

# Decision Trees

CMSC 422

SOHEIL FEIZI

[sfeizi@cs.umd.edu](mailto:sfeizi@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 loss


# Ideas for Final Project

Talk tomorrow 11-12 same classroom

**Deep Learning Foundations: Interpretability, Robustness and Generative Models**

Soheil Feizi - University of Maryland

3117 Computer Science Instructional Center (CSI)

Friday, February 8, 2019, 11:00 am-12:00 pm 

# Machine Learning as Function Approximation

## Problem setting

- Set of possible instances  $X$
- Unknown target function  $f: X \rightarrow Y$
- Set of function hypotheses  $H = \{h \mid h: X \rightarrow Y\}$

## Input

- Training examples  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$  of unknown target function  $f$

## Output

- Hypothesis  $h \in H$  that best approximates target function  $f$

# Formalizing induction: Loss Function

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

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

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 sample from it: our training data

# Formalizing induction: Expected loss

- $h$  should make good predictions
  - as measured by loss  $l$
  - on **future** examples that are also drawn from  $D$
- Formally
  - $\varepsilon$ , the expected loss of  $f$  over  $D$  should be small

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

# 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  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$
- All we can compute is the training error

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

Empirical Risk Minimization (ERM)



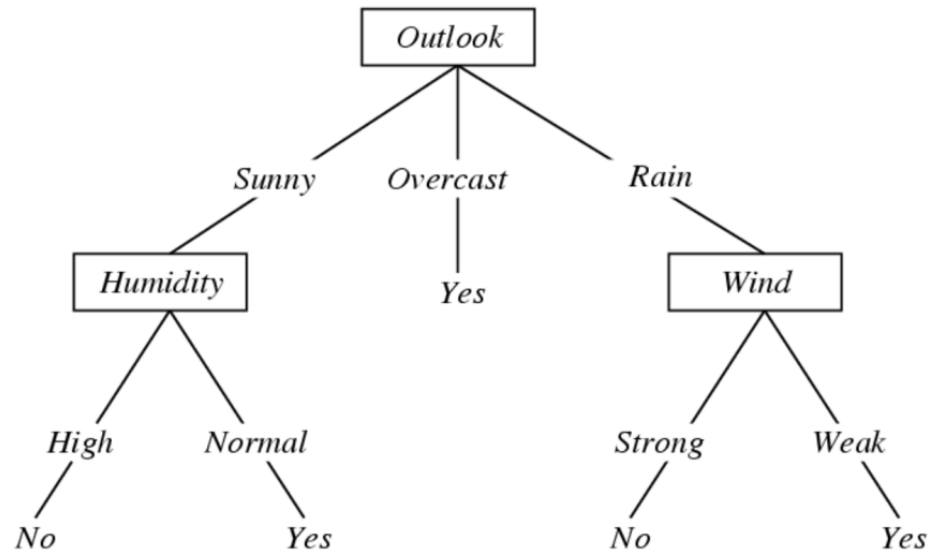
# Today: Decision Trees

- **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> → PlayTennis?

# Exercise

- How would you represent the following Boolean functions with decision trees?
  - AND
  - OR
  - XOR

# 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  $x = [x_1, \dots, x_D]$
- 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  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(N)}, y^{(N)})\}$  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):

1. Find F, the “best” decision feature for next node
2. 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 W, H, T

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
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
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

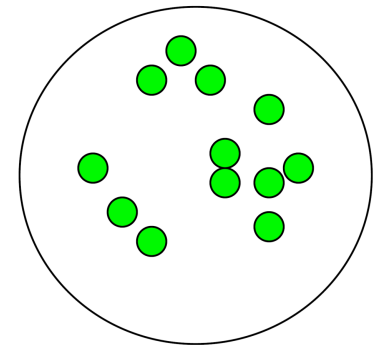
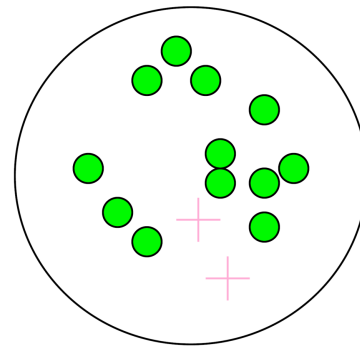
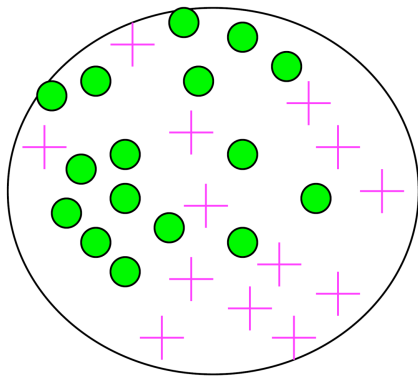
# Partitioning examples:

H = Normal and W = Strong

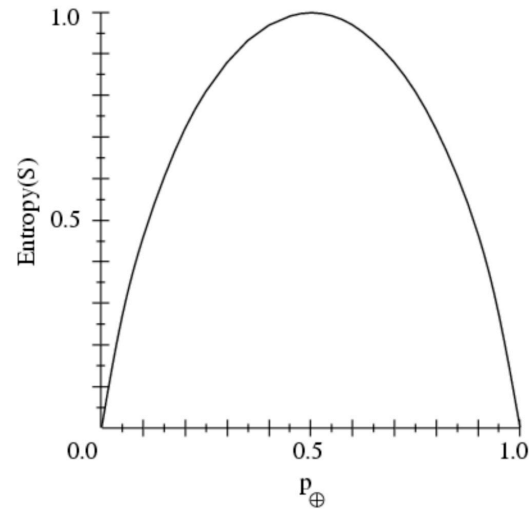
Day	Outlook	Temperature	Humidity	Wind	PlayTennis?
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal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

# 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



# Sample Entropy



- $S$  is a sample of training examples
- $p_+$  is the proportion of positive examples in  $S$
- $p_-$  is the proportion of negative examples in  $S$
- Entropy measures the impurity of  $S$

$$H(S) \equiv -p_+ \log_2 p_+ - p_- \log_2 p_-$$



# Entropy

# of possible values for X

Entropy  $H(X)$  of a random variable  $X$

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

$H(X)$  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  $-\log_2 P(X=i)$  bits to encode the message  $X=i$
- So, expected number of bits to code one random  $X$  is:

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

# Conditional Entropy

Conditional Entropy  $H(Y|X)$  of a random variable  $Y$  conditioned on a random variable  $X$

$$H(Y|X) = - \sum_{j=1}^v P(X = x_j) \sum_{i=1}^k P(Y = y_i | X = x_j) \log_2 P(Y = y_i | X = x_j)$$

Example:

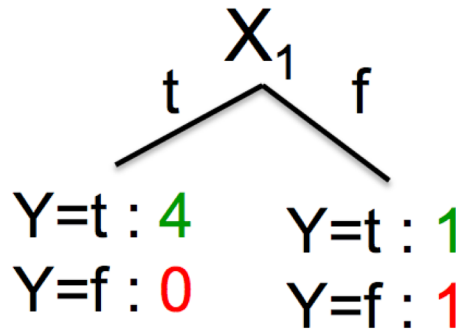
$$P(X_1=t) = 4/6$$

$$P(X_1=f) = 2/6$$

$$H(Y|X_1) = - 4/6 (1 \log_2 1 + 0 \log_2 0)$$

$$- 2/6 (1/2 \log_2 1/2 + 1/2 \log_2 1/2)$$

$$= 2/6$$



$X_1$	$X_2$	$Y$
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F

# Information gain

- Decrease in entropy (uncertainty) after splitting

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

In our running example:

$$\begin{aligned} IG(X_1) &= H(Y) - H(Y|X_1) \\ &= 0.65 - 0.33 \end{aligned}$$

$IG(X_1) > 0 \rightarrow$  we prefer the split!

$X_1$	$X_2$	Y
T	T	T
T	F	T
T	T	T
T	F	T
F	T	T
F	F	F

# Today: Decision Trees

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

# Inductive bias in decision tree learning

CurrentNode = Root

DTtrain(examples for CurrentNode, features at CurrentNode):

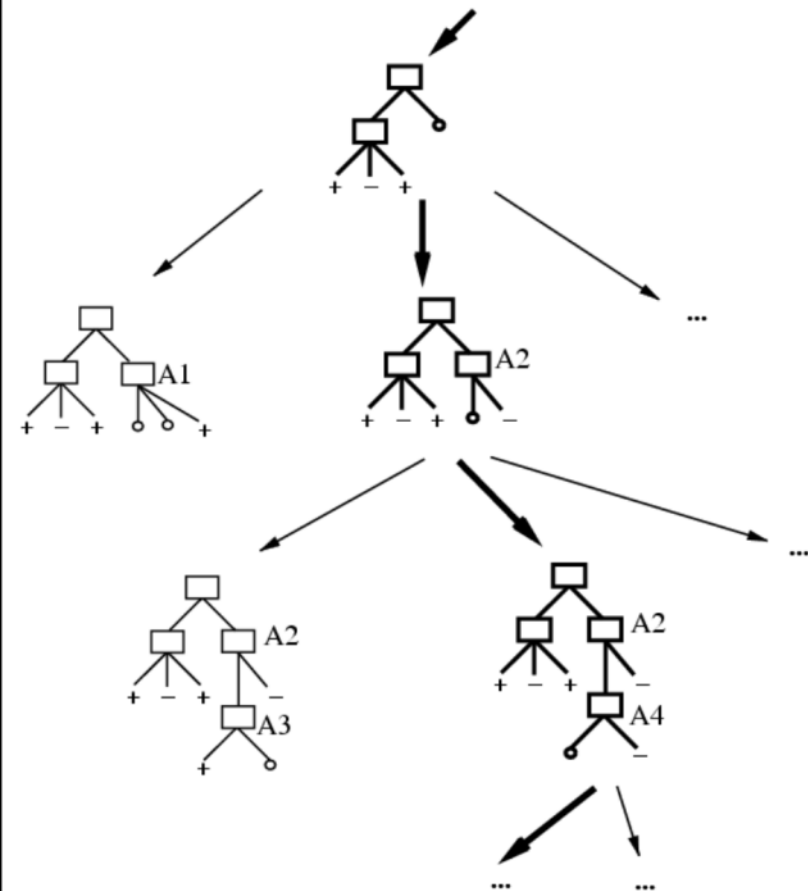
1. Find F, the “best” decision feature for next node
2. For each value of F, create new descendant of node
3. Sort training examples to leaf nodes
4. If training examples perfectly classified

    Stop

Else

    Recursively apply DTtrain over new leaf nodes

# Inductive bias in decision tree learning



- Our learning algorithm performs heuristic search through space of decision trees
- It stops at smallest acceptable tree
- Occam's razor: prefer the simplest hypothesis that fits the data

# Why prefer short hypotheses?

- Pros
  - Fewer short hypotheses than long ones
    - A short hypothesis that fits the data is less likely to be a statistical coincidence
- Cons
  - What's so special about short hypotheses?

# Evaluating the learned hypothesis $h$

- Assume
  - we've learned a tree  $h$  using the top-down induction algorithm
  - It fits the training data perfectly
- Are we done? Can we guarantee we have found a good hypothesis?



# Recall: 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  $l$ .

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

# Training error is not sufficient

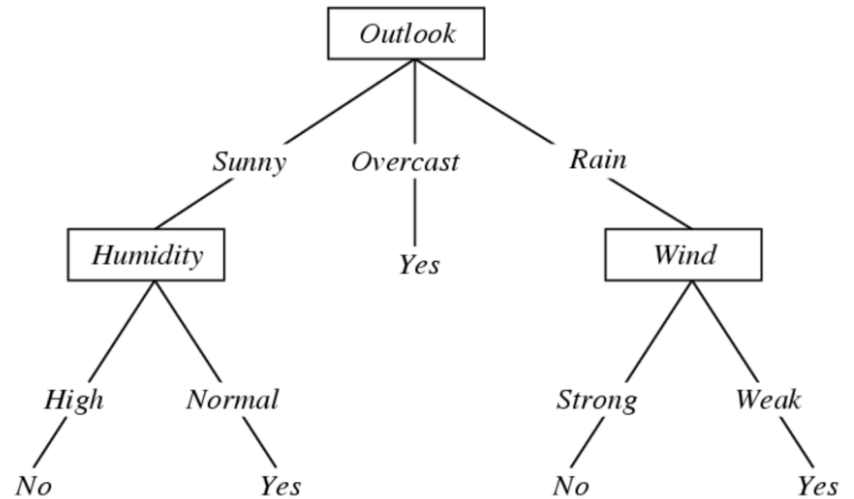
- We care about **generalization** to new examples
- A tree can classify training data perfectly, yet classify new examples incorrectly
  - Because training examples are only a sample of data distribution
    - a feature might correlate with class by coincidence
  - Because training examples could be noisy
    - e.g., accident in labeling

# Let's add a noisy training example.

## How does this affect the learned decision tree?

**Day Outlook Temperature Humidity Wind**

D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	Yes
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	Yes
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	Yes
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
<b>D15</b>	<b>Sunny</b>	<b>Hot</b>	<b>Normal</b>	<b>Strong</b>	<b>No</b>



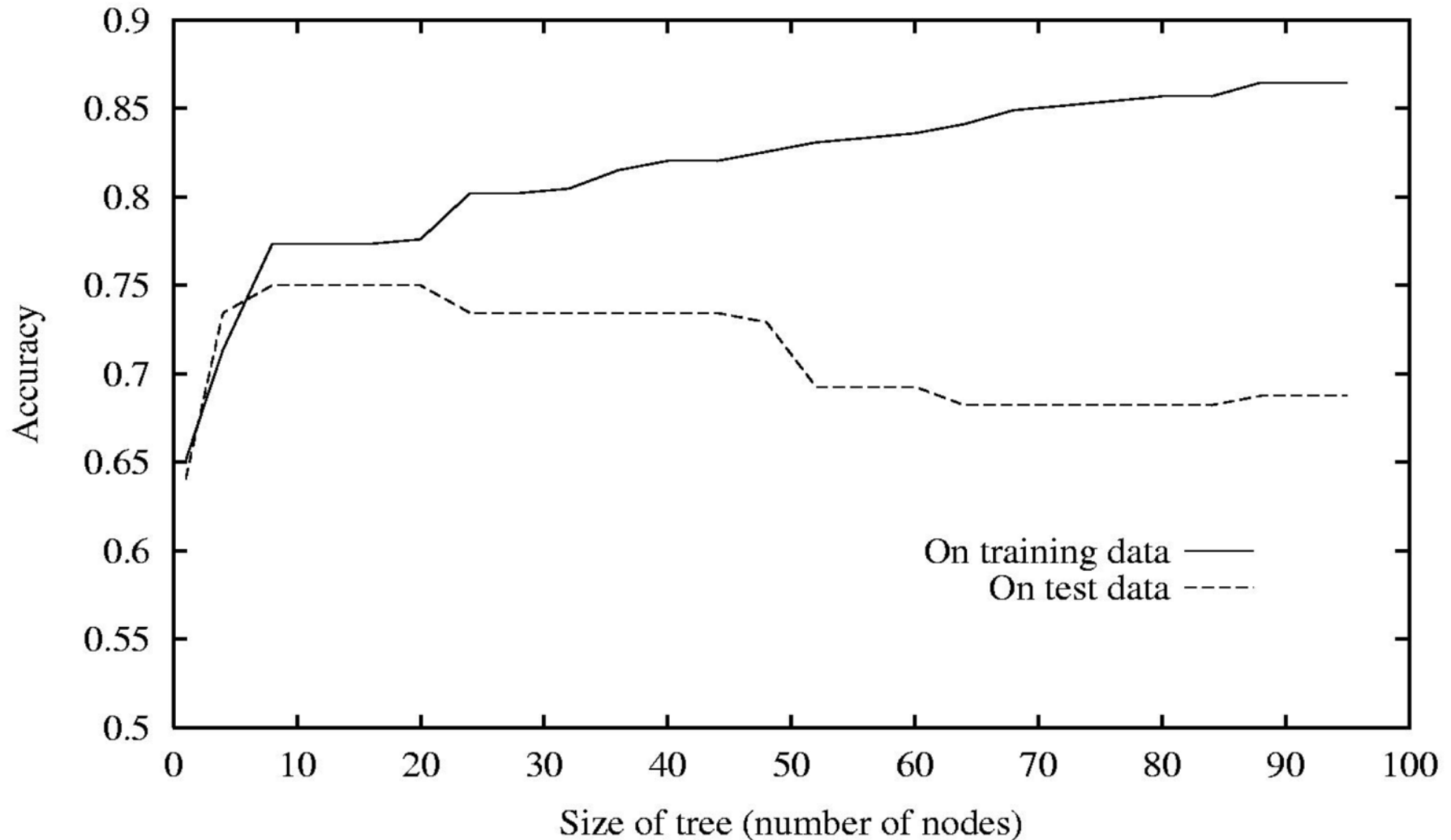
# Overfitting

- Consider a hypothesis  $h$  and its:
  - Error rate over training data  $error_{train}(h)$
  - True error rate over all data  $error_{true}(h)$
- We say  $h$  overfits the training data if
$$error_{train}(h) < error_{true}(h)$$
- Amount of overfitting =
$$error_{true}(h) - error_{train}(h)$$

# Evaluating on test data

- Problem: we don't know  $error_{true}(h)$ !
- Solution:
  - we set aside a test set
    - some examples that will be used for evaluation
  - we don't look at them during training!
  - after learning a decision tree, we calculate  $error_{test}(h)$

# Measuring effect of overfitting in decision trees



# Underfitting/Overfitting

- Underfitting
  - Learning algorithm had the opportunity to learn more from training data, but didn't
- Overfitting
  - Learning algorithm paid too much attention to learn noisy part of the training data; the resulting tree doesn't generalize
- What we want:
  - A decision tree that neither underfits nor overfits
  - Because it is expected to do best in the future

# Today: Decision Trees

- What is a decision tree?
- How to learn a decision tree from data?
  - Top-down induction to minimize classification error
- What is the inductive bias?
  - Occam's razor: preference for short trees
- Generalization?
  - Overfitting can be an issue