**Outline**
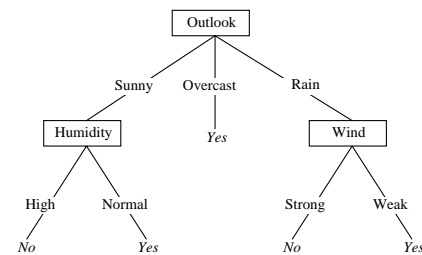
- Decision tree representation
- ID3 learning algorithm
- Entropy, Information gain
- Overfitting

---

**Decision Tree for** *PlayTennis*

---

**Decision Trees**

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

How would we represent:
- $\wedge, \vee$, XOR
- $(A \wedge B) \vee (C \wedge \neg D \wedge E)$
- $M$ of $N$

---

**When to Consider Decision Trees**

- Instances describable by attribute–value pairs
- Target function is discrete valued
- Disjunctive hypothesis may be required
- Possibly noisy training data

Examples:
- Medical diagnosis
- Credit risk analysis
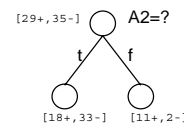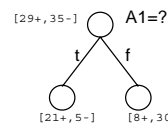- Modeling calendar scheduling preferences

---

**Top-Down Induction of Decision Trees**

Main loop:
1. $A \leftarrow$ the "best" decision attribute for next *node*
2. Assign $A$ as decision attribute for *node*
3. For each value of $A$, create new descendant of *node*
4. Sort training examples to leaf nodes
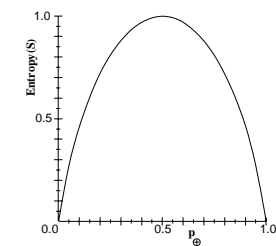5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

Which attribute is best?

---

**Entropy**



- $S$ is a sample of training examples
- $p_{\oplus}$ is the proportion of positive examples in $S$
- $p_{\ominus}$ is the proportion of negative examples in $S$
- Entropy measures the impurity of $S$

$$Entropy(S) \equiv -p_{\oplus} \log_2 p_{\oplus} - p_{\ominus} \log_2 p_{\ominus}$$

## Entropy

$Entropy(S)$ = expected number of bits needed to encode class ($\oplus$ or $\ominus$) 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 $\oplus$ or $\ominus$ of random member of $S$:

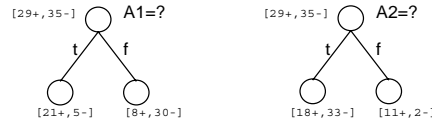$$p_{\oplus}(-\log_2 p_{\oplus}) + p_{\ominus}(-\log_2 p_{\ominus})$$

$$Entropy(S) \equiv -p_{\oplus}\log_2 p_{\oplus} - p_{\ominus}\log_2 p_{\ominus}$$

---

## Information Gain

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

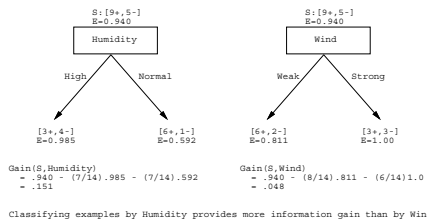$$Gain(S, A) \equiv Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

[29+,35-] ◯ A1=?
   t   f

◯ [21+,5-]   ◯ [8+,30-]

[29+,35-] ◯ A2=?
   t   f

◯ [18+,33-]   ◯ [11+,2-]

---

## Training Examples

| 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 |

---

## Selecting the Next Attribute

**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

Classifying examples by Humidity provides more information gain than by Wind

---

**Partially Learned Decision Tree**

{D1, D2, ..., D14}
[9+,5-]

Outlook

Sunny    Overcast    Rain

{D1,D2,D8,D9,D11}    {D3,D7,D12,D13}    {D4,D5,D6,D10,D14}
[2+,3-]    [4+,0-]    [3+,2-]

?    Yes    ?

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

---

## Hypothesis Space Search by ID3

A1     A2     ...
A2     A2
A3     A4
...

## Hypothesis Space Search by ID3

- Hypothesis space is complete!
  - Target function surely in there...
- Outputs a single hypothesis (which one?)
  - Can't play 20 questions...
- No back tracking
  - Local minima...
- Statisically-based search choices
  - Robust to noisy data...
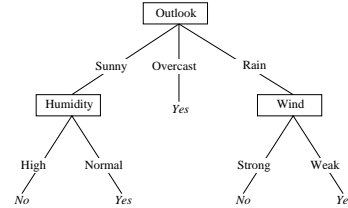- Inductive bias: approx "prefer shortest tree"

13

## Overfitting in Decision Trees

Consider adding noisy training example #15:

$Sunny,\ Hot,\ Normal,\ Strong,\ PlayTennis = No$

What effect on earlier tree?

```
                    Outlook
           Sunny   Overcast    Rain
        Humidity      Yes       Wind
                      Yes
     High    Normal         Strong   Weak
      No       Yes            No       Yes
```

14

## Overfitting

Consider error of hypothesis $h$ over

- training data: $error_{train}(h)$
- entire distribution $\mathcal{D}$ of data: $error_{\mathcal{D}}(h)$

Hypothesis $h \in H$ **overfits** training data if there is an alternative hypothesis $h' \in H$ such that
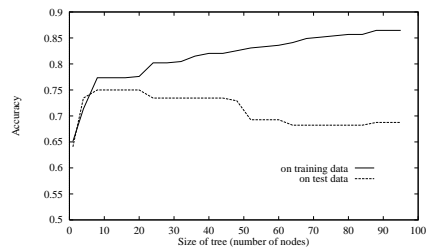
$$error_{train}(h) < error_{train}(h')$$

but

$$error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h')$$

15

## Overfitting in Decision Tree Learning



16