# Decision Tree Learning

A method for approximating discrete-valued functions which is
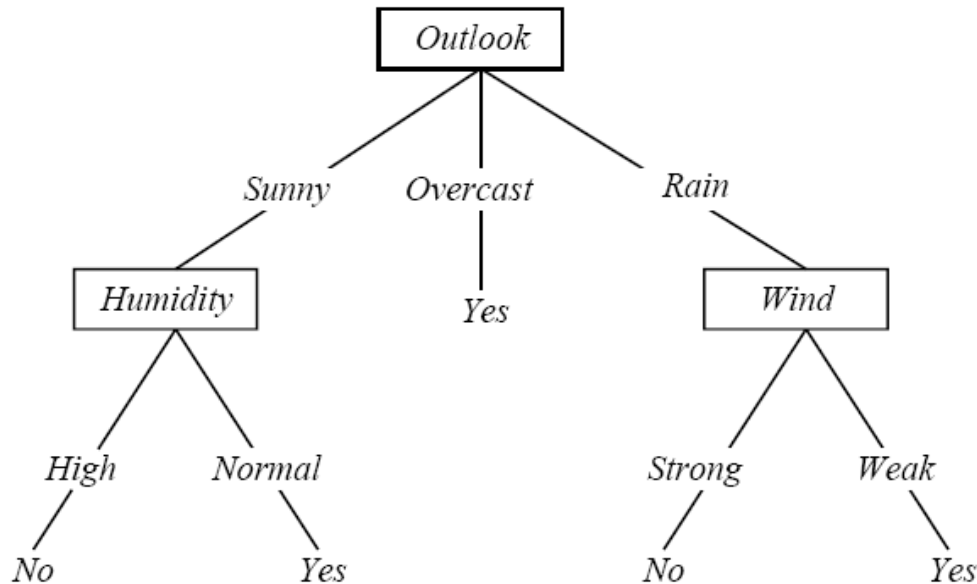
1. robust to noisy data

2. capable of learning disjunctive expressions

   ➔ Searches a completely expressive hypothesis space

3. easy to understand by humans (especially its learned results!)

Inductive Bias: "Prefer small trees over large trees"

- Decision tree representation

- ID3 learning algorithm

- Entropy, Information gain

- Overfitting

# Training Examples

| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| 1 | Sunny | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast | Hot | High | Weak | Yes |
| 4 | Rain | Mild | High | Weak | Yes |
| 5 | Rain | Cool | Normal | Weak | Yes |
| 6 | Rain | Cool | Normal | Strong | No |
| 7 | Overcast | Cool | Normal | Strong | Yes |
| 8 | Sunny | Mild | High | Weak | No |
| 9 | Sunny | Cool | Normal | Weak | Yes |
| 10 | Rain | Mild | Normal | Weak | Yes |
| 11 | Sunny | Mild | Normal | Strong | Yes |
| 12 | Overcast | Mild | High | Strong | Yes |
| 13 | Overcast | Hot | Normal | Weak | Yes |
| 14 | Rain | Mild | High | Strong | No |

lecture slides for textbook Machine Learning, T. Mitchell, McGraw Hill, 1997

How to convert tree into disjunctions of conjunctions of constraints on attribute values of instances?
➔ Convert into set of if-then rules

➔ Each path from the root node to a leaf corresponds to a conjunctions of attribute tests

➔ Tree is a disjunctions of conjunctions of constraints on attribute values

1. How Would you classify <outlook=Sunny, Temperature=Hot, Humidity=High, Wind=Strong>?

# Typical Datamining Task

Given:

- 9714 patient records each describing a pregnancy and birth
- Each patient record contains 215 features

Learn to predict:

- Classes of future patients at high risk for Emergency Cesarean Section

---

## Data:

| Patient103 time=1 | → | Patient103 time=2 | ... → | Patient103 time=n |

Age: 23
FirstPregnancy: no
Anemia: no
Diabetes: no
PreviousPrematureBirth: no
Ultrasound: ?

Elective C−Section: ?
Emergency C−Section: ?
...

Age: 23
FirstPregnancy: no
Anemia: no
Diabetes: YES
PreviousPrematureBirth: no
Ultrasound: abnormal

Elective C−Section: no
Emergency C−Section: ?
...

Age: 23
FirstPregnancy: no
Anemia: no
Diabetes: no
PreviousPrematureBirth: no
Ultrasound: ?

Elective C−Section: no
**Emergency C−Section: Yes**
...

lecture slides for textbook Machine Learning, T. Mitchell, McGraw Hill, 1997

Learned from medical records of 1000 women, <u>neg. examples are C-sections</u>:

[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-
        » Birth_Weight < 3349: [201+,10.6-] .95+ .05-
        » Birth_Weight >= 3349: [133+,36.4-] .78+ .22-
      - 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-

- Decision tree representation:
  - Each internal node tests an attribute
  - Each branch corresponds to an attribute value
  - Each leaf node assigns a classification

- How would we represent:
  - AND, OR, XOR
  - $(A \land B) \lor (C \land \neg D \land E)$
  - *M* of *N*

# A ∧ B

# A ∨ B
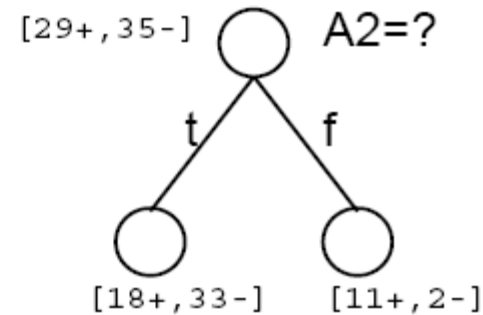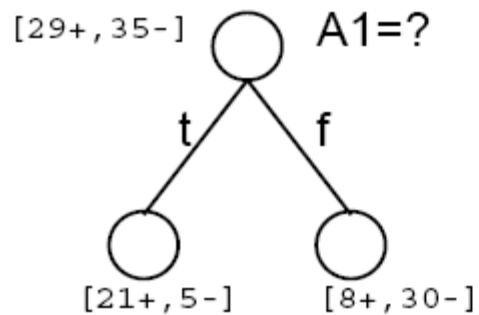
$$(A \land B) \lor (C \land \neg D \land E)$$

- Instances described by attribute values

- Target function is discrete valued

- Disjunctive hypotheses may be required

- Possibly noisy training data

- Training data may contain missing attribute values


- Examples:
  - Equipment or medical diagnosis
  - Credit risk analysis
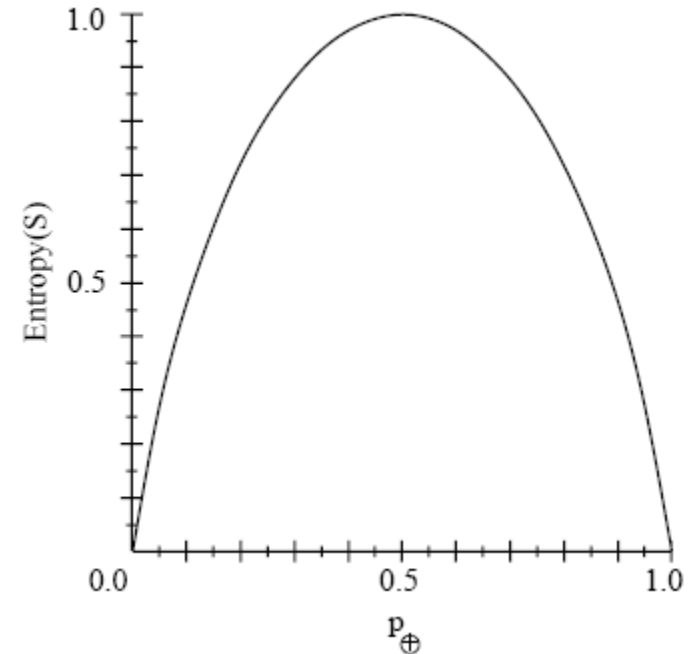  - Modeling calendar scheduling preferences

- Main loop:
  - *A* ← the "best" decision attribute for next *node*
  - Assign *A* as decision attribute for *node*
  - For each value of *A*, create new descendant of *node*
  - Sort training examples to leaf nodes
  - If training examples perfectly classified, then STOP, else iterate over new leaf nodes

- ## Which attribute is best?

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

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

*Entropy(S)* = expected number of bits needed to encode class + or - of randomly drawn member of *S* (under the optimal, shortest-length code)

Why?

*Information theory*: Optimal length code assigns $-\log_2 p$ bits to message having probability *p*.

So, expected number of bits to encode + or - of random member of *S*:

$$p_+(-\log_2 p_+) + p_-(-\log_2 p_-)$$

$$Entropy(S) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

lecture slides for textbook Machine Learning, T. Mitchell, McGraw Hill, 1997

$$Entropy(S) \equiv \sum_{i=1}^{c} - p_i \log_2 p_i$$

- $p_i$ = proportion of *S* belonging to class *i*
- Maximum value = $\log_2 c$

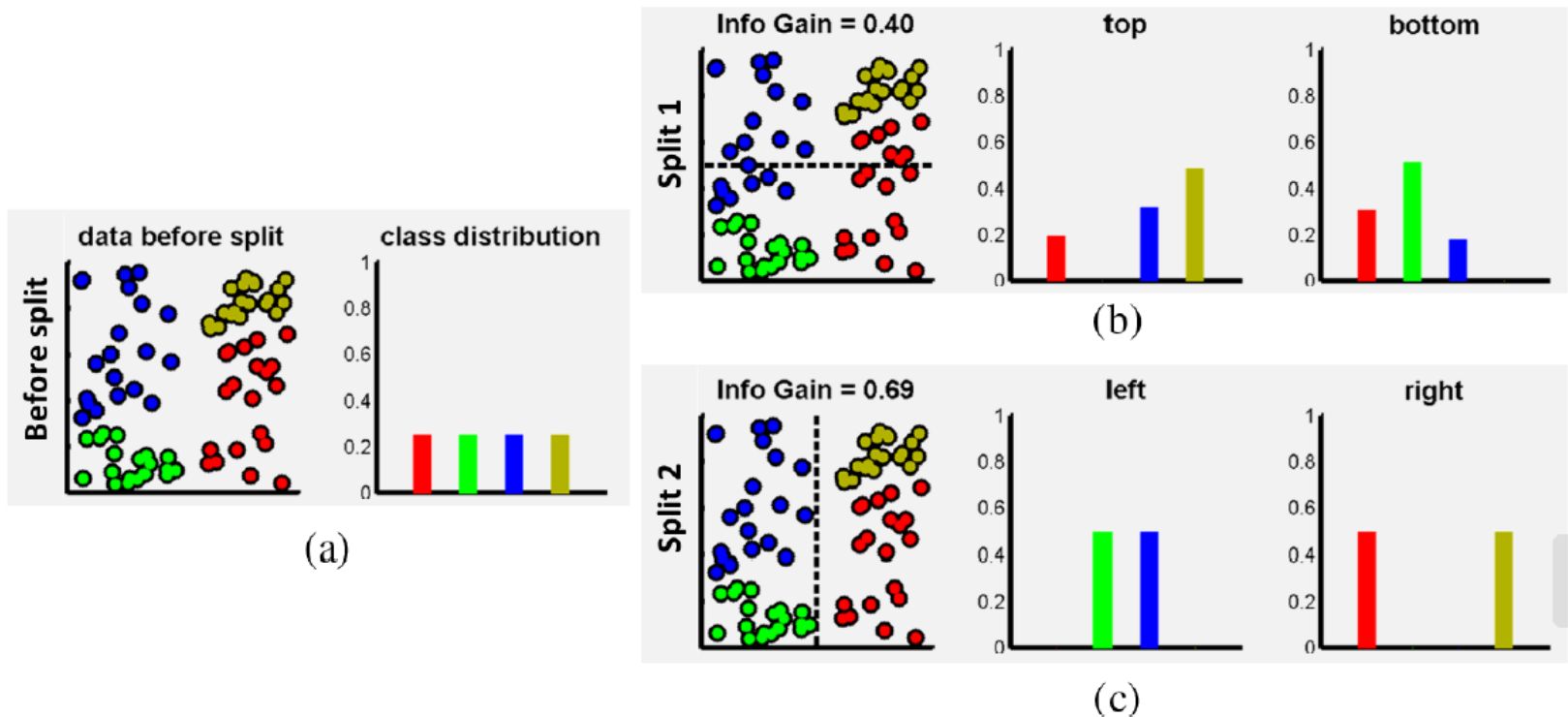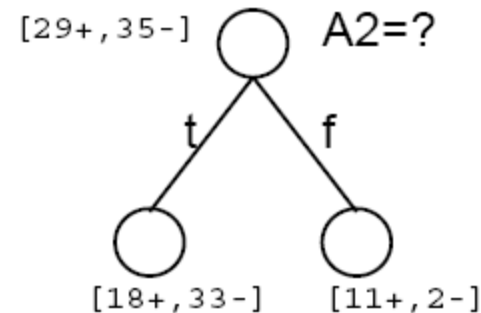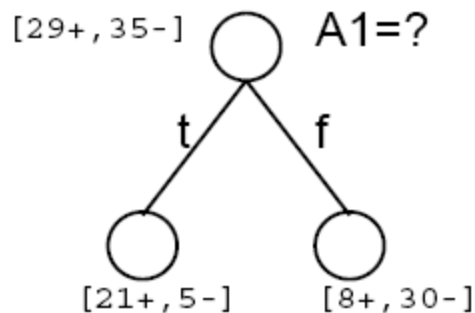- Measures the expected encoding length measured in *bits*

Fig. 2.5 Information gain for discrete, non-parametric distributions. (a) Dataset $\mathcal{S}$ before a split. (b) After a horizontal split. (c) After a vertical split. In this example the vertical split produces purer class distributions in the child nodes. Classes are colour coded.

- *Gain(S,A)* = expected reduction in entropy due to sorting on *A*

Original entropy

Entropy after partitioning

$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

[29+,35-]   A1=?

t   f

[21+,5-]   [8+,30-]

[29+,35-]   A2=?

t   f

[18+,33-]   [11+,2-]

➔ Let's compute all Entropy and Gain values

lecture slides for textbook Machine Learning, T. Mitchell, McGraw Hill, 1997
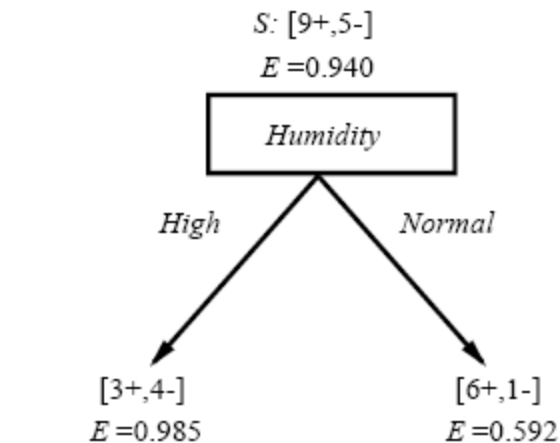
# Training Examples

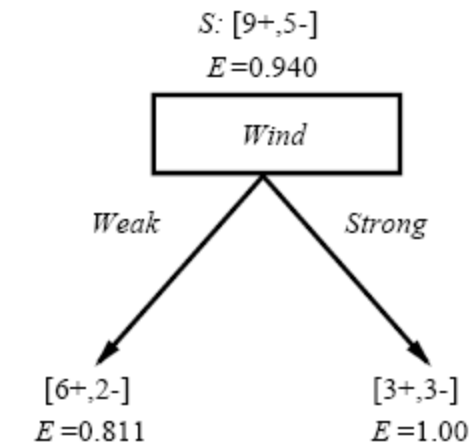| Day | Outlook | Temperature | Humidity | Wind | PlayTennis |
|-----|---------|-------------|----------|------|------------|
| 1 | Sunny | Hot | High | Weak | No |
| 2 | Sunny | Hot | High | Strong | No |
| 3 | Overcast | Hot | High | Weak | Yes |
| 4 | Rain | Mild | High | Weak | Yes |
| 5 | Rain | Cool | Normal | Weak | Yes |
| 6 | Rain | Cool | Normal | Strong | No |
| 7 | Overcast | Cool | Normal | Strong | Yes |
| 8 | Sunny | Mild | High | Weak | No |
| 9 | Sunny | Cool | Normal | Weak | Yes |
| 10 | Rain | Mild | Normal | Weak | Yes |
| 11 | Sunny | Mild | Normal | Strong | Yes |
| 12 | Overcast | Mild | High | Strong | Yes |
| 13 | Overcast | Hot | Normal | Weak | Yes |
| 14 | Rain | Mild | High | Strong | No |

lecture slides for textbook Machine Learning, T. Mitchell, McGraw Hill, 1997

- Which attribute is the best classifier?



S: [9+,5-]
E = 0.940

Humidity

High          Normal

[3+,4-]          [6+,1-]
E = 0.985          E = 0.592

Gain (S, Humidity)
= .940 - (7/14).985 - (7/14).592
= .151

S: [9+,5-]
E = 0.940

Wind

Weak          Strong

[6+,2-]          [3+,3-]
E = 0.811          E = 1.00

Gain (S, Wind)
= .940 - (8/14).811 - (6/14)1.0
= .048

• Which attribute should be tested here?

- $S_{sunny} = \{D1, D2, D8, D9, D11\}$

  $Gain\ (S_{sunny}\ ,\ Humidity)\ =\ .970\ -\ (3/5)\ 0.0\ -\ (2/5)\ 0.0\ =\ .970$

  $Gain\ (S_{sunny}\ ,\ Temperature)\ =\ .970\ -\ (2/5)\ 0.0\ -\ (2/5)\ 1.0\ -\ (1/5)\ 0.0\ =\ .570$

  $Gain\ (S_{sunny}\ ,\ Wind)\ =\ .970\ -\ (2/5)\ 1.0\ -\ (3/5)\ .918\ =\ .019$

- ## Humidity has highest information gain

- **Hypothesis space is complete!**
  - Target function surely is in there...

- **Outputs a single hypothesis (which one?)**
  - Can't play 20 questions...

- **No back tracking**
  - Local minima...

- **Statistically-based search choices**
  - Robust to noisy data...

- **Inductive bias:**
  - Approximately: „prefer shortest tree"

Note *H* is the power set of instances *X*

– Unbiased?

Not really...

- Preference for short trees, and for those with high information gain attributes near the root

- Bias is a *preference* for some hypotheses, rather than a *restriction* of hypothesis space *H*

- Occam's razor: prefer the shortest hypothesis that fits the data

- In favor:
  - There are fewer short hypotheses than long hypotheses
  - A short hypothesis that fits the data is unlikely to be coincidence
  - A long hypothesis that fits the data might be coincidence

- Opposed:
  - There are many ways to define small sets of hypotheses
    e.g., all trees with a prime number of nodes that use attributes beginning with "Z"
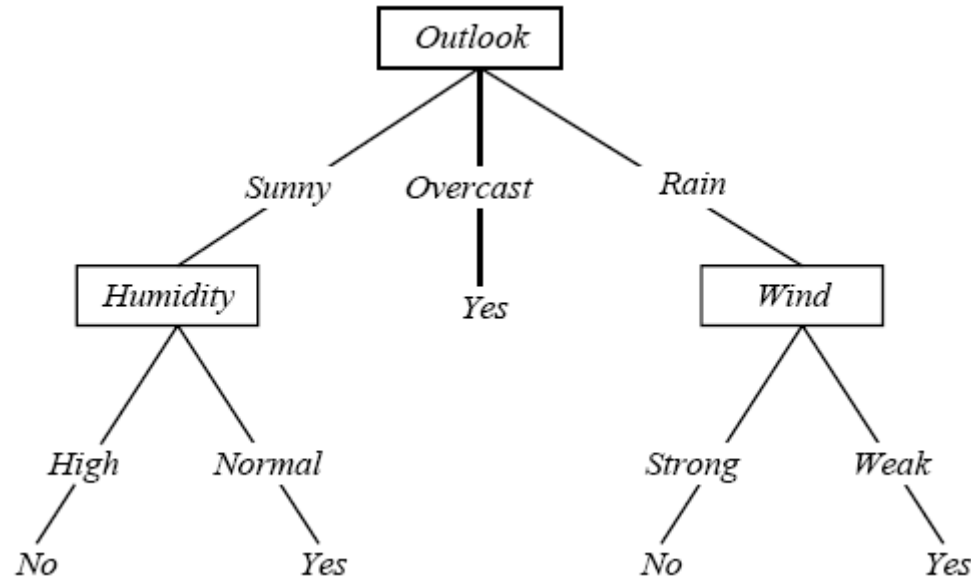  - What is short depends on language

# Consider adding noisy training example #15:

(Outlook, Temperature, Humidity, Wind)
*Sunny, Hot, Normal, Strong, PlayTennis=No*

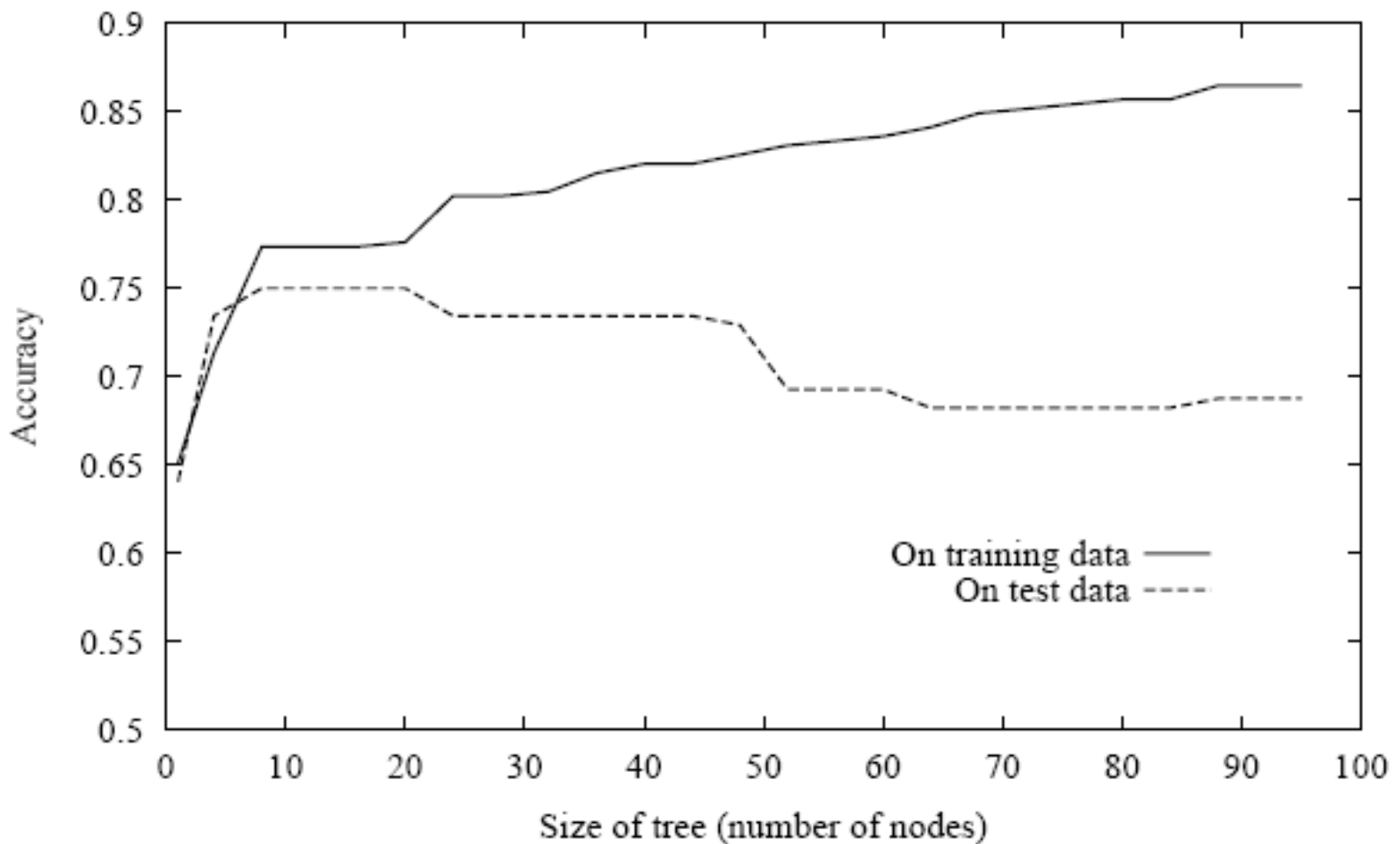# What is the effect on earlier tree?

- Consider error of hypothesis $h$ over
  - Training data: $error_{train}(h)$
  - Entire distribution of data: $error_D(h)$

- Hypothesis $h \in H$ overfits training data if there is an alternative hypothesis $h' \in H$ such that

$$error_{train}(h) < error_{train}(h')$$
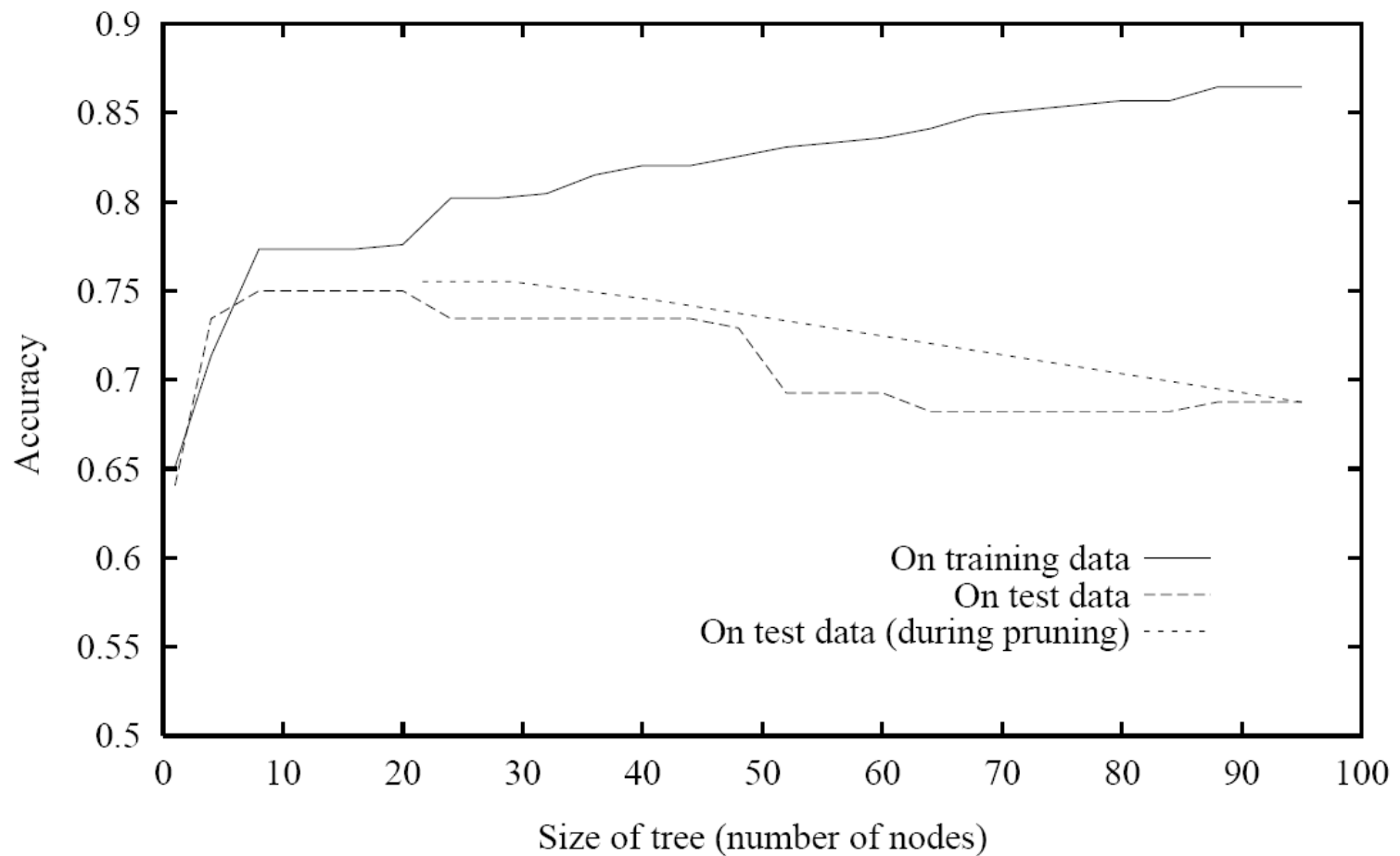
and

$$error_D(h) > error_D(h').$$

- How can we avoid overfitting?

  - stop growing when data split is not statistically significant

  - grow full tree, then post-prune

- How to select „best" tree:

  - Measure performance over training data

  - Measure performance over separate validation data set

  - Minimum Description Length: minimize *size(tree) + size(misclassifications(tree))*
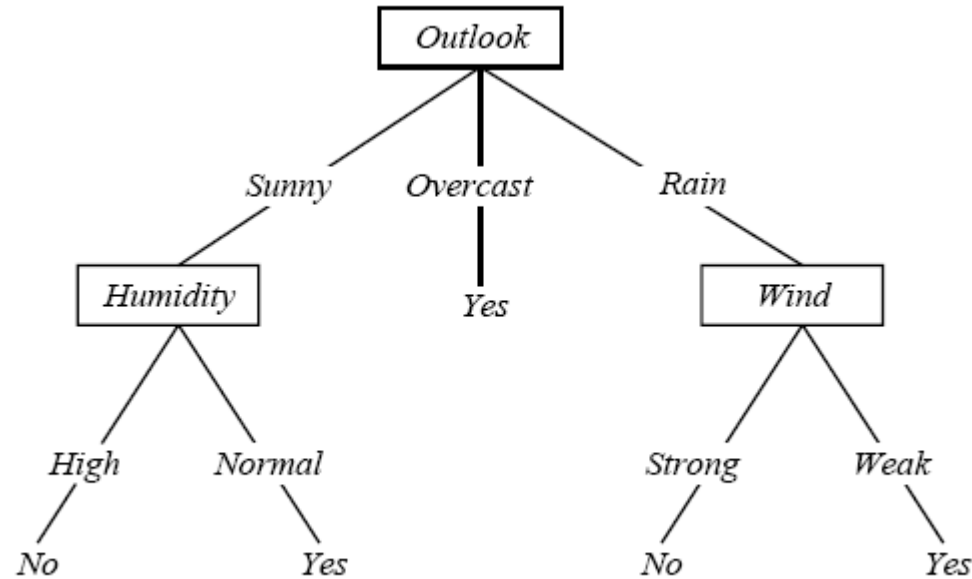
- Split data into *training* and *validation* set

- Do until further pruning is harmful:

  1. Evaluate impact on *validation* set of pruning each possible node (plus those below it)

  2. Greedily remove the one that most improves *validation* set accuracy

- IF (*Outlook=Sunny*) and (*Humidity=High*)
  THEN *PlayTennis=No*

- IF (*Outlook=Sunny*) and (*Humidity=Normal*)
  THEN *PlayTennis=Yes*

- …

- Create a discrete attribute to test continuous
  - Temperature = 82.5
  - (Temperature > 54) = true ELSE false

$$
\begin{array}{llcccccc}
Temperature: & 40 & 48 & | & 60 & 72 & 80 & | & 90 \\
PlayTennis: & No & No & | & Yes & Yes & Yes & | & No
\end{array}
$$

- **Problem:** Gain measure has a natural bias towards attributes with many values
  - if attribute has many values, *Gain* will select it
  - Imagine using Date = Jun 3 1996 as attribute
- Use alternative method for selecting attributes
- One approach: use GainRatio instead:

$$GainRatio(S, A) \equiv \frac{Gain(S, A)}{SplitInformation(S, A)}$$

$$SplitInformation(S, A) \equiv -\sum_{i=1}^{c} \frac{|S_i|}{|S|} \log_2 \frac{|S_i|}{|S|}$$

$\log_2 c$

where $S_i$ is subset of $S$ for which $A$ has value $v_i$

Discourages the selection of attributes with many uniformly distributed values

E.g. two-class problem in 2D ($x \in \mathbb{R}^2$):
$$f(x) = I_{x_2 > x_1}$$

This linear, but not axis aligned decision function needs to be approximated by **many** nodes. It cannot be efficiently learned. It also needs many training examples in order to be approximately learned.

What if some examples have missing values of *A*?

Use training example anyway, sort through tree

1. If node *n* tests *A*, assign most common value of *A* among other examples sorted to node *n*

2. Assign most common value of *A* among other examples with same target value

3. Assign probability $p_i$ to each possible value $v_i$ of *A*
   – assign fraction $p_i$ of example to each descendant in tree

Classify new examples in same fashion.

Model $n$ -dimensional data $\{(x, y)\}, x \in \mathbb{R}^n$ by $n$-variate Gaussian:

$$G = ..$$

In the continuous case we define the differential entropy

$$H(S) = -\int_y p(y)\, log\big(p(y)\big)\, dy$$

which becomes for the multi-variate Gaussian

$$H(S) = \frac{1}{2}\log\big((2\pi e)^n\big) \cdot \det(\Sigma)$$
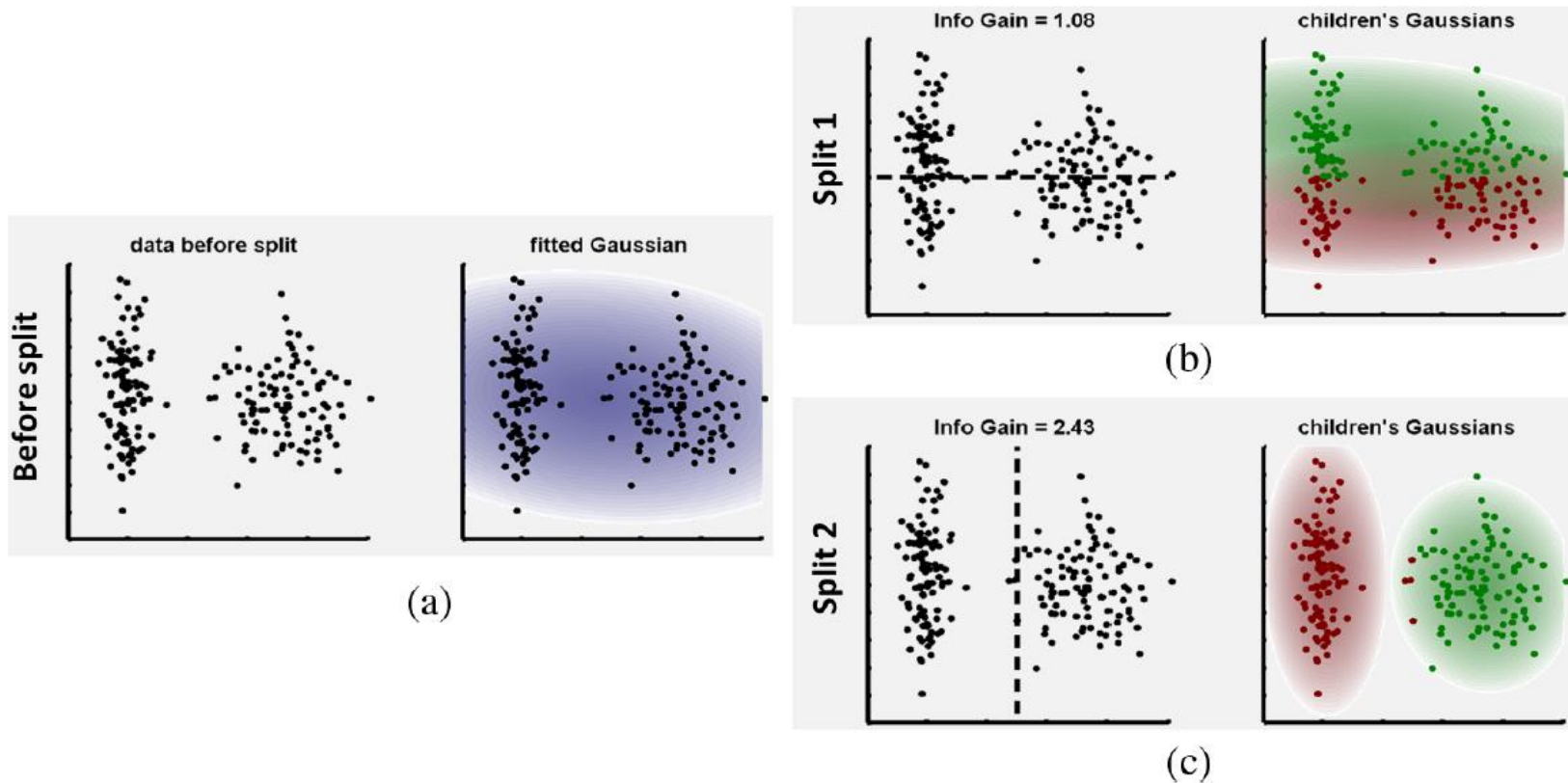
Fig. 2.6 Information gain for continuous, parametric densities. (a) Dataset $\mathcal{S}$ before a split. (b) After a horizontal split. (c) After a vertical split. A vertical split produces better separation and a correspondingly higher information gain.