

DS 4400

Machine Learning and Data Mining I
Spring 2024

David Liu

Khoury College of Computer Science
Northeastern University

February 27 2024

Announcements

- Final Project
 - Project proposals due this Friday on Gradescope
 - Each member should submit a copy.
 - Expect further details on final project report (6-8 pages) and video (~4 minutes) after spring break.
- Homework 3
 - Released on February 26, due on March 18
- Midterm grades will be released after Spring Break.

Where We Are

Regression

Simple, Multiple,
Polynomial

Classification

Logistic Regression,
kNN

Generative Models

LDA, Naïve Bayes

Ensemble Learning

Bagging, Boosting

Neural Networks

Multi-layer Perceptron,
Convolutional Neural
Networks

Tools

Optimization:

- Closed form solution
- Gradient descent
- Backpropagation

Preventing Overfitting:

- Regularization
- Cross Validation

Evaluation of ML:

- Precision, Recall, AUC

Outline

- Generative classifiers
- Naïve Bayes classifiers [focus for today]
 - Naïve Bayes assumption
 - Laplace smoothing
 - Comparison to LDA
- Decision trees [preview]
 - Information gain / entropy measures
 - Training algorithm

Recap: ML Concepts

- Supervised vs unsupervised learning
- Classification vs regression
- Linear vs non-linear classifiers
- Generative vs discriminative classifiers
- Loss functions
- Metrics for evaluation

Generative vs Discriminative

- **Generative model**
 - Given X and Y , learns the joint probability $P(X, Y)$
 - Can generate more examples from distribution
 - Examples: LDA, Naïve Bayes, language models (GPT-2, GPT-3, BERT)
- **Discriminative model**
 - Given X and Y , learns a decision function for classification

Generative classifiers based on Bayes Theorem

$$P(A | B) = \frac{P(B | A) \times P(A)}{P(B)}$$

- Exactly the process we just used
- The most important formula in probabilistic machine learning

(Super Easy) Derivation:

$$P(A \wedge B) = P(A | B) \times P(B)$$

$$P(B \wedge A) = P(B | A) \times P(A)$$

these are the same

Just set equal...

$$P(A | B) \times P(B) = P(B | A) \times P(A)$$

and solve...



Bayes, Thomas (1763) An essay towards solving a problem in the doctrine of chances. *Philosophical Transactions of the Royal Society of London*, 53:370-418

LDA Training and Testing

Given training data $(x_i, y_i), i = 1, \dots, n, y_i \in \{1, \dots, K\}$

1. Estimate mean and variance

$$\begin{aligned}\hat{\mu}_k &= \frac{1}{n_k} \sum_{i:y_i=k} x_i \\ \hat{\sigma}^2 &= \frac{1}{n-K} \sum_{k=1}^K \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2\end{aligned}$$

2. Estimate prior

$$\hat{\pi}_k = n_k / n.$$

Given testing point x , predict k that maximizes:

$$\hat{\delta}_k(x) = x \cdot \frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log(\hat{\pi}_k)$$

Naïve Bayes Classifier

Idea: Use the training data to estimate

$$P(X | Y) \text{ and } P(Y) .$$

Then, use Bayes rule to infer $P(Y|X_{\text{new}})$ for new data

Easy to estimate
from data

Impractical, but necessary

$$P[Y = k|X = x] = \frac{P[Y = k]P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d|Y = k]}{P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d]}$$

Unnecessary, as it turns out

Learning Joint Distributions

Step 1:

Build a JD table for your attributes in which the probabilities are unspecified

A	B	C	Prob
0	0	0	?
0	0	1	?
0	1	0	?
0	1	1	?
1	0	0	?
1	0	1	?
1	1	0	?
1	1	1	?

Step 2:

Then, fill in each row with:








$$\hat{P}(\text{row}) = \frac{\text{records matching row}}{\text{total number of records}}$$

A	B	C	Prob
0	0	0	0.30
0	0	1	0.05
0	1	0	0.10
0	1	1	0.05
1	0	0	0.05
1	0	1	0.10
1	1	0	0.25
1	1	1	0.10

Fraction of all records in which
A and B are true but C is false

Example – Learning Joint Probability Distribution

This Joint PD was obtained by learning from three attributes in the UCI “Adult” Census Database [Kohavi 1995]

gender	hours_worked	wealth		
Female	v0:40.5-	poor	0.253122	
		rich	0.0245895	
	v1:40.5+	poor	0.0421768	
		rich	0.0116293	
Male	v0:40.5-	poor	0.331313	
		rich	0.0971295	
	v1:40.5+	poor	0.134106	
		rich	0.105933	

Naïve Bayes Classifier

Problem: estimating the joint density isn't practical

However, if we make the assumption that the attributes are independent given the class label, estimation is easy!

Naïve Bayes Classifier

Problem: estimating the joint PD or CPD isn't practical
– Severely overfits, as we saw before

However, if we make the assumption that the attributes are independent given the class label, estimation is easy!

$$P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d | Y = k] = \prod_{j=1}^d P[X_j = x_j | Y = k]$$

- In other words, we assume all attributes are *conditionally independent* given Y
- Often this assumption is violated in practice, but more on that later...

Using the Naïve Bayes Classifier

- Now, we have

$$P[Y = k|X = x] = \frac{P[Y = k]P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d|Y = k]}{P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d]}$$

This is constant for a given instance,
and so irrelevant to our prediction

Using the Naïve Bayes Classifier

- Now, we have

$$P[Y = k|X = x] = \frac{P[Y = k]P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d|Y = k]}{P[X_1 = x_1 \wedge \cdots \wedge X_d = x_d]}$$

This is constant for a given instance,
and so irrelevant to our prediction

- In practice, we use log-probabilities to prevent underflow

- To classify a new point \mathbf{x} ,

$$h(\mathbf{x}) = \arg \max_{y_k} P(Y = k) \prod_{j=1}^d P(X_j = x_j | Y = k)$$

$\underbrace{x_j}_{j^{\text{th}} \text{ attribute value of } \mathbf{x}}$

Naïve Bayes Classifier

TRAIN

- For each class label k
 1. Estimate prior $\pi_k = P[Y = k]$ from the data
 2. For each value v of attribute X_j
 - Estimate $P[X_j = v | Y = k]$

TEST on INPUT $x = (x_1, \dots, x_d)$

- For every k , compute the probabilities
 - $p_k = P[Y = k] \prod_{j=1}^d P[X_j = x_j | Y = k]$
- Classify x to the class k that maximizes p_k

Training Naïve Bayes

Estimate $P[X_j = x_j | Y = k]$ and $P[Y = k]$ directly from the training data by counting!

<u>Sky</u>	<u>Temp</u>	<u>Humid</u>
sunny	warm	normal
sunny	cold	high
rainy	cold	high
sunny	warm	high

<u>Play?</u>
yes
yes
no
yes

Prior: $P(\text{Play} = \text{Yes}) =$

$P(\text{Play} = \text{No}) =$

Conditional feature distributions

$P(\text{Sky} = \text{sunny} | \text{Play} = \text{Yes}) =$

$P(\text{Temp} = \text{warm} | \text{Play} = \text{Yes}) =$

$P(\text{Humid} = \text{high} | \text{Play} = \text{Yes}) =$

$P(\text{Sky} = \text{sunny} | \text{Play} = \text{No}) =$

$P(\text{Temp} = \text{warm} | \text{Play} = \text{No}) =$

$P(\text{Humid} = \text{high} | \text{Play} = \text{No}) =$

Training Naïve Bayes

Estimate $P[X_j = x_j | Y = k]$ and $P[Y = k]$ directly from the training data by counting!

<u>Sky</u>	<u>Temp</u>	<u>Humid</u>
sunny	warm	normal
sunny	cold	high
rainy	cold	high
sunny	warm	high

<u>Play?</u>
yes
yes
no
yes

Classify new point:

X: Sky = sunny, Temp = cold, Humid = high

Laplace Smoothing

- Notice that some probabilities estimated by counting might be zero
 - Possible overfitting!

Laplace Smoothing

- Notice that some probabilities estimated by counting might be zero
 - Possible overfitting!
- Fix by using Laplace smoothing:
 - Adds 1 to each count

$$P(X_j = v \mid Y = k) = \frac{c_v + 1}{\sum_{v' \in \text{values}(X_j)} c_{v'} + |\text{values}(X_j)|}$$

where

- c_v is the count of training instances with a value of v for attribute j and class label k
- $|\text{values}(X_j)|$ is the number of values X_j can take on

Naïve Bayes Classifier

TRAIN

- For each class label k
 1. Estimate prior $\pi_k = P[Y = k]$ from the data
 2. For each value v of attribute X_j
 - Estimate $P[X_j = v | Y = k]$ with Laplace smoothing

TEST on INPUT $x = (x_1, \dots, x_d)$

- For every k , compute the probabilities
 - $P[Y = k] \prod_{j=1}^d P[X_j = x_j | Y = k]$
- Classify x to the class k that maximizes the above product

Continuous Features

- Naïve Bayes can be extended to continuous features
- Gaussian Naïve Bayes
 - Here an additional assumption is that each distribution $P[X_j|Y = k]$ is Gaussian $N(\mu_j, \sigma_j)$
 - It estimates the mean and standard deviation from training data
- This leads to a linear classifier

Comparison to LDA

- **Similarity to LDA**

- Both are generative models
- They both estimate:

$$P[X = x \text{ and } Y = k] = P[X = x|Y = k]P[Y = k]$$

Using Bayes Theorem

- **Difference from LDA**

- Naïve Bayes can handle discrete data
- LDA uses multi-variate normal
- LDA assumes same variances for all classes
- Naïve Bayes make the conditional independence assumption
- LDA is linear, while Naïve Bayes is usually not linear

Naïve Bayes Summary

Advantages:

- Fast to train (single scan through data)
- Fast to classify
- Not sensitive to irrelevant features
- Handles real and discrete data
- Handles streaming data well

Disadvantages:

- Assumes independence of features

Outline

- Generative classifiers
- Naïve Bayes classifiers
 - Naïve Bayes assumption
 - Laplace smoothing
 - Comparison to LDA
- Decision trees
 - Information gain / entropy measures
 - Training algorithm

Sample Dataset

- Columns denote features X_i
- Rows denote labeled instances $\langle x_i, y_i \rangle$
- Class label denotes whether a tennis game was played

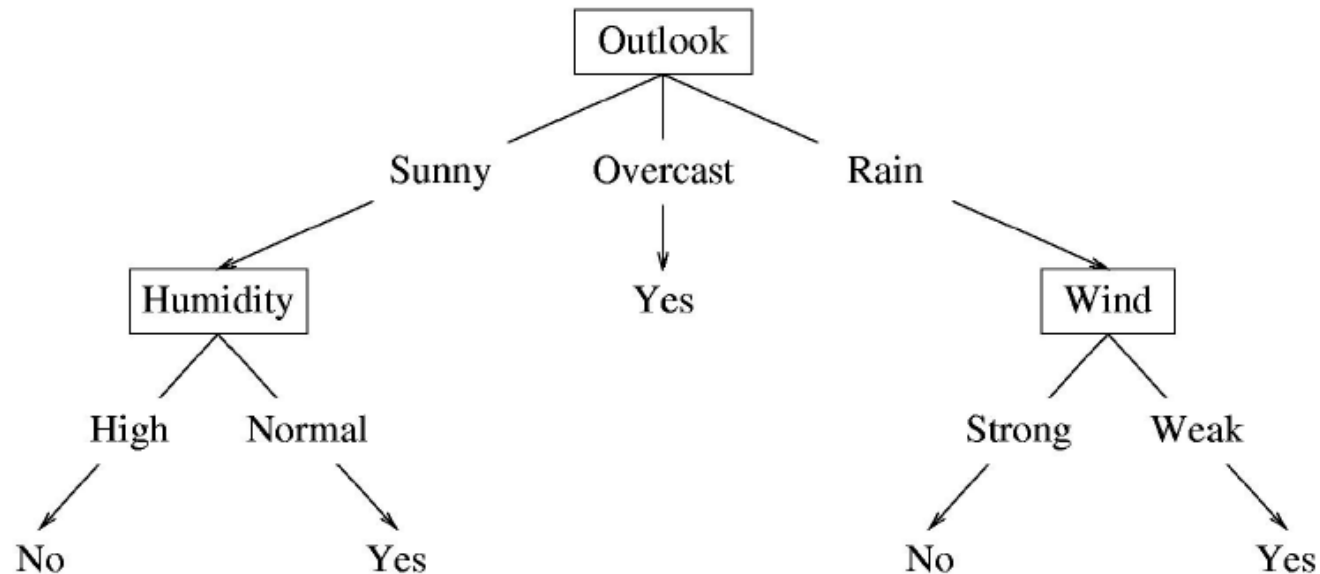
$\langle x_i, y_i \rangle$

Categorical
data

Predictors				Response
Outlook	Temperature	Humidity	Wind	Class
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

Decision Tree

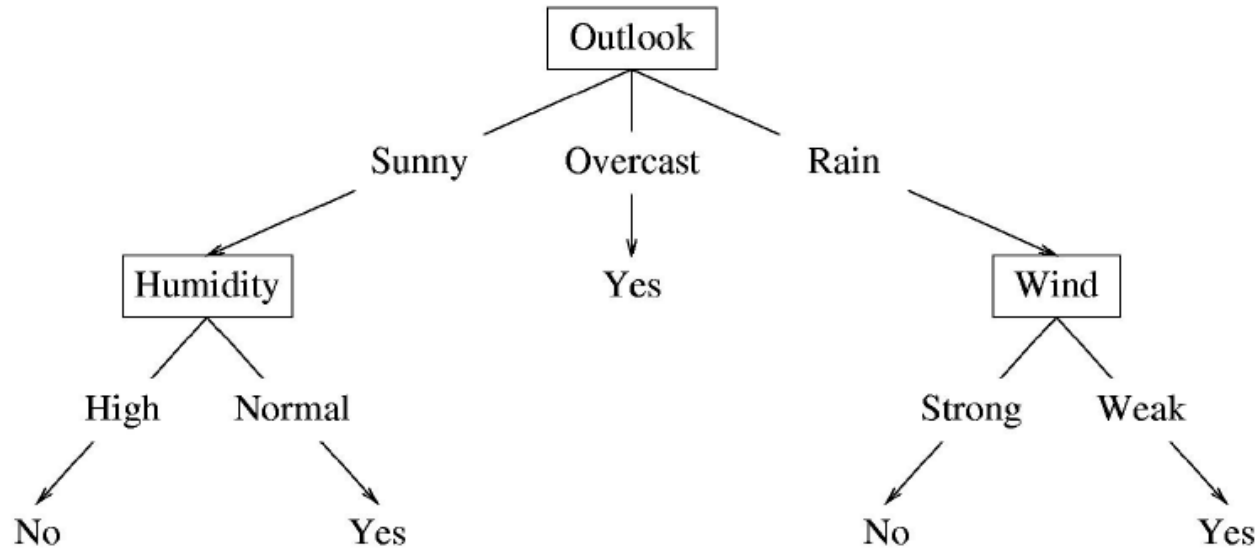
- A possible decision tree for the data:



- Each internal node: test one attribute X_i
- Each branch from a node: selects one value for X_i
- Each leaf node: predict Y (or $p(Y \mid x \in \text{leaf})$)

Decision Tree

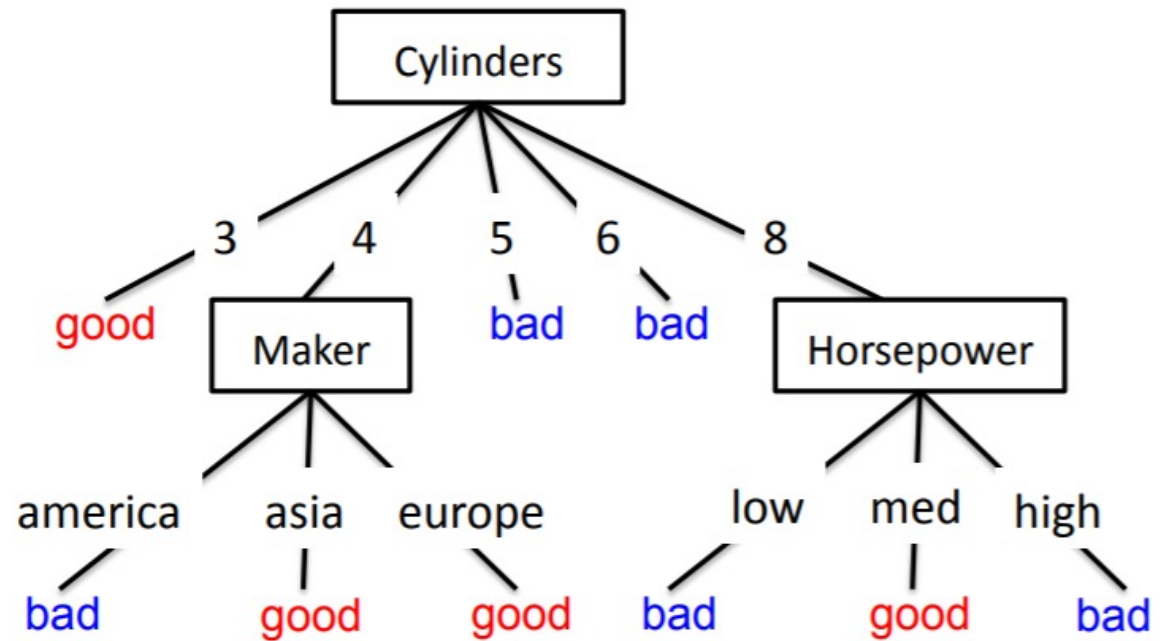
- A possible decision tree for the data:



- What prediction would we make for
<outlook=sunny, temperature=hot, humidity=high, wind=weak> ?

Interpretability

- Each internal node tests an attribute x_i
- One branch for each possible attribute value $x_i=v$
- Each leaf assigns a class y
- To classify input x : traverse the tree from root to leaf, output the labeled y

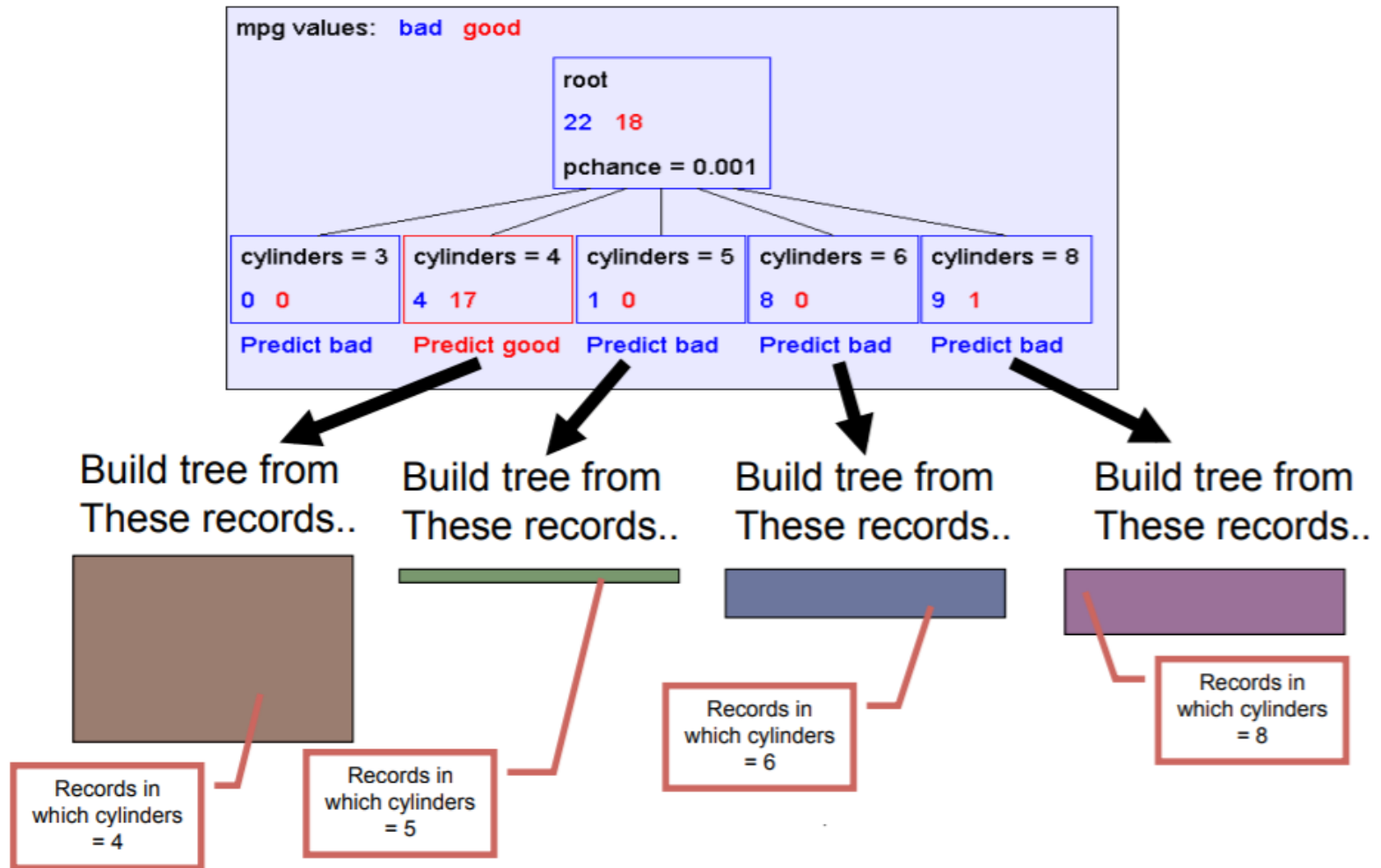


Human interpretable!

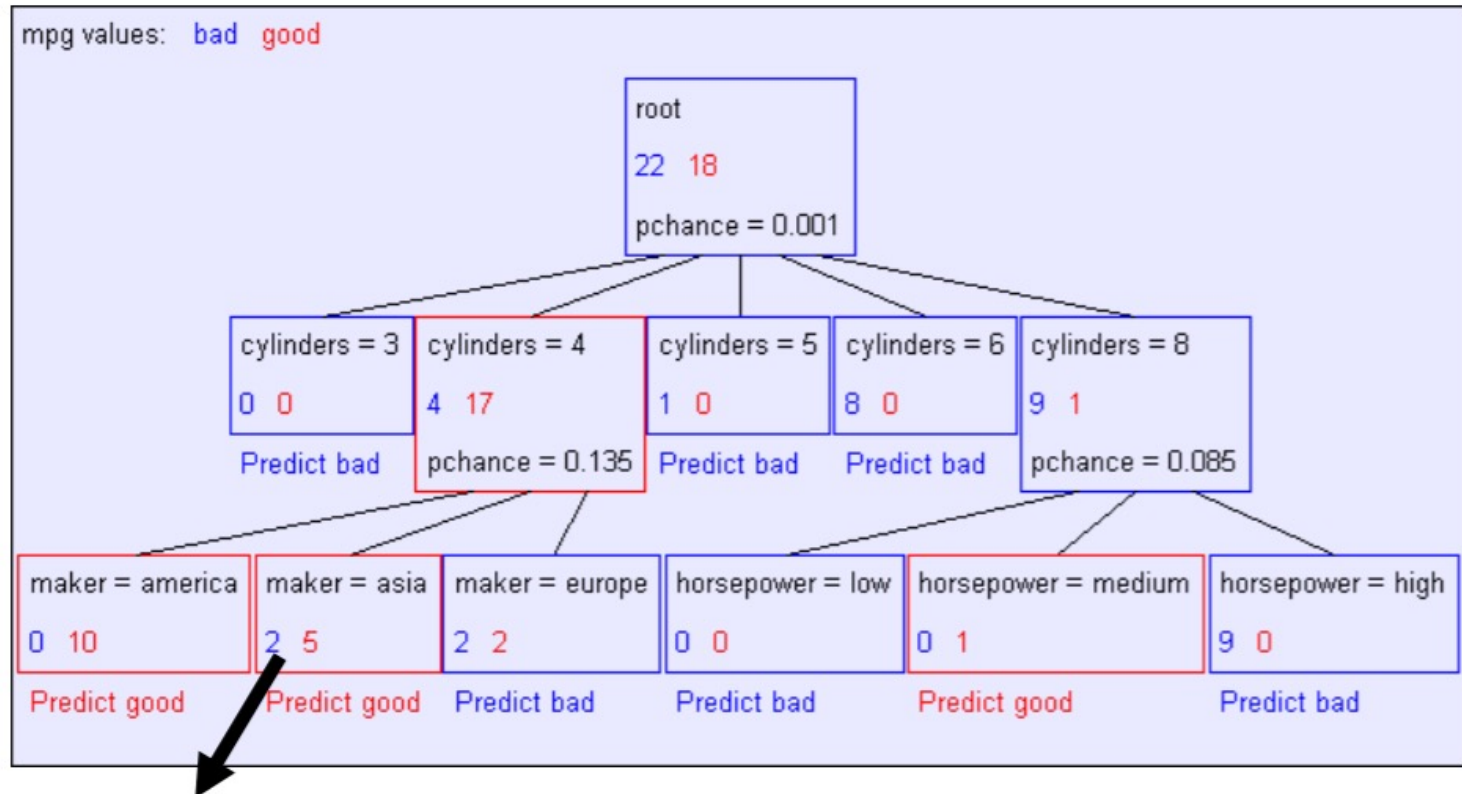
Learning Decision Trees

- Learning the simplest (smallest) decision tree is an NP-complete problem [Hyafil & Rivest '76]
- Resort to a greedy heuristic:
 - Start from empty decision tree
 - Split on **next best attribute (feature)**
 - Recurse

Key Idea: Use Recursion Greedily



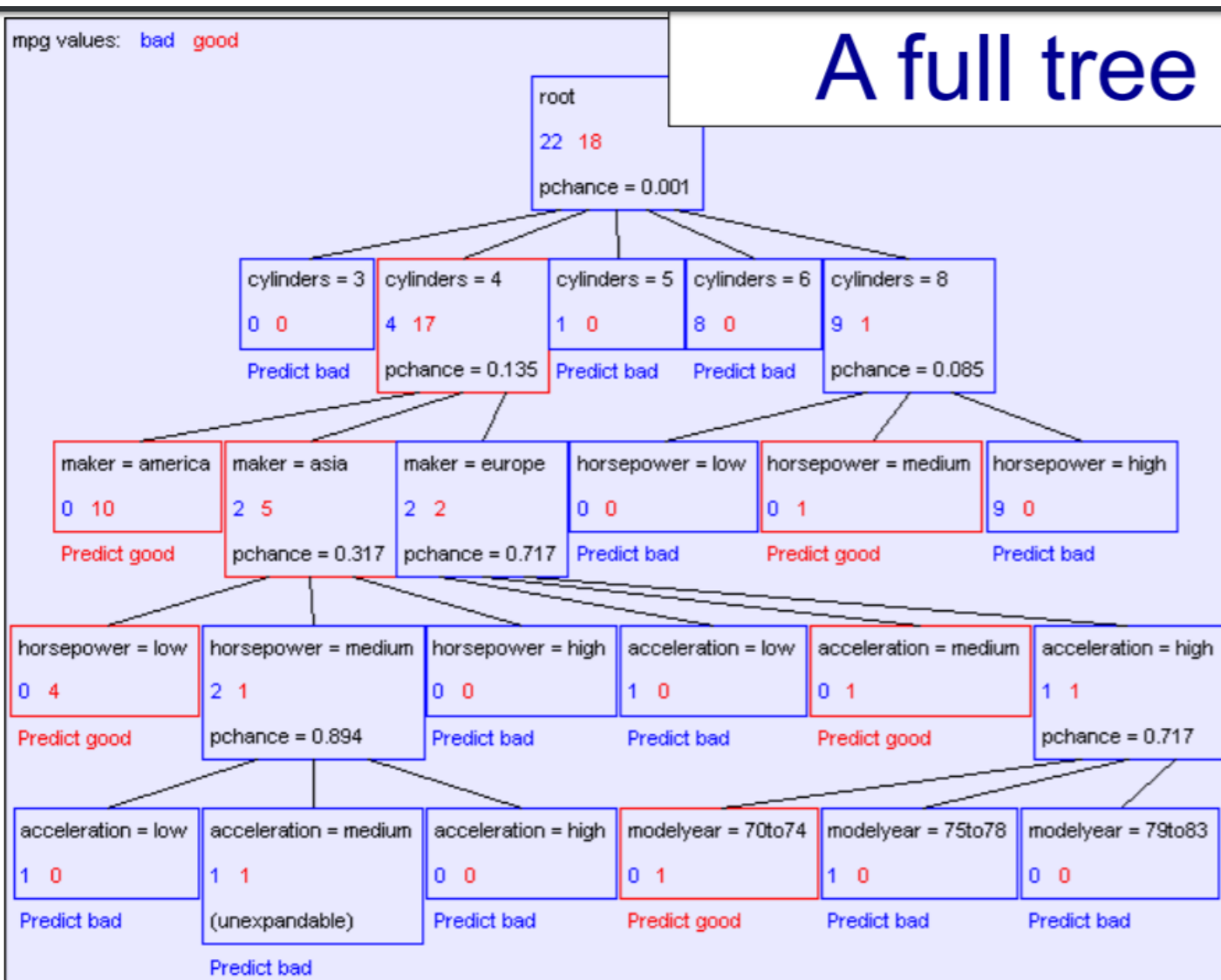
Second Level



Recursively build a tree from the seven records in which there are four cylinders and the maker was based in Asia

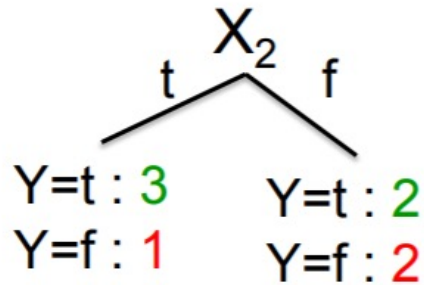
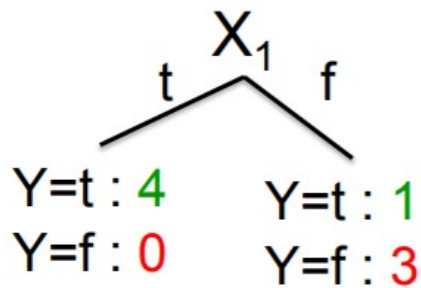
(Similar recursion in the other cases)

Full Tree



Splitting

Would we prefer to split on X_1 or X_2 ?



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

Entropy

Suppose X can have one of m values... V_1, V_2, \dots, V_m

$P(X=V_1) = p_1$	$P(X=V_2) = p_2$	$P(X=V_m) = p_m$
------------------	------------------	------	------------------

What's the smallest possible number of bits, on average, per symbol, needed to transmit a stream of symbols drawn from X 's distribution? It's

$$\begin{aligned} H(X) &= -p_1 \log_2 p_1 - p_2 \log_2 p_2 - \dots - p_m \log_2 p_m \\ &= -\sum_{j=1}^m p_j \log_2 p_j \end{aligned}$$

$H(X)$ = The entropy of X

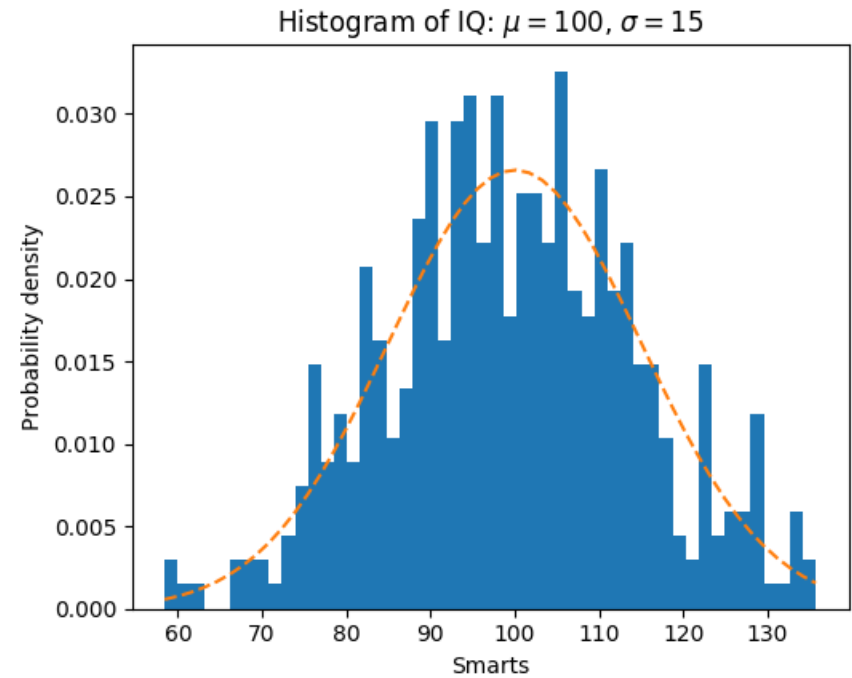
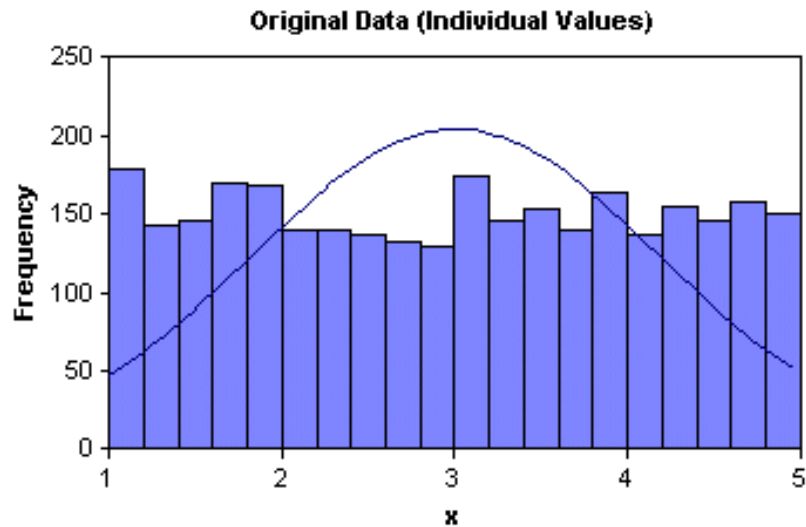
- "High Entropy" means X is from a uniform (boring) distribution
- "Low Entropy" means X is from varied (peaks and valleys) distribution

Entropy Examples

Entropy Examples

High/Low Entropy

Which distribution has high entropy?



Conditional Entropy

Suppose I'm trying to predict output Y and I have input X

X = College Major

Y = Likes "Gladiator"

Let's assume this reflects the true probabilities

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

Conditional Entropy

X = College Major

Y = Likes "Gladiator"

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

Definition of Specific Conditional Entropy:

$H(Y|X=v)$ = **The entropy of Y among only those records in which X has value v**

Example:

- $H(Y|X=Math) =$
- $H(Y|X=History) =$
- $H(Y|X=CS) =$

Conditional Entropy

X = College Major

Y = Likes "Gladiator"

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

Definition of Conditional Entropy:

$H(Y|X)$ = The average specific conditional entropy of Y

= if you choose a record at random what will be the conditional entropy of Y , conditioned on that row's value of X

= Expected number of bits to transmit Y if both sides will know the value of X

$$= \sum_j \text{Prob}(X=v_j) H(Y | X = v_j)$$

Conditional Entropy

X = College Major

Y = Likes "Gladiator"

Definition of Conditional Entropy:

$H(Y|X)$ = The average conditional entropy of Y

$$= \sum_j \text{Prob}(X=v_j) H(Y | X = v_j)$$

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

Example:

v_j	$\text{Prob}(X=v_j)$	$H(Y X = v_j)$
Math		
History		
CS		

Information Gain

X = College Major

Y = Likes "Gladiator"

Definition of Information Gain:

$IG(Y|X) =$ **I must transmit Y .**

How many bits on average would it save me if both ends of the line knew X ?

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

X	Y
Math	Yes
History	No
CS	Yes
Math	No
Math	No
CS	Yes
History	No
Math	Yes

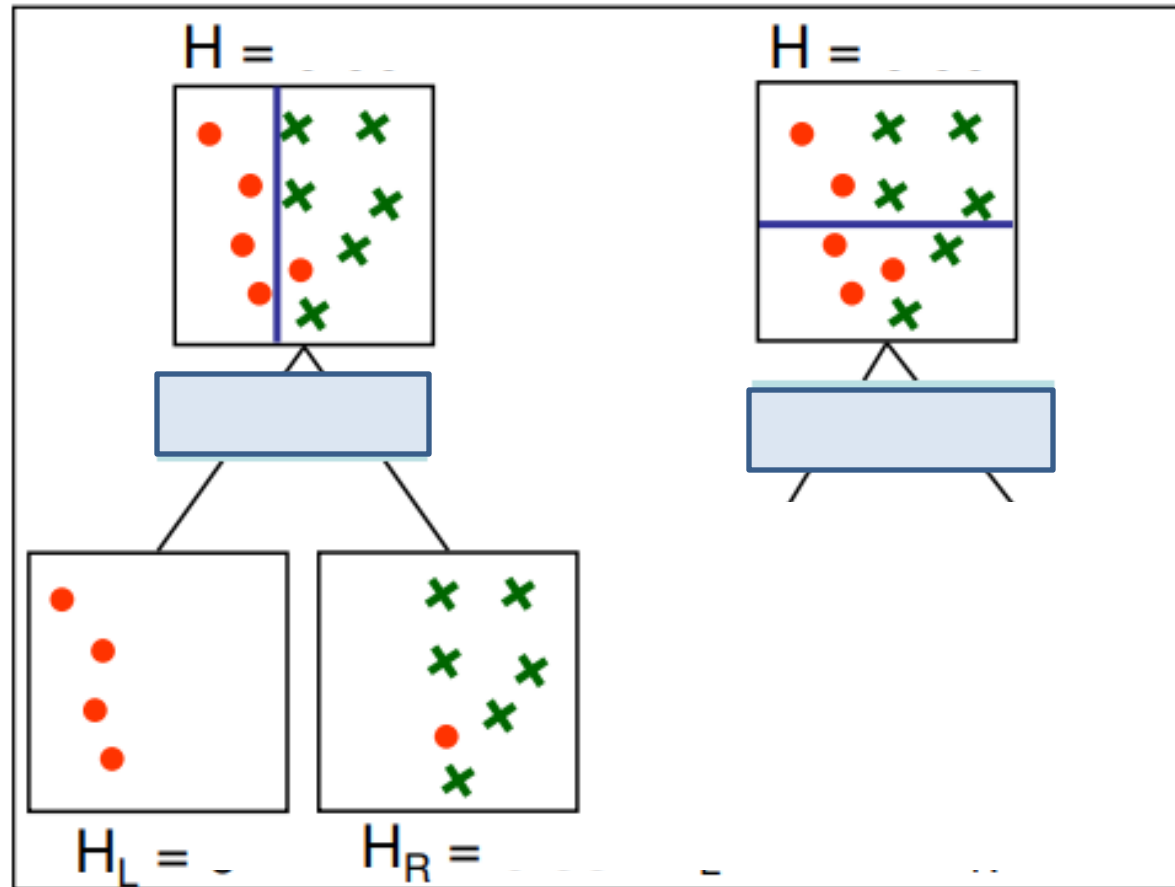
Example:

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

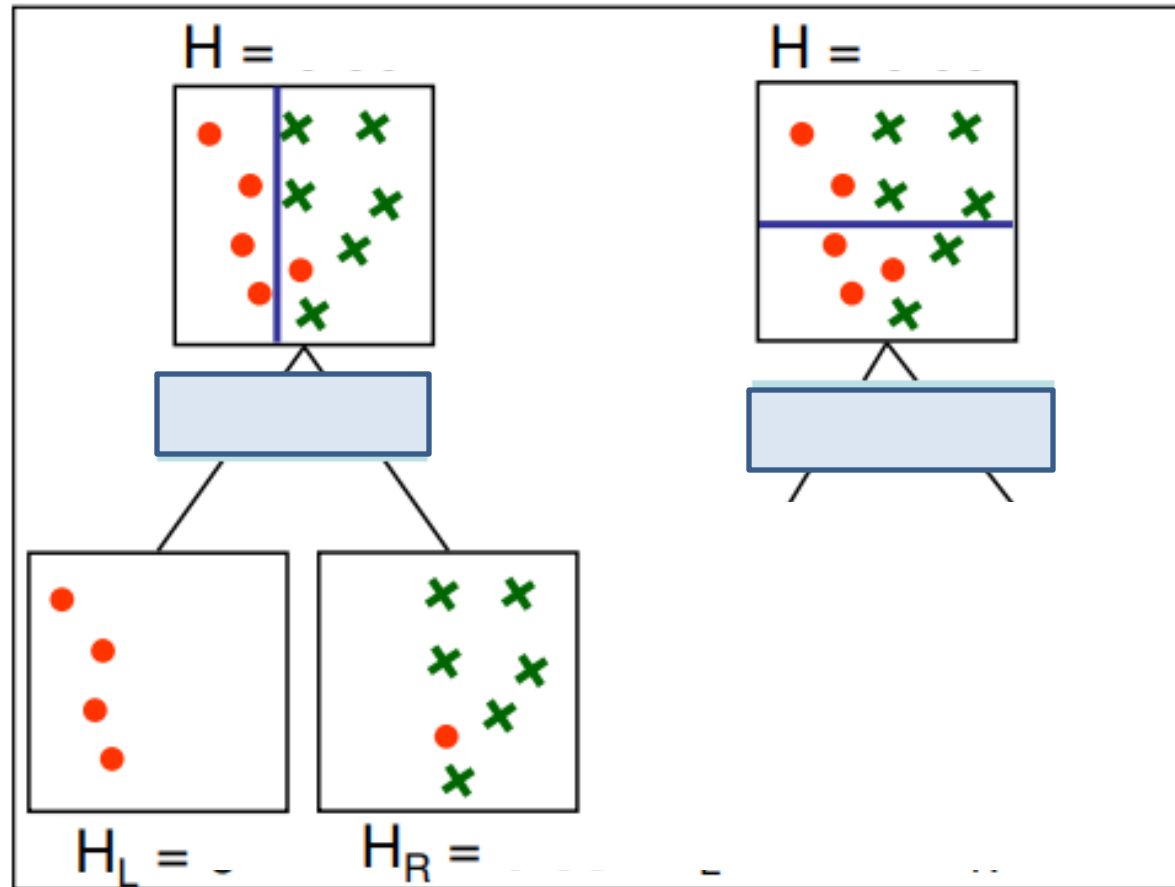
Relevance for decision trees

- Multiple features X_1, \dots, X_d
- Label Y : Initial entropy $H(Y)$
- How much each feature X_i helps explain uncertainty in Y
 - Compute Information gain
$$IG(Y|X_i) = H(Y) - H(Y|X_i)$$
- Select feature of maximum IG and split
- Recurse

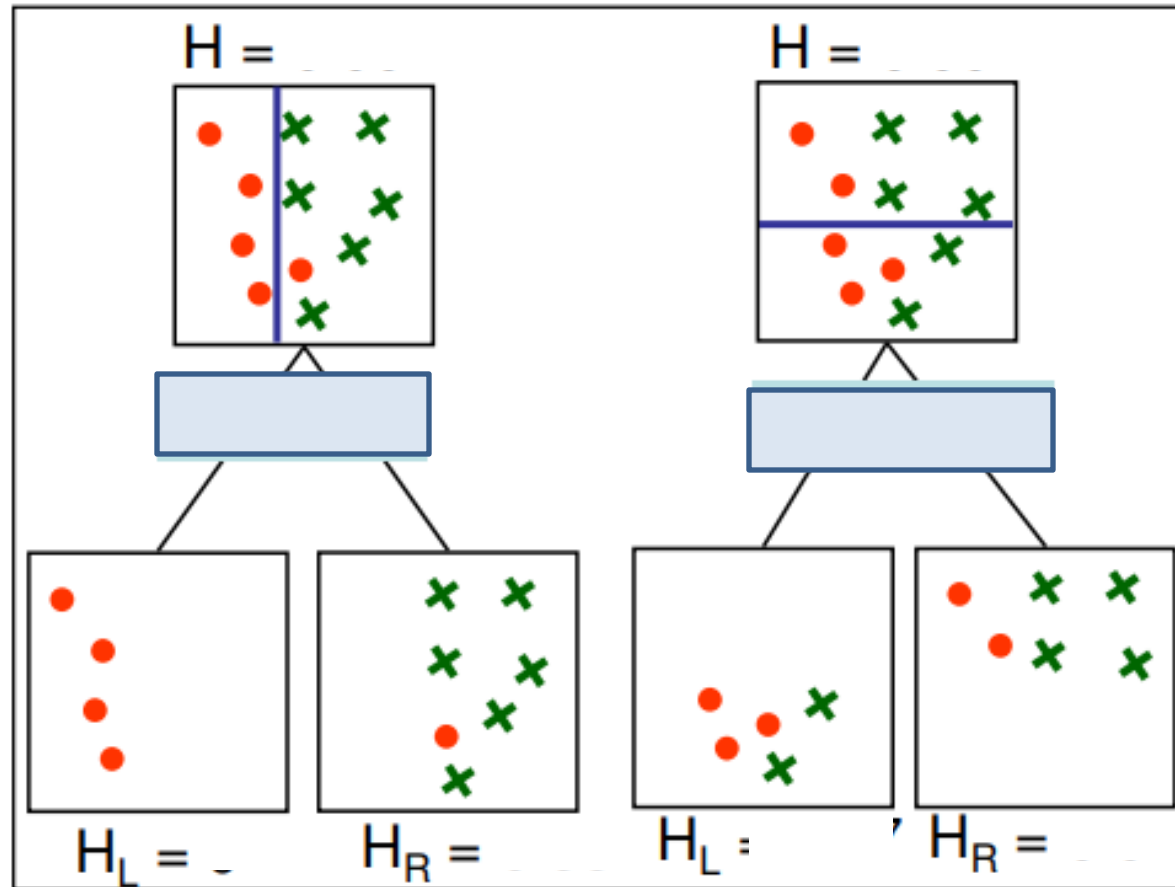
Example Information Gain



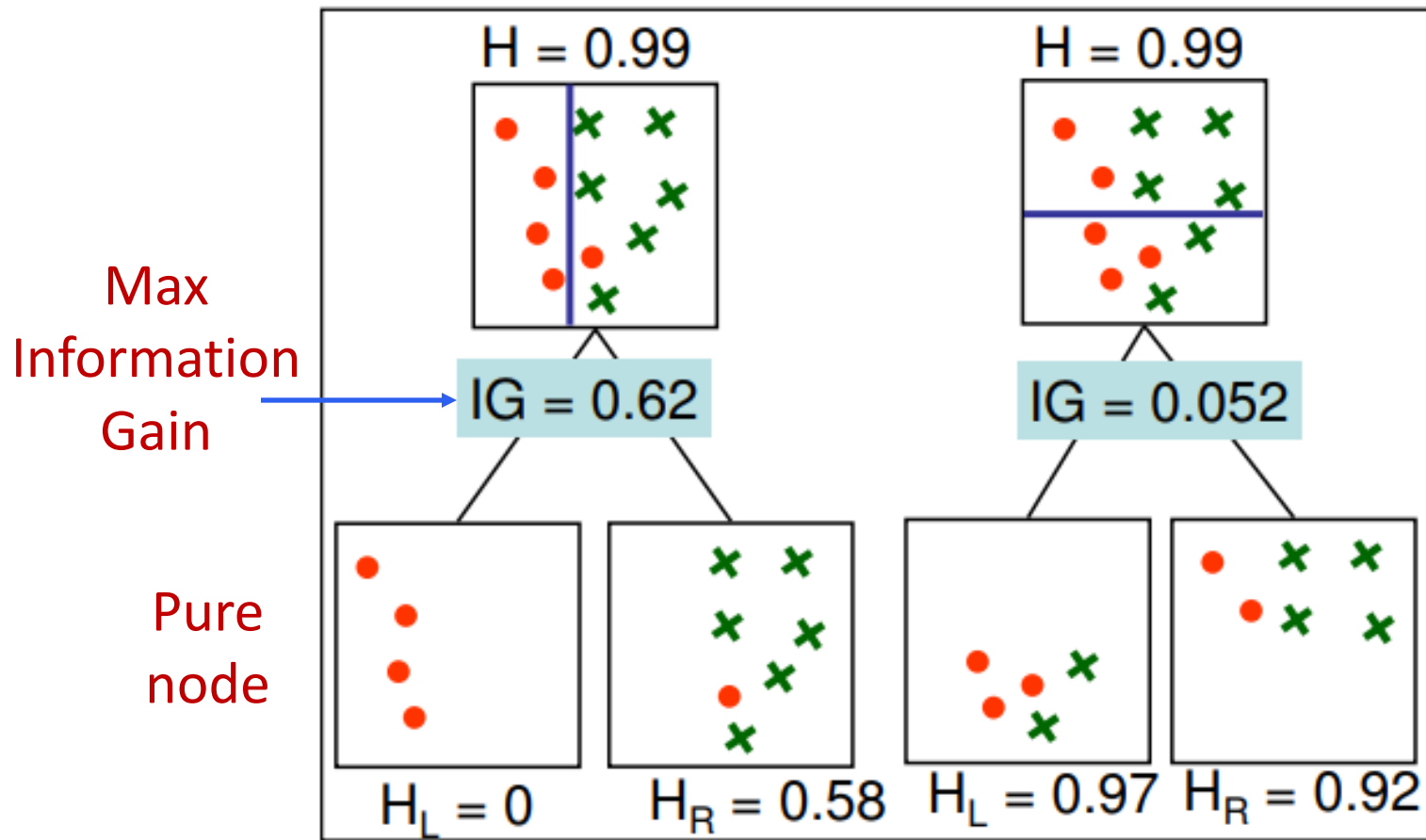
Example Information Gain



Example Information Gain



Example Information Gain



Learning Decision Trees

- Start from empty decision tree
- Split on **next best attribute (feature)**
 - Use, for example, information gain to select attribute:

$$\arg \max_i IG(X_i) = \arg \max_i H(Y) - H(Y | X_i)$$

- Recurse

ID3 algorithm uses Information Gain
Information Gain reduces uncertainty on Y

Impurity Metrics

Split a node according to max reduction of impurity

1. Entropy

2. Gini Index

- For binary case with prob p_0, p_1 :

$$I(p_0, p_1) = 2p_0p_1 = 2p_0(1 - p_0)$$

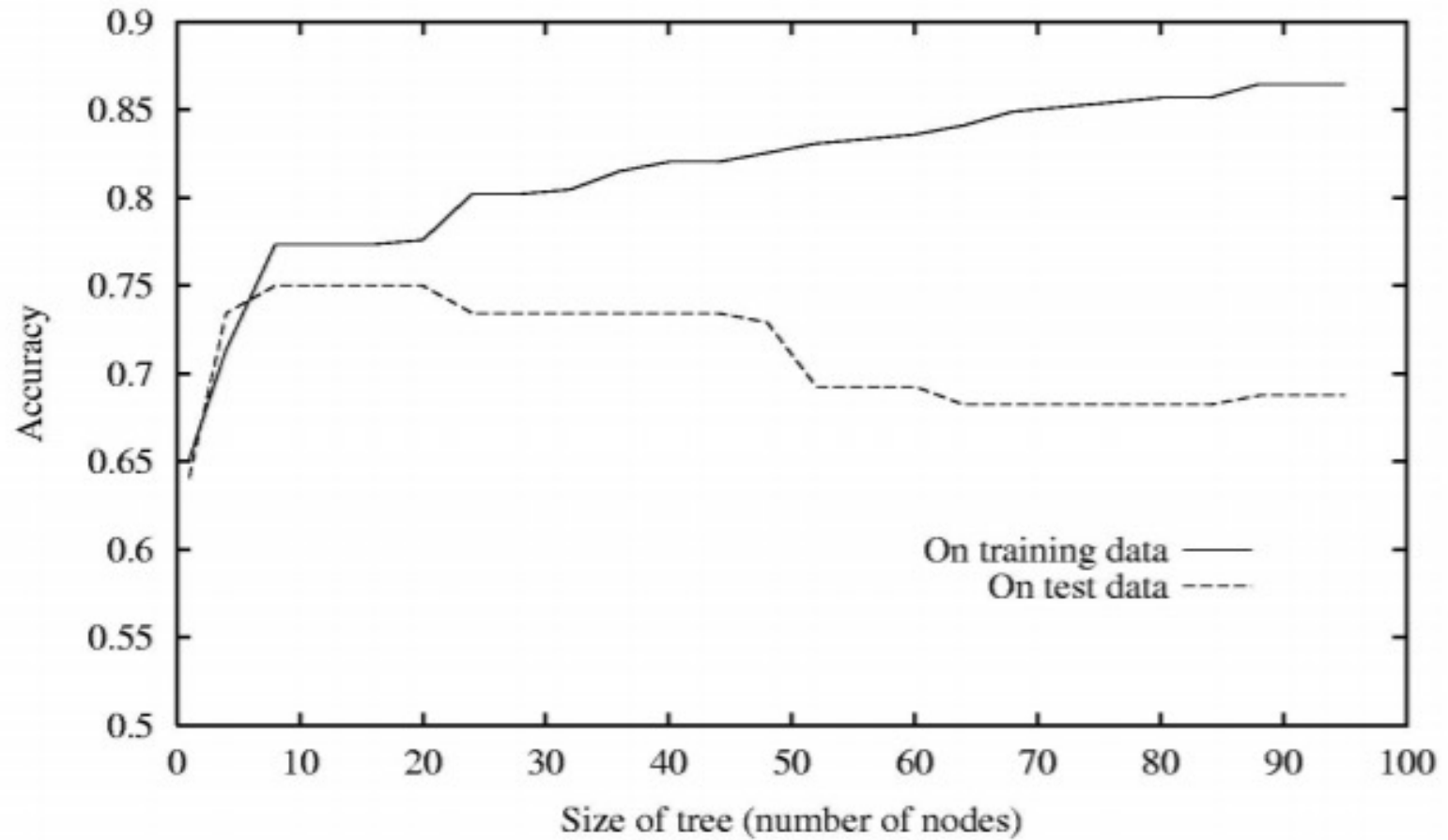
- For multi-class with prob p_1, \dots, p_K :

$$I(p_1, \dots, p_K) = 1 - \sum_{k=1}^K p_k^2$$

- Properties

- Impurity metrics have value 0 for pure nodes
- Impurity metrics are maximized for uniform distribution (nodes with most uncertainty)

Overfitting



Solutions against Overfitting

- Standard decision trees have no learning bias
 - Training set error is always zero!
 - (If there is no label noise)
 - Lots of variance
 - Must introduce some bias towards simpler trees
- Many strategies for picking simpler trees
 - Fixed depth
 - Minimum number of samples per leaf
- Pruning
 - Remove branches of the tree that increase error using cross-validation

Summary Decision Trees

- Greedy method for training
 - Not based on optimization or probabilities
- Uses impurity metric (e.g., information gain or Gini index) for splitting
- Advantages
 - Interpretability of decisions
- Limitations
 - Decision trees are prone to overfitting
 - Can be addressed by pruning or using ensembles of decision trees (next lecture)