# Administrative Details

- HW 3 is posted

  - Due on March 4

  - Available until March 6

- MidTerm on March 6

  - Syllabus:

    - Everything covered in class up to Feb 28

  - Cheat Sheet:

    - One page written on both sides

# Decision Tree

# Supervised Learning: find $f$

- Given: Training set $\{(x_i, y_i) \mid i = 1 \ldots n\}$
- Find: A good approximation to $f : X \rightarrow Y$

# Supervised Learning: find $f$

- Given: Training set $\{(x_i, y_i) \mid i = 1 \ldots n\}$

- Find: A good approximation to $f: X \rightarrow Y$

  Examples: what are $X$ and $Y$ ?

  - Spam Detection
    - Map email to {Spam,Ham}
  - Digit recognition
    - Map pixels to {0,1,2,3,4,5,6,7,8,9}
  - Stock Prediction
    - Map new, historic prices, etc. to $\mathbb{R}$ (the real numbers)

# A Supervised Learning Problem

- Consider a simple, Boolean dataset:
    - $f : X \rightarrow Y$
    - $X = \{0,1\}^4$
    - $Y = \{0,1\}$

- Question 1: How should we pick the *hypothesis space*, the set of possible functions $f$?

- Question 2: How do we find the best $f$ in the hypothesis space?

Dataset:

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---------|-------|-------|-------|-------|-----|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

# Most General Hypothesis Space

Consider all possible boolean functions over four input features!

Dataset:

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|-----|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

# Most General Hypothesis Space

Consider all possible boolean functions over four input features!

- $2^{16}$ possible hypotheses

- $2^9$ are consistent with our dataset

- How do we choose the best one?

Dataset:

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|-------|-------|-------|-------|-----|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

# Most General Hypothesis Space

Consider all possible boolean functions over four input features!

- $2^{16}$ possible hypotheses

- $2^{9}$ are consistent with our dataset

- How do we choose the best one?

**Dataset:**

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | ? |
| 0 | 0 | 0 | 1 | ? |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | ? |
| 1 | 0 | 0 | 0 | ? |
| 1 | 0 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | ? |
| 1 | 0 | 1 | 1 | ? |
| 1 | 1 | 0 | 0 | 0 |
| 1 | 1 | 0 | 1 | ? |
| 1 | 1 | 1 | 0 | ? |
| 1 | 1 | 1 | 1 | ? |

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
|---|---|---|---|---|---|
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

# A Restricted Hypothesis Space

Consider all conjunctive boolean functions.

- 16 possible hypotheses

- None are consistent with our dataset

- How do we choose the best one?

Dataset:

| Rule | Counterexample |
| --- | --- |
| $\Rightarrow y$ | 1 |
| $x_1 \Rightarrow y$ | 3 |
| $x_2 \Rightarrow y$ | 2 |
| $x_3 \Rightarrow y$ | 1 |
| $x_4 \Rightarrow y$ | 7 |
| $x_1 \wedge x_2 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_3 \wedge x_4 \Rightarrow y$ | 4 |
| $x_1 \wedge x_2 \wedge x_3 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |
| $x_1 \wedge x_2 \wedge x_3 \wedge x_4 \Rightarrow y$ | 3 |

| Example | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $y$ |
| --- | --- | --- | --- | --- | --- |
| 1 | 0 | 0 | 1 | 0 | 0 |
| 2 | 0 | 1 | 0 | 0 | 0 |
| 3 | 0 | 0 | 1 | 1 | 1 |
| 4 | 1 | 0 | 0 | 1 | 1 |
| 5 | 0 | 1 | 1 | 0 | 0 |
| 6 | 1 | 1 | 0 | 0 | 0 |
| 7 | 0 | 1 | 0 | 1 | 0 |

# Simple Training Data 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 for
f: <Outlook, Temperature, Humidity, Wind> → PlayTennis?



Each internal node: test one discrete-valued attribute $X_i$

Each branch from a node: selects one value for $X_i$

Each leaf node: predict Y (or P(Y|X ∈ leaf))

# A Decision tree for

f: <Outlook, Temperature, Humidity, Wind> → PlayTennis?

# Decision Tree Learning

**Problem Setting**:

- Set of possible instances $X$
  - each instance $x$ in $X$ is a feature vector
  - e.g., *<Humidity=low, Wind=weak, Outlook=rain, Temp=hot>*
- Unknown target function $f : X \rightarrow Y$
  - Y=1 if we play tennis on this day, else 0
- Set of function hypotheses $H=\{ h \mid h : X \rightarrow Y \}$
  - each hypothesis $h$ is a decision tree
  - trees sorts $x$ to leaf, which assigns $y$

# Decision Tree Learning

**Problem Setting**:

- Set of possible instances $X$

  – each instance $x$ in $X$ is a feature vector

  $x = < x_1, x_2 \ldots x_n>$

- 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 $\{<x^{(i)}, y^{(i)}>\}$ of unknown target function $f$

**Output**:

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

# Choosing the Best Attribute—An Example

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

| 4 | 4 |
|---|---|

x1

| 1 | 3 |
|---|---|

| 3 | 1 |
|---|---|

J=2

| 4 | 4 |
|---|---|

x2

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

| 4 | 4 |
|---|---|

x3

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

# Choosing the Best Attribute—An Example

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

| 4 | 4 |
|---|---|

x1

| 1 | 3 |   | 3 | 1 |
|---|---|---|---|---|

J=2

| 4 | 4 |
|---|---|

x2

| 2 | 2 |   | 2 | 2 |
|---|---|---|---|---|

J=4

| 4 | 4 |
|---|---|

x3

| 2 | 2 |   | 2 | 2 |
|---|---|---|---|---|

J=4

# Choosing the Best Attribute—An Example

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

| 4 | 4 |
|---|---|

x1

| 1 | 3 | | 3 | 1 |
|---|---|---|---|---|

J=2

| 4 | 4 |
|---|---|

x2

| 2 | 2 | | 2 | 2 |
|---|---|---|---|---|

J=4

| 4 | 4 |
|---|---|

x3

| 2 | 2 | | 2 | 2 |
|---|---|---|---|---|

J=4

# Choosing the Best Attribute—An Example

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

| 4 | 4 |
|---|---|

x1

| 1 | 3 |
|---|---|

| 3 | 1 |
|---|---|

J=2

| 4 | 4 |
|---|---|

x2

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

| 4 | 4 |
|---|---|

x3

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

# Choosing the Best Attribute—An Example

| $x_1$ | $x_2$ | $x_3$ | $y$ |
|-------|-------|-------|-----|
| 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 |
| 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |

| 4 | 4 |
|---|---|

x1

| 1 | 3 |
|---|---|

| 3 | 1 |
|---|---|

J=2

| 4 | 4 |
|---|---|

x2

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

| 4 | 4 |
|---|---|

x3

| 2 | 2 |
|---|---|

| 2 | 2 |
|---|---|

J=4

# Choosing the Best Attribute (3)

Unfortunately, this measure does not always work well, because it does not detect cases where we are making "progress" toward a good tree.

# A Better Heuristic From Information Theory

Let $V$ be a random variable with the following probability distribution:

| $P(V = 0)$ | $P(V = 1)$ |
|:---:|:---:|
| 0.2 | 0.8 |

The *surprise*, $S(V = v)$ of each value of $V$ is defined to be

$$S(V = v) = -\lg P(V = v).$$

An event with probability 1 gives us zero surprise.

An event with probability 0 gives us infinite surprise!

# A Better Heuristic From Information Theory

Let $V$ be a random variable with the following probability distribution:

| $P(V = 0)$ | $P(V = 1)$ |
|:---:|:---:|
| 0.2 | 0.8 |

The *surprise*, $S(V = v)$ of each value of $V$ is defined to be

$$S(V = v) = -\lg P(V = v).$$

An event with probability 1 gives us zero surprise.

An event with probability 0 gives us infinite surprise!

It turns out that the surprise is equal to the number of bits of information that need to be transmitted to a recipient who knows the probabilities of the results.

# A Better Heuristic From Information Theory

Let $V$ be a random variable with the following probability distribution:

| $P(V = 0)$ | $P(V = 1)$ |
|:----------:|:----------:|
| 0.2 | 0.8 |

The *surprise*, $S(V = v)$ of each value of $V$ is defined to be

$$S(V = v) = -\lg P(V = v).$$

An event with probability 1 gives us zero surprise.

An event with probability 0 gives us infinite surprise!

It turns out that the surprise is equal to the number of bits of information that need to be transmitted to a recipient who knows the probabilities of the results.
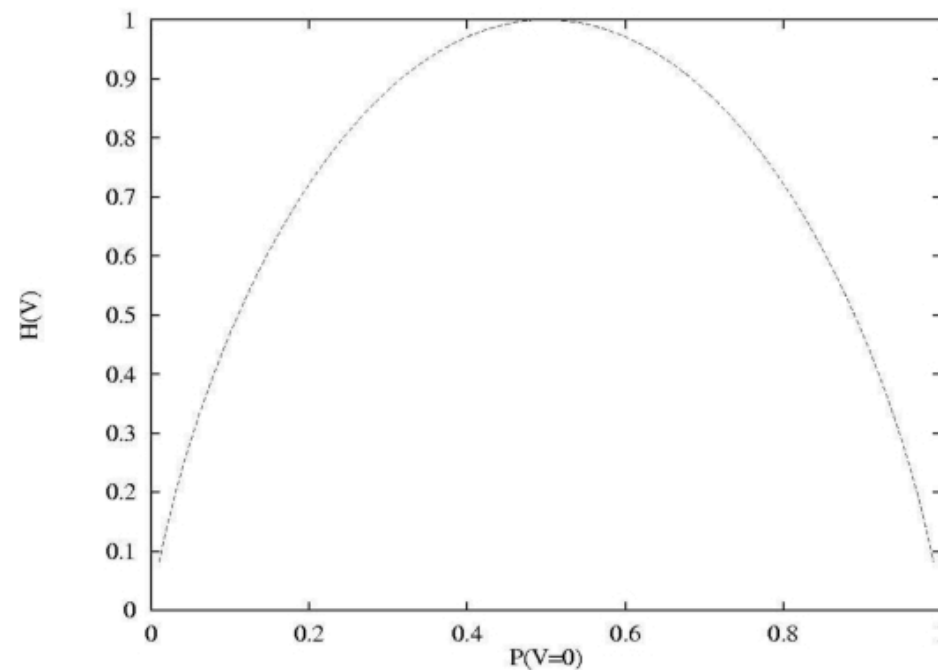
This is also called the *description length* of $V = v$.

Fractional bits only make sense if they are part of a longer message (e.g., describe a whole sequence of coin tosses).

# Entropy

The *entropy* of $V$, denoted $H(V)$ is defined as follows:

$$H(V) = \sum_{v=0}^{1} -P(H = v) \lg P(H = v).$$

This is the average surprise of describing the result of one "trial" of $V$ (one coin toss).



Entropy can be viewed as a measure of uncertainty.