

Concept Learning

Decision Trees

- 1 Concepts and Hypotheses
 - Definitions
 - Example
 - Hypotheses
- 2 Decision Trees
 - Using Trees
 - Learning
- 3 Unpredictability
 - Entropy
 - Entropy for datasets
 - Information Gain
- 4 Improvements

- 1 Concepts and Hypotheses
 - Definitions
 - Example
 - Hypotheses
- 2 Decision Trees
 - Using Trees
 - Learning
- 3 Unpredictability
 - Entropy
 - Entropy for datasets
 - Information Gain
- 4 Improvements

Concept Learning

Concept Learning

Learning of a **boolean function** from examples

Categories

- “Nice weather”
- “Dog”
- “Motor vehicle”
- “Criminal offence”

Subsets of a superset X

Terminology

c The concept to learn

$$c(x) \rightarrow \text{True/False}, \quad x \in X$$

h Hypothesis, Result of the learning ("guessed c ")

$$h(x) \rightarrow \text{True/False}, \quad x \in X$$

H Hypotheses space, All conceivable hypotheses
(before data arrives)

$$h \in H$$

D Set of available training data

$$D \subseteq X$$

Example of a *concept*

"Nice Weather"

Let each "weather instance" x_i be composed of four **attributes**:

$$x_1 = \langle \text{Sunny, Warm, Windy, Dry} \rangle$$

$$x_2 = \langle \text{Cloudy, Warm, Calm, Dry} \rangle$$

$$x_3 = \dots$$

Generally: $Sky \times Temperature \times Wind \times Humidity$

Terminology

Two kinds of training examples

Positive example:

$$x : c(x) = \text{True}, \quad x \in D$$

Negative example:

$$x : c(x) = \text{False}, \quad x \in D$$

Assume that the attributes can only take on certain discrete values:

$$\text{Sky} \in \{ \text{Sunny, Cloudy, Rainy} \}$$

$$\text{Temp} \in \{ \text{Warm, Cold} \}$$

$$\text{Wind} \in \{ \text{Windy, Calm} \}$$

$$\text{Humid} \in \{ \text{Humid, Dry} \}$$

Number of possible weathers: $|X| = 3 \cdot 2 \cdot 2 \cdot 2 = 24$

Typical training samples

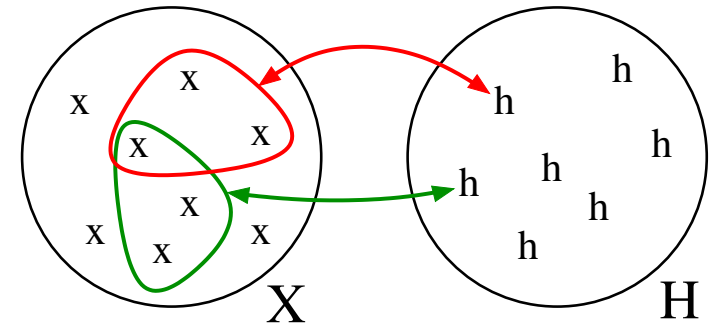
$x_1 = \langle \text{Sunny, Warm, Windy, Dry} \rangle \rightarrow \text{Nice}$
 $x_2 = \langle \text{Sunny, Warm, Windy, Humid} \rangle \rightarrow \text{Nice}$
 $x_3 = \langle \text{Rainy, Cold, Windy, Humid} \rangle \rightarrow \text{Bad}$
 $x_4 = \langle \text{Sunny, Warm, Calm, Humid} \rangle \rightarrow \text{Nice}$

How many hypotheses can we choose from?
How many subsets does X have?

$$|H| = 2^{|X|}$$

$$|H| = 2^{24} = 16777216$$

What does the hypotheses space H look like?

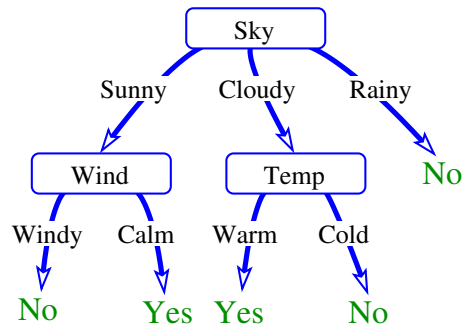


Each hypothesis h corresponds to one **subset** of X

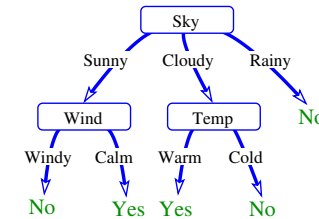
- 1 Concepts and Hypotheses
 - Definitions
 - Example
 - Hypotheses
- 2 Decision Trees
 - Using Trees
 - Learning
- 3 Unpredictability
 - Entropy
 - Entropy for datasets
 - Information Gain
- 4 Improvements

Decision Tree

- Test the attributes sequentially
- Choose attributes to test depending on earlier attribute values



The results (classifications) are coded by the *leaves*



What does the tree encode?

$$(\text{Sunny} \wedge \text{Calm}) \vee (\text{Cloudy} \wedge \text{Warm})$$

Works as a *disjunction of conjunctions*

Normal Form for boolean functions

Arbitrary boolean functions can be represented!

How can a decision tree be constructed automatically?

- 1 Choose an attribute to test
- 2 Branches with a unique class become leaves
- 3 Other branches are extended recursively

Remaining question: how do we choose attributes?

Greedy approach:

Choose the attribute which *tells us most* about the answer

1 Concepts and Hypotheses

- Definitions
- Example
- Hypotheses

2 Decision Trees

- Using Trees
- Learning

3 Unpredictability

- Entropy
- Entropy for datasets
- Information Gain

4 Improvements

Entropy

Entropy — measure of **unpredictability**

$$\text{Entropy} = \sum_i -p_i \log_2 p_i$$

p_i probability for outcome i

Entropy

Example: tossing a coin

$$p_{\text{head}} = 0.5; \quad p_{\text{tail}} = 0.5$$

$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= -0.5 \log_2 0.5 + -0.5 \log_2 0.5 = -0.5 \underbrace{\log_2 0.5}_{-1} + -0.5 \underbrace{\log_2 0.5}_{-1} = \\ &= 1 \end{aligned}$$

The result of a coin-toss has **1 bit** of information

Entropy

Example: rolling a dice

$$p_1 = \frac{1}{6}; \quad p_2 = \frac{1}{6}; \dots \quad p_6 = \frac{1}{6}$$

$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= 6 \times -\frac{1}{6} \log_2 \frac{1}{6} = \\ &= -\log_2 \frac{1}{6} = \log_2 6 \approx 2.58 \end{aligned}$$

The result of a dice-roll has **2.58 bit** of information

Entropy

Example: rolling a **fake dice**

$$p_1 = 0.1; \dots \quad p_5 = 0.1; \quad p_6 = 0.5$$

$$\begin{aligned} \text{Entropy} &= \sum_i -p_i \log_2 p_i = \\ &= -5 \cdot 0.1 \log_2 0.1 - 0.5 \log_2 0.5 = \\ &\approx 2.16 \end{aligned}$$

A real dice is **more unpredictable** (2.58 bit) than a fake (2.16 bit)

Entropy

Unpredictability of a **dataset**

- 100 examples, 42 positive

$$-\frac{58}{100} \log_2 \frac{58}{100} - \frac{42}{100} \log_2 \frac{42}{100} = 0.981$$

- 100 examples, 3 positive

$$-\frac{97}{100} \log_2 \frac{97}{100} - \frac{3}{100} \log_2 \frac{3}{100} = 0.194$$

Back to the decision trees

Smart idea:

Ask about the attribute which maximizes the expected reduction of the entropy.

Information Gain

Assume that we ask about attribute A for a dataset S

$$\text{Gain} = \underbrace{\text{Ent}(S)}_{\text{before}} - \underbrace{\sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Ent}(S_v)}_{\text{weighted average}}$$

What is the entropy for this dataset?

$$\text{Ent} = -\frac{12}{25} \log_2 \frac{12}{25} - \frac{13}{25} \log_2 \frac{13}{25} \approx \mathbf{0.9988}$$

$$A = \bullet: \frac{6}{12} \text{ positive} \rightarrow 1.0$$

$$A = \circ: \frac{6}{13} \text{ positive} \rightarrow 0.9957$$

$$\text{Expected: } \frac{12}{25} \cdot 1.0 + \frac{13}{25} \cdot 0.9957 \approx \mathbf{0.9977}$$

$$B = \bullet: \frac{9}{11} \text{ positive} \rightarrow 0.684$$

$$B = \circ: \frac{3}{14} \text{ positive} \rightarrow 0.750$$

$$\text{Expected: } \mathbf{0.721}$$

$$C = \bullet: \frac{3}{6} \text{ positive} \rightarrow 1.0$$

$$C = \circ: \frac{9}{19} \text{ positive} \rightarrow 0.9980$$

$$\text{Expected: } \mathbf{0.9985}$$

$$D = \bullet: \frac{3}{5} \text{ positive} \rightarrow 0.9710$$

$$D = \circ: \frac{9}{20} \text{ positive} \rightarrow 0.9928$$

$$\text{Expected: } \mathbf{0.9884}$$

A	B	C	D	
•	•	○	○	+
○	•	•	○	+
○	○	○	○	
•	○	○	•	+
○	•	○	○	+
•	○	•	○	
•	•	○	○	+
○	○	○	○	
○	○	•	○	
•	•	○	○	+
○	○	○	•	+
•	○	○	○	
•	•	•	○	+
○	•	○	○	+
○	○	○	○	
•	○	○	•	+
○	•	•	○	+
○	○	○	○	
•	○	○	○	

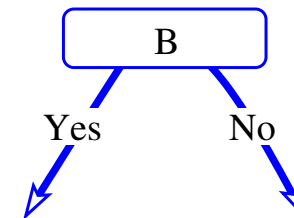
$$\text{Gain}(A) = 0.9988 - 0.9977 = \mathbf{0.0011}$$

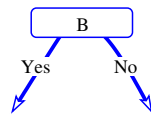
$$\text{Gain}(B) = 0.9988 - 0.7210 = \mathbf{0.2778}$$

$$\text{Gain}(C) = 0.9988 - 0.9985 = \mathbf{0.0003}$$

$$\text{Gain}(D) = 0.9988 - 0.9884 = \mathbf{0.0104}$$

Attribute B gives most information



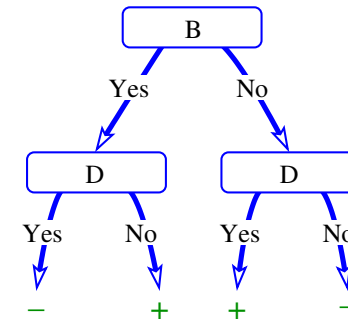


Examples where
 $B = \bullet$

A	B	C	D	
•	•	○	○	+
○	•	•	○	+
○	•	○	○	+
•	•	○	○	+
•	•	○	○	+
•	•	•	○	+
○	•	○	•	-
•	•	○	•	-
○	•	○	○	+
○	•	○	○	+
○	•	•	○	+

Examples where
 $B = \circ$

A	B	C	D	
○	○	○	○	
•	○	○	•	+
•	○	•	○	
○	○	○	○	
○	○	•	○	
○	○	○	•	+
•	○	○	○	
○	○	○	○	
○	○	○	○	
•	○	○	○	
○	○	•	○	
•	○	○	•	+
○	○	○	○	
•	○	○	○	



Possible ways of improving the decision trees

- Avoid overfitting
 - Limit the tree's height
 - Pruning
- Attributes with graded values
- Missing attribute values
- Variable cost for different attributes