
1	0	0
1	0	1
1	0	1
1	1	1
1	1	1
1	1	1
1	1	1

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



Mutual Information: $H(Y) - \left(-\frac{2}{8} \log_2 \frac{2}{8} - \frac{6}{8} \log_2 \frac{6}{8} \right) - \frac{1}{2}(1) - \frac{1}{2}(0) \approx 0.31$

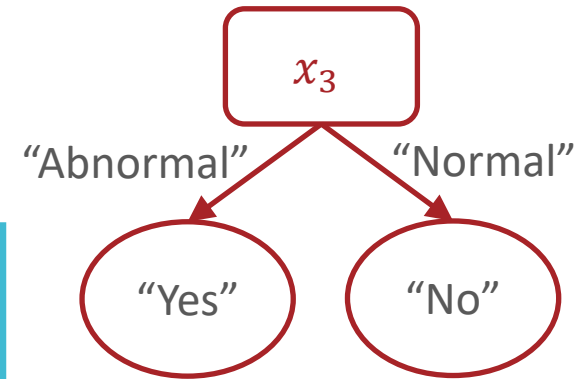
Decision Stumps: Questions

1. How can we pick which feature to split on?
2. Why stop at just one feature?

From Decision Stump

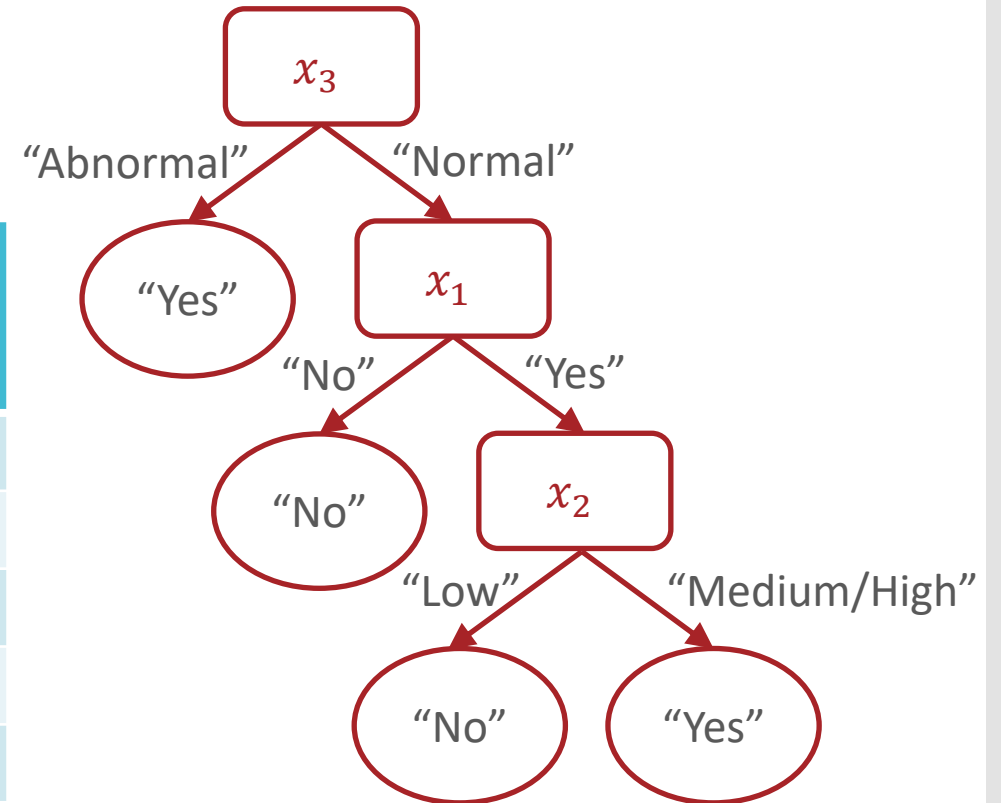
...

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



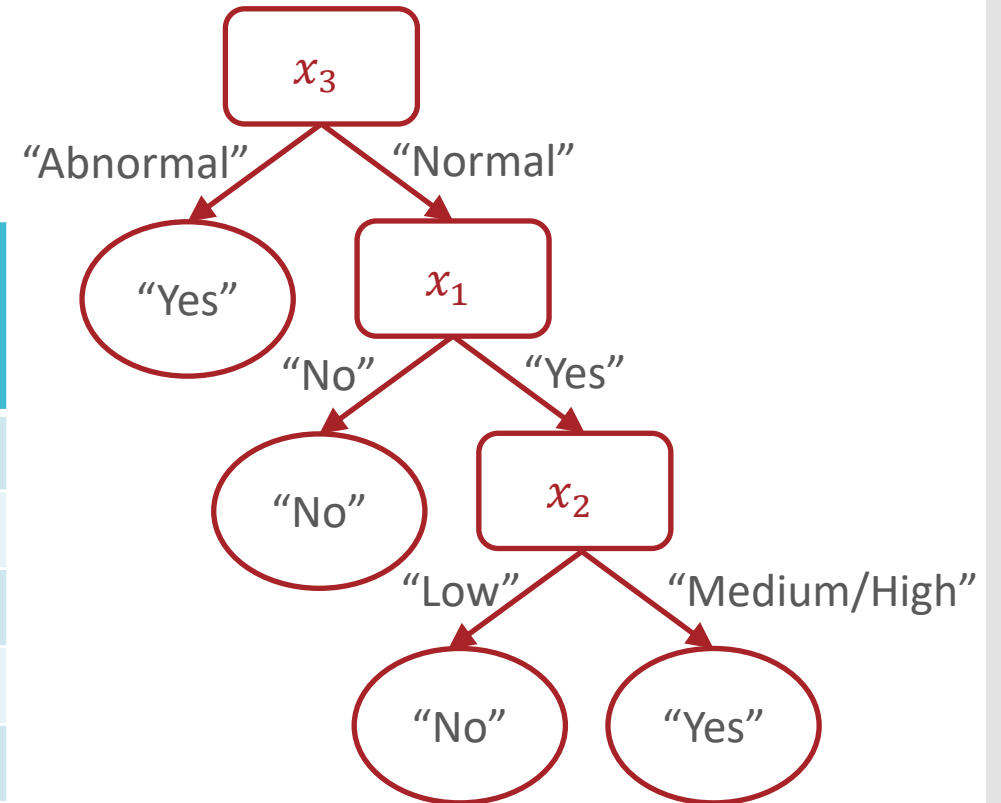
From Decision Stump to Decision Tree

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



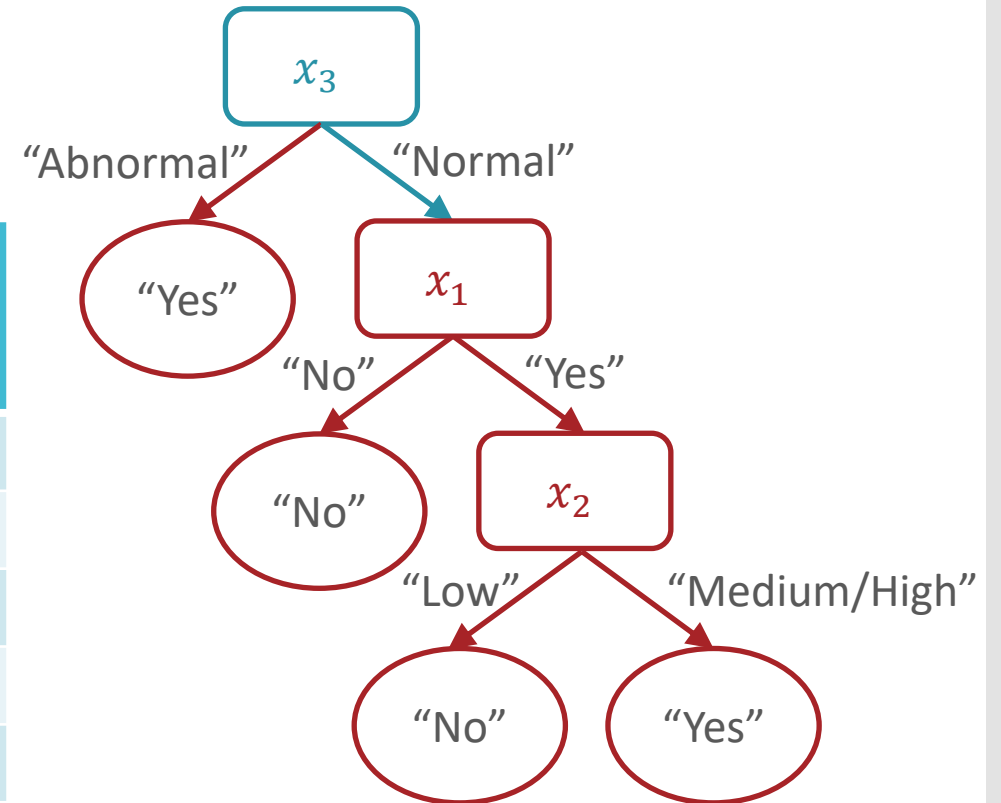
From Decision Stump to Decision Tree

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y 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



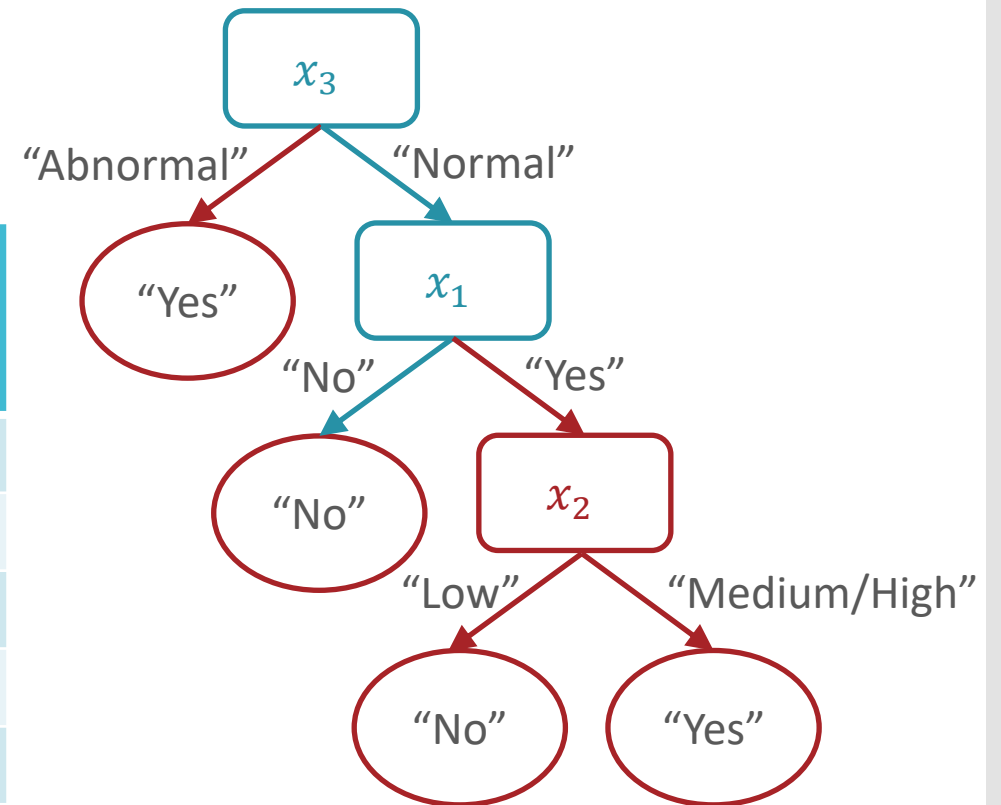
From Decision Stump to Decision Tree

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y 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



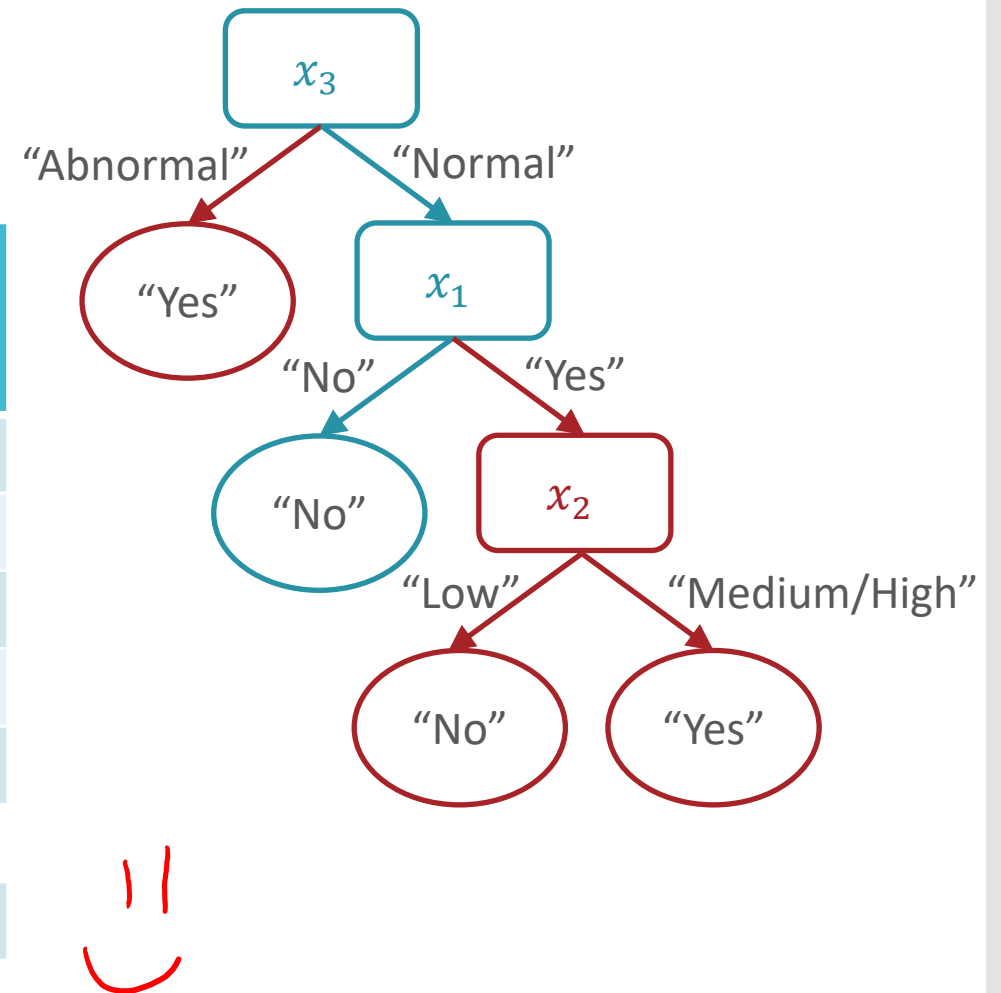
From Decision Stump to Decision Tree

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y 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



From Decision Stump to Decision Tree

x_1 Family History	x_2 Resting Blood Pressure	x_3 Cholesterol	y 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



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-

→ 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-

→ Previous_Csection = 1: [55+,35-] .61+ .39-

{ - Fetal_Presentation = 2: [3+,29-] .11+ .89-

- Fetal_Presentation = 3: [8+,22-] .27+ .73-

↖ leaf nodes

Decision Tree: Pseudocode

```
def predict( $\underline{x'}$ ):
```

```
- walk from root node to a leaf node
```

```
while (true)
```

```
    if current node is not a leaf
```

```
        check the associated feature,  $x_d$ 
```

```
        go down the branch corresponding to  $x_d$ 
```

```
    else
```

```
        return the label stored at that leaf
```

Decision Tree: Pseudocode

```
def train( $\mathcal{D}$ ):
```

Decision Tree: Pseudocode

```
def train( $\mathcal{D}$ ):  
    store root = tree_recurse( $\mathcal{D}$ )  
  
def tree_recurse( $\mathcal{D}'$ ):  
    q = new node()  
    base case – if (SOME CONDITION):  
    recursion – else:
```

find the best attribute to split on, x_d
 $q.\text{split} = x_d$
for v in $V(x_d)$
 $\mathcal{D}_v = \{(x^{(n)}, y^{(n)}) \in \mathcal{D} \mid x_d^{(n)} = v\}$
 $q.\text{children}(v) = \text{tree_recurse}(\mathcal{D}_v)$
return q

Decision Tree: Pseudocode

```
def train( $\mathcal{D}$ ):
```

```
    store root = tree_recurse( $\mathcal{D}$ )
```

```
def tree_recurse( $\mathcal{D}'$ ):
```

```
    q = new node()
```

```
    base case - if ( $\mathcal{D}'$  is empty or
```

all labels in \mathcal{D} are the same or
all feature values in \mathcal{D} are the
same or MI of all features is low...)

q.prediction = majority_vote(\mathcal{D}')

```
    recursion - else:
```

```
    return q
```

x_1	x_2
1	0
1	0
1	0

Decision Trees: Pros & Cons

- Pros

- Interpretable
- Efficient
- Regression or Classification / work on real-valued or categorical

- Cons

- Greedy \Rightarrow no optimality guarantee in terms of minimizing # of splits
- Overfitting!

Decision Trees: Inductive Bias

- The **inductive bias** of a machine learning algorithm is the principal by which it generalizes to unseen examples
- What is the inductive bias of the ID3 algorithm?

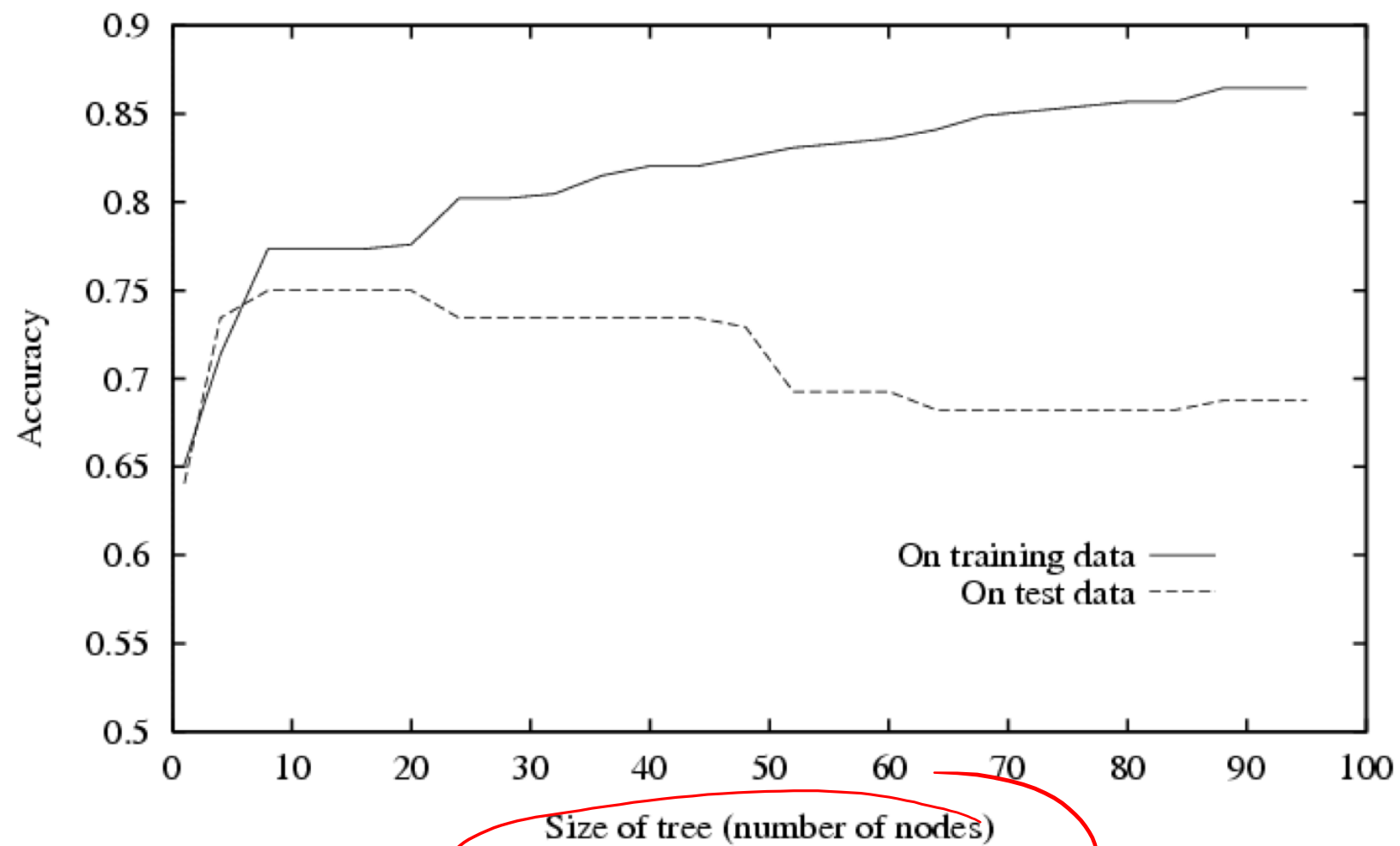
*Greedily find the smallest decision tree
that is consistent w/ the training data
and has high mutual information features at the top.*

- Occam's razor: try to find the "simplest" (e.g., smallest decision tree) classifier that explains the training dataset

Overfitting

- Overfitting occurs when the classifier (or model)...
 - is too complex
 - fits noise or “outliers” in the training dataset as opposed to the actual pattern of interest
 - doesn’t have enough inductive bias pushing it to generalize
- Underfitting occurs when the classifier (or model)...
 - is too simple
 - can’t capture the actual pattern of interest in the training dataset
 - has too much inductive bias

Overfitting in Decision Trees



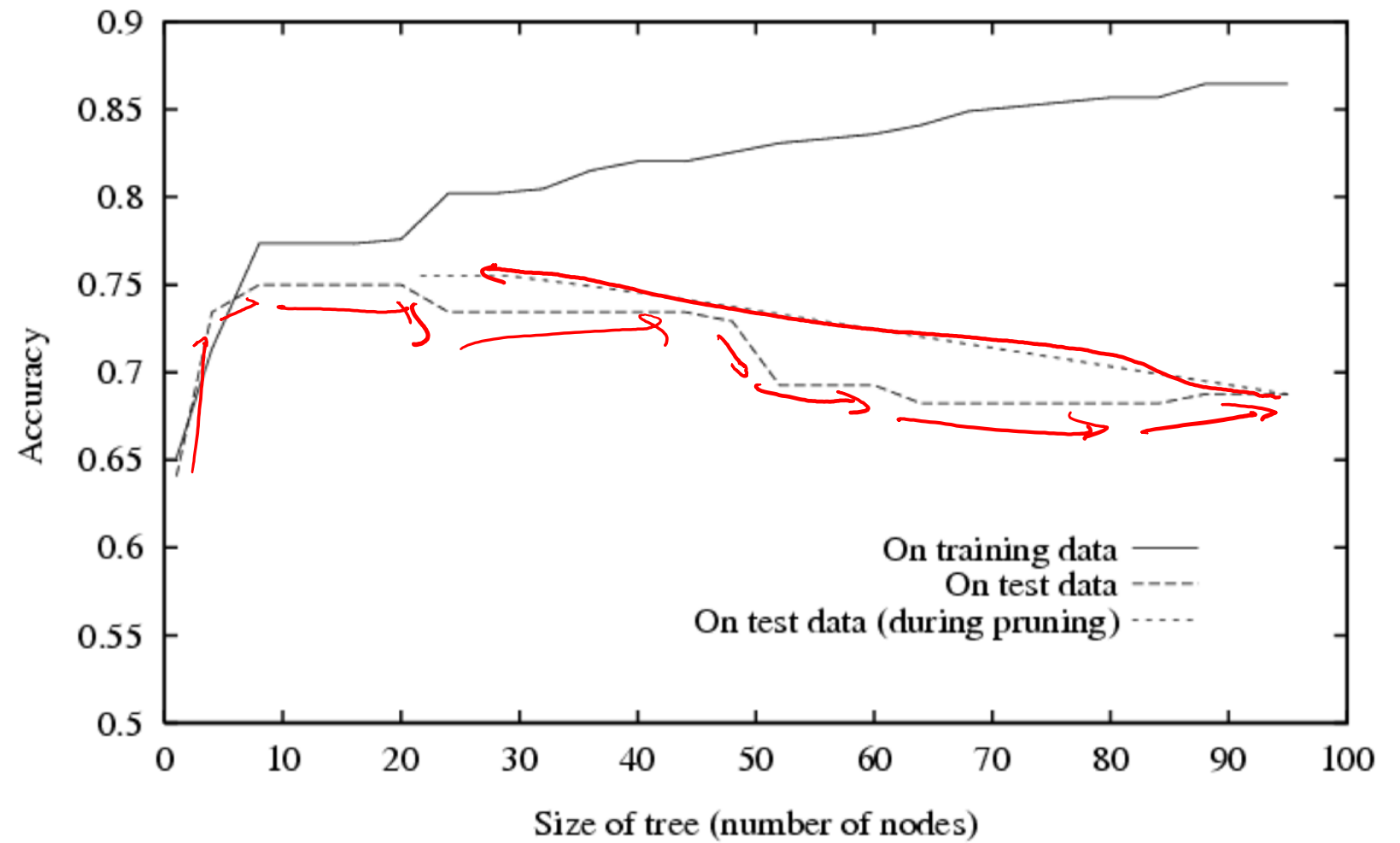
Combatting Overfitting in Decision Trees

- Heuristics:
 - Do not split leaves past a fixed depth, δ
 - Do not split leaves with fewer than C data points
 - Do not split leaves where the maximal information gain is less than τ
 - Take a majority vote in impure leaves

Combatting Overfitting in Decision Trees

- Pruning:
 - First, learn a decision tree
 - Then, evaluate each split using a “validation” dataset by comparing the validation error rate with and without that split
 - Greedily remove the split that most decreases the validation error rate
 - Stop if no split is removed

Pruning Decision Trees



Key Takeaways

- Mutual information as a splitting criterion for decision stumps/trees
- Decision tree algorithm via recursion
- Inductive bias of decision trees
- Overfitting vs. Underfitting
- How to combat overfitting in decision trees