# LECTURE 10: Decision tree learning

Were there some restriction in concept learning?

Example 1: Recall most specific hypothesis to these entries:

| | | | | | | Enjoy Sport |
|---|---|---|---|---|---|---|
| 1. Sunny | Warm | Normal | Strong | Cool | Change | Yes |
| 2. Cloudy | Warm | Normal | Strong | Cool | Change | Yes |
| 3. Rainy | Warm | Normal | Strong | Cool | Change | No |

$S_1$: $\langle$Sunny, Warm, Normal, Strong, Cool, Change$\rangle$

$S_2$: $\langle$ ? , Warm, Normal, Strong, Cool, Change$\rangle$

$S_3$: $\emptyset$

$\Rightarrow$ More __expressive__ hypothesis space needed!

E.g.

$\langle$Sunny, ?, ?, ?, ?, ?$\rangle$ $\vee$ $\langle$Cloudy, ?, ?, ?, ?, ?$\rangle$

In concept learning variables can be fixed or free, but this is not sufficient for all tasks (Example 1). Allowing disjunctions ($\vee$) does not work in concept learning, and therefore, we need new __learning model__ and __method__!

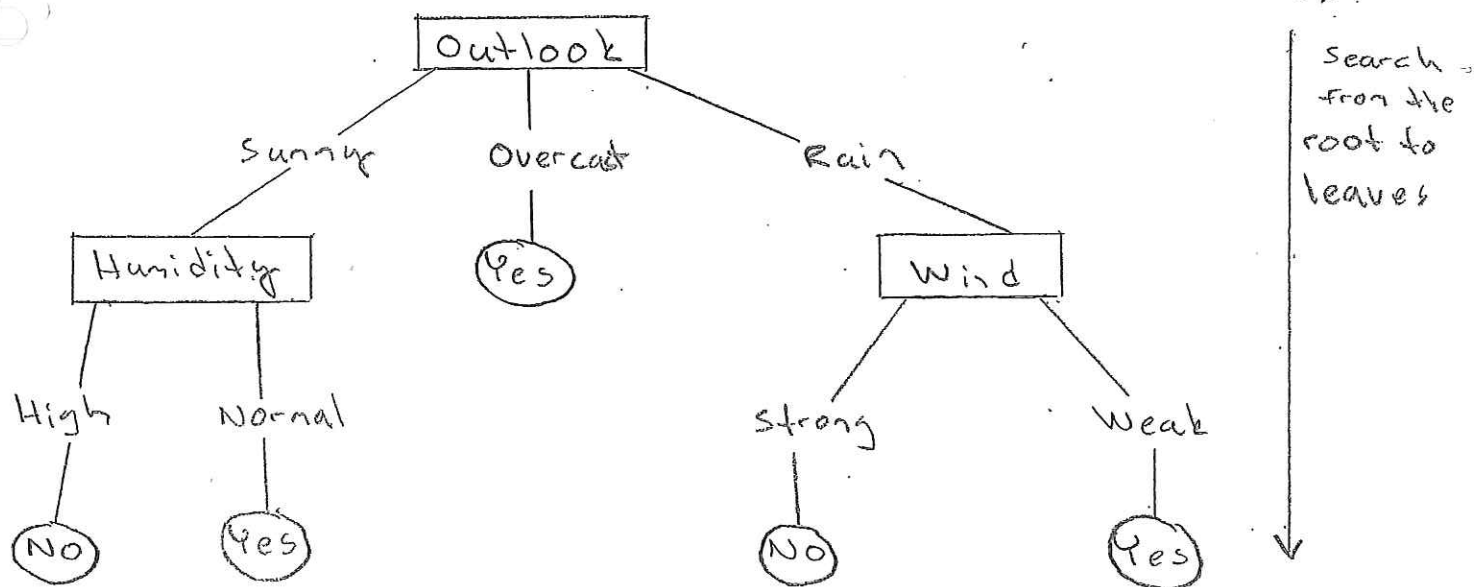1. Decision tree learning model (representation)

&lt;SLIDE&gt;

Decisio tree represents knowledge (concept) by sorting values down the tree from the root to some leaf node, which provides information if the concept is satisfied or not.

Node: Test of some attribute

Branch: Possible value of an attribute

Example 2   PlayTennis decision tree (from TABLE 3.2 in Mitchell (SLIDE))



Search from the root to leaves

Decision trees represent learnt information as a disjunction of conjunctions of constraints on the attribute values of instances

Example 3   PlayTennis as disjunctions of conjunctions (left-first search)

(Outlook == Sunny ∧ Humidity == norma)

∨ (Outlook == Overcast)

∨ (Outlook == Rain ∧ Wind == Weak)

2.ii. Decision tree learning (ID3 algorithm)

It is straightforward to implement a learning method if a single question can be answered; <SLIDE

"Which attribute should be tested at the root of the tree?"

If this can be solved, then the ID3 algorithm works (see Table 3.1 and note recursive structure)!

Top-Down induction: <SLIDE>

iii. Selecting the best attribute to be tested

We should seek answer from information theory: which attribute provides largest information gain?

⟹ Wanted information is division of example instances to two classes

A perfect attribute would be the one which divides examples to exactly positive and negative examples. The worst attribute holds equally for the both (having them completely mixed.

iii.1 Measure of homogeneity: entropy

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

$p_{\oplus}$ proportion or positive examples

$p_{\ominus}$ proportion of negative examples

Example 4    14 samples including 9 positive and 5 negative ([9+,5-]). Compute entropy (Table 3.2)

$$Entropy([9+,5-]) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

What does entropy tell? What do maximum and minimum values denote?

$P_\oplus$ and $P_\ominus$ may have values in $[0,1]$ such that $P_\oplus + P_\ominus = 1$ always.

Min. $Entropy(S) = 0$, when
$P_\oplus = 1$ or $P_\ominus = 1$ $\left(-1 \cdot \underbrace{\log_2(1)}_{=0} - 0 \cdot \log_2(0)\right)$

Max. $Entropy(S) = 1$, when
$\oplus = P_\ominus = \frac{1}{2}$ $\left(\underbrace{-\frac{1}{2} \cdot \log_2 \frac{1}{2} - \frac{1}{2} \cdot \log_2(\frac{1}{2})}_{=1}\right)$



Entropy in general represents how many bits on average are required to represent information.

For two classes ($\oplus$ and $\ominus$) the maximum is 1. For $C$ classes $\lceil \log_2 C \rceil$ and the entropy is computed as

$$Entropy(S) = \sum_{i=1}^{C} -p_i \log_2 p_i \quad \left[\begin{array}{l}\text{note addition to} \\ \text{probabilities}\end{array}\right]$$

## III.2 Information gain

Entropy measures non-homogeneity of data, i.e. how mixed the data is. Our goal in best attribute selection is to "un-mix" data. We want to reduce entropy.

  Information gain $==$ reduction in entropy caused by partitioning examples according to given attribute
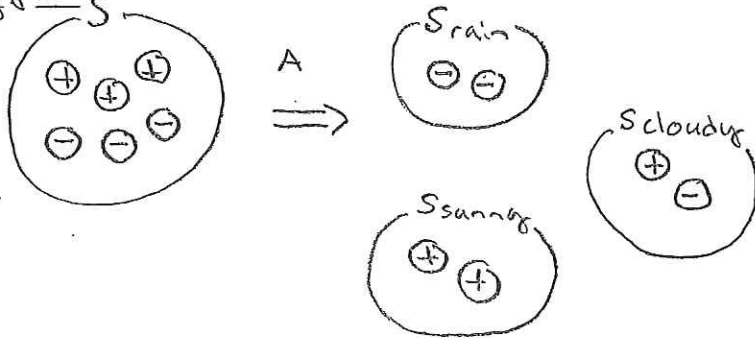
Definition 1

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

where Values(A) is all possible values of the attribute A.
and $S_v$ subset of S for which A has value $v$.

Example 5



$$Entropy(S) = 1.0$$
$$Entropy(S_{rain}) = 0.0$$
$$Entropy(S_{sunny}) = 0.0$$
$$Entropy(S_{cloudy}) = 1.0$$

$$Gain(S,A) = 1.0 - 0.0 - 0.0 - \frac{2}{6} \cdot 1.0 = \frac{4}{6} = \frac{2}{3}$$

Example 6    Computing Gain() (Cont. Example 4)

$$Values(Wind) = \{weak, strong\}$$
$$S_{weak} \leftarrow [6+, 2-]$$
$$S_{strong} \leftarrow [3+, 3-]$$

$$Gain(S, Wind) = Entropy(S) - \sum_{v \in \{weak, strong\}} \frac{|S_v|}{|S|} Entropy(S_v)$$

$$= Entropy(S) - \frac{8}{14} Entropy(S_{weak}) - \frac{6}{14} Entropy(S_{strong})$$

$$= 0.940 - \frac{8}{14} \cdot 0.811 - \frac{6}{14} \cdot 1.00 = \underline{0.048}$$

Finally, ID3 algorithm can be implemented using
Gain(S,A)-function    <SLIDE>

Example 7    ID3 for the data in TABLE 3.2

First step uses whole data S:

$Gain(S, Outlook) = 0.246$, $Gain(S, Humidity) = 0.151$, $Gain(S, Wind) = 0.048$
$Gain(S, Temp) = 0.029$ ... See Figures 3.3 and 3.4

Real example: C-Section Risk   <SLIDES>

## 3. Issues in DT Learning   <SLIDES>

Reduced-error pruning — not very essential

# RANDOMISATION IN ML (RANDOM FORESTS)

1. Intro &lt;SLIDES&gt;

2. Decision trees + Randomisation = random forests
   &lt;SLIDES&gt;