

# Einführung in die Computerlinguistik

## Decision Trees

Hinrich Schütze & Robert Zangenfeind

Centrum für Informations- und Sprachverarbeitung, LMU München

2015-11-23

This lecture is based on  
Russell and Norvig's introduction  
to artificial intelligence.

# Take-away today

- Introduction to decision trees
- (Almost) fully understand one complex machine learning model
- Basis for hands-on training and applying this model in next practical exercise
- Exam: questions testing basic understanding of decision trees, but no formulas

# Overview

1 Decision trees

2 NLTK

# Outline

1 Decision trees

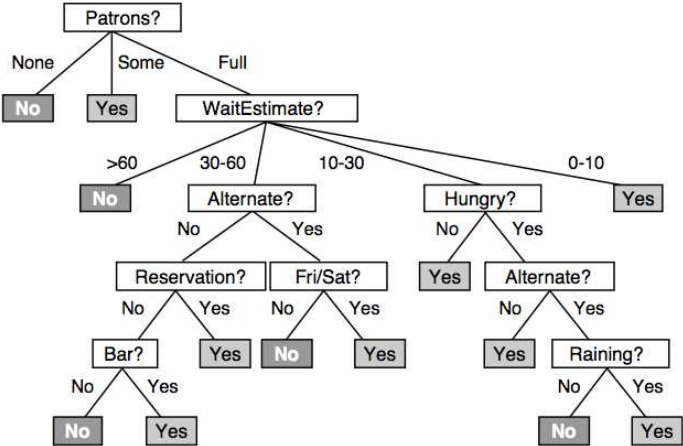
2 NLTK

# Attributes

As an example, we will build a decision tree to decide whether to wait for a table at a restaurant. The aim here is to learn a definition for the **goal predicate** *WillWait*. First we list the attributes that we will consider as part of the input:

1. *Alternate*: whether there is a suitable alternative restaurant nearby.
2. *Bar*: whether the restaurant has a comfortable bar area to wait in.
3. *Fri/Sat*: true on Fridays and Saturdays.
4. *Hungry*: whether we are hungry.
5. *Patrons*: how many people are in the restaurant (values are *None*, *Some*, and *Full*).
6. *Price*: the restaurant's price range (\$, \$\$, \$\$\$).
7. *Raining*: whether it is raining outside.
8. *Reservation*: whether we made a reservation.
9. *Type*: the kind of restaurant (French, Italian, Thai, or burger).
10. *WaitEstimate*: the wait estimated by the host (0–10 minutes, 10–30, 30–60, or >60).

# Decision tree for deciding whether to wait



**Figure 18.2** A decision tree for deciding whether to wait for a table.

(explain

one path)

# Expressiveness of decision trees

## 18.3.2 Expressiveness of decision trees

A Boolean decision tree is equivalent to a logical expression of the form

$$Goal \Leftrightarrow (Path_1 \vee Path_2 \vee \dots),$$

where each  $Path_i$  has the form

$$Path = (A_i = a_i \wedge A_j = a_j \wedge \dots),$$

that is, the goal is true if and only if there is a path through the tree that ends in a positive result. Since this is equivalent to disjunctive normal form, that means that any function in propositional logic can be expressed as a decision tree. As an example, the rightmost path in Figure 18.2 is

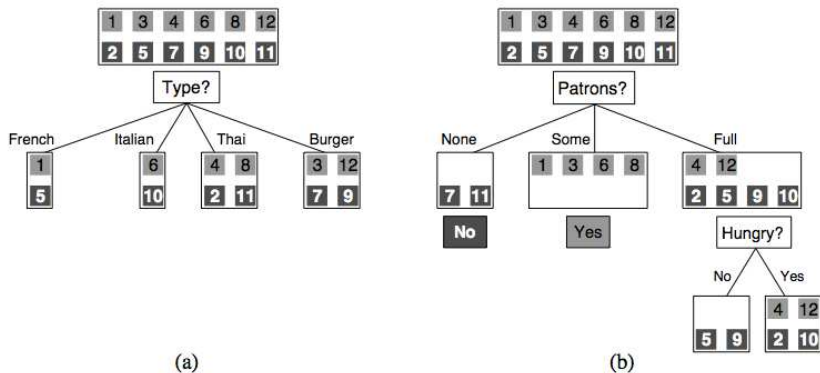
$$Path = (Patrons = Full \wedge WaitEstimate = 0-10).$$



# Training set

Example	Input Attributes										Goal
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
$x_1$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0-10</i>	$y_1 = \text{Yes}$
$x_2$	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30-60</i>	$y_2 = \text{No}$
$x_3$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Some</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_3 = \text{Yes}$
$x_4$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Thai</i>	<i>10-30</i>	$y_4 = \text{Yes}$
$x_5$	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>&gt;60</i>	$y_5 = \text{No}$
$x_6$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Italian</i>	<i>0-10</i>	$y_6 = \text{Yes}$
$x_7$	<i>No</i>	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>0-10</i>	$y_7 = \text{No}$
$x_8$	<i>No</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$</i>	<i>Yes</i>	<i>Yes</i>	<i>Thai</i>	<i>0-10</i>	$y_8 = \text{Yes}$
$x_9$	<i>No</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>Yes</i>	<i>No</i>	<i>Burger</i>	<i>&gt;60</i>	$y_9 = \text{No}$
$x_{10}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>Italian</i>	<i>10-30</i>	$y_{10} = \text{No}$
$x_{11}$	<i>No</i>	<i>No</i>	<i>No</i>	<i>No</i>	<i>None</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>0-10</i>	$y_{11} = \text{No}$
$x_{12}$	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30-60</i>	$y_{12} = \text{Yes}$

# Usefulness of attributes



**Figure 18.4** Splitting the examples by testing on attributes. At each node we show the positive (light boxes) and negative (dark boxes) examples remaining. (a) Splitting on *Type* brings us no nearer to distinguishing between positive and negative examples. (b) Splitting on *Patrons* does a good job of separating positive and negative examples. After splitting on *Patrons*, *Hungry* is a fairly good second test.

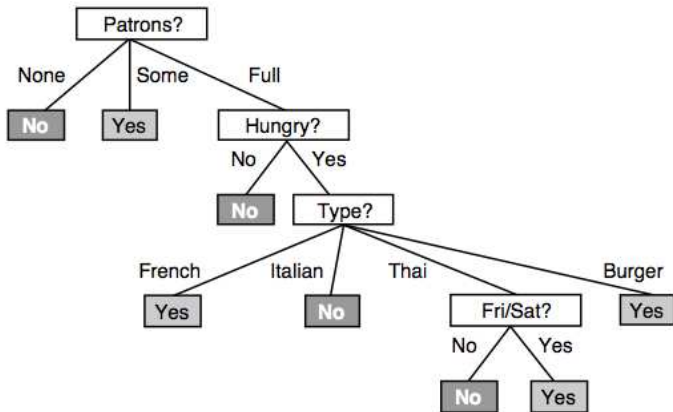
# Decision tree learning : Importance?

```
function DECISION-TREE-LEARNING(examples, attributes, parent_examples) returns  
a tree  
  
if examples is empty then return PLURALITY-VALUE(parent_examples)  
else if all examples have the same classification then return the classification  
else if attributes is empty then return PLURALITY-VALUE(examples)  
else  
  attr  $\leftarrow$   $\operatorname{argmax}_{a \in \text{attributes}}$  IMPORTANCE(a, examples)  
  tree  $\leftarrow$  a new decision tree with root test attr  
  for each value  $v_i$  of attr do  
    exs  $\leftarrow$  {e : e  $\in$  examples and e.attr =  $v_i$ }  
    subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes - attr, examples)  
    add a branch to tree with label (attr =  $v_i$ ) and subtree subtree  
  return tree
```

---

**Figure 18.5** The decision-tree learning algorithm. The function IMPORTANCE is described in Section 18.3.4. The function PLURALITY-VALUE selects the most common output value among a set of examples, breaking ties randomly.

# Induced decision tree



**Figure 18.6** The decision tree induced from the 12-example training set.

# Hand-designed vs. induced trees

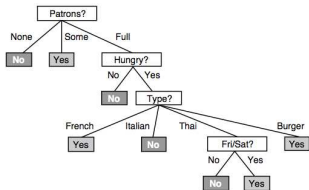


Figure 18.6 The decision tree induced from the 12-example training set.

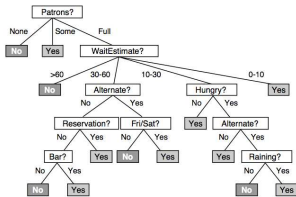


Figure 18.2 A decision tree for deciding whether to wait for a table.

# Decision tree learning : Importance?

```
function DECISION-TREE-LEARNING(examples, attributes, parent_examples) returns  
a tree  
  
if examples is empty then return PLURALITY-VALUE(parent_examples)  
else if all examples have the same classification then return the classification  
else if attributes is empty then return PLURALITY-VALUE(examples)  
else  
  attr  $\leftarrow$   $\operatorname{argmax}_{a \in \text{attributes}}$  IMPORTANCE(a, examples)  
  tree  $\leftarrow$  a new decision tree with root test attr  
  for each value  $v_i$  of attr do  
    exs  $\leftarrow$  {e : e  $\in$  examples and e.attr =  $v_i$ }  
    subtree  $\leftarrow$  DECISION-TREE-LEARNING(exs, attributes - attr, examples)  
    add a branch to tree with label (attr =  $v_i$ ) and subtree subtree  
  return tree
```

---

**Figure 18.5** The decision-tree learning algorithm. The function IMPORTANCE is described in Section 18.3.4. The function PLURALITY-VALUE selects the most common output value among a set of examples, breaking ties randomly.

# Entropy

$$\begin{aligned} H(V) &= \sum_i P(v_i) \log_2 \frac{1}{P(v_i)} \\ &= - \sum_i P(v_i) \log_2 P(v_i) \end{aligned}$$

## Entropy: Restaurant example

$$\begin{aligned} H(V) &= \sum_i P(v_i) \log_2 \frac{1}{P(v_i)} \\ &= - \sum_i P(v_i) \log_2 P(v_i) \end{aligned}$$

$$H(\text{Goal}) = -\left(\frac{p}{p+n} \log_2 \frac{p}{p+n} + \frac{n}{p+n} \log_2 \frac{n}{p+n}\right)$$

$p$  is number of positive examples (“will wait”),

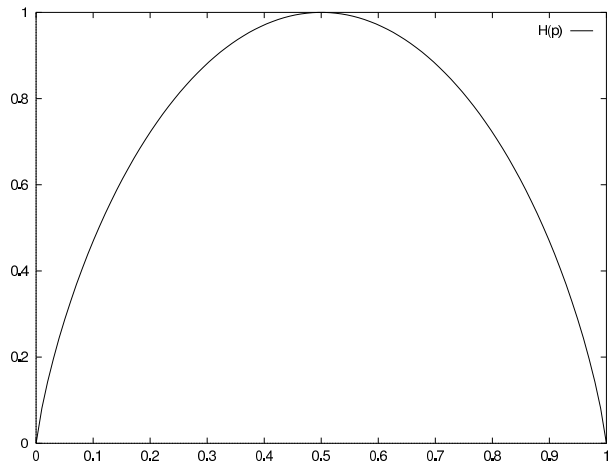
$n$  is number of negative examples (“will not wait”)



## Entropy: Restaurant example

$$\begin{aligned}H(\text{Goal}) &= -\left(\frac{p}{p+n} \log_2 \frac{p}{p+n} + \frac{n}{p+n} \log_2 \frac{n}{p+n}\right) \\&= -\left(\frac{6}{6+6} \log_2 \frac{6}{6+6} + \frac{6}{6+6} \log_2 \frac{6}{6+6}\right) \\&= -2 \frac{6}{6+6} \log_2 \frac{6}{6+6} \\&= \log_2 2 \\&= 1\end{aligned}$$

# Plot of entropy



## Notation

$$H(T(q)) = -q \log_2 q - (1 - q) \log_2(1 - q)$$

Remainder = Remaining uncertainty

$$\text{Remainder}(A) = \sum_{i=1}^{|\nu(A)|} \frac{p_i + n_i}{p + n} H\left(T\left(\frac{p_i}{p_i + n_i}\right)\right)$$

$p_i$  is the number of positive examples that have attribute value  $v_i$  ( $A = v_i$ ) and  $n_i$  is the number of negative examples that have attribute value  $v_i$  ( $A = v_i$ ).

## Information gain

$$\text{Gain}(A) = H\left(T\left(\frac{p}{p+n}\right)\right) - \text{Remainder}(A)$$

Information gain for Pt = Patrons and Tp = Type

$$\begin{aligned}\text{Gain}(T_p) &= 1 - \left[ \frac{2}{12}H\left(T\left(\frac{1}{2}\right)\right) + \frac{2}{12}H\left(T\left(\frac{1}{2}\right)\right) + \frac{4}{12}H\left(T\left(\frac{2}{4}\right)\right) + \frac{4}{12}H\left(T\left(\frac{2}{4}\right)\right) \right] \\ &= 1 - \left[ \frac{2}{12} + \frac{2}{12} + \frac{4}{12} + \frac{4}{12} \right] \\ &= 0 \text{ bits}\end{aligned}$$

$$\begin{aligned}\text{Gain}(P_t) &= 1 - \left[ \frac{2}{12}H\left(T\left(\frac{0}{2}\right)\right) + \frac{4}{12}H\left(T\left(\frac{4}{4}\right)\right) + \frac{6}{12}H\left(T\left(\frac{2}{6}\right)\right) \right] \\ &= 1 - \left[ 0 + 0 + \frac{1}{2}H\left(T\left(\frac{1}{3}\right)\right) \right] \\ &\approx 0.541 \text{ bits}\end{aligned}$$

# Entropy exercise (goal)

example	decision	type	day of week	colleague?
x <sub>1</sub>	yes	french	saturday	"let's stay"
x <sub>2</sub>	no	thai	friday	"let's go"
x <sub>3</sub>	yes	burger	saturday	"let's stay"
x <sub>4</sub>	yes	thai	sunday	"let's stay"
x <sub>5</sub>	no	french	friday	"let's go"
x <sub>6</sub>	yes	italian	sunday	"let's stay"
x <sub>7</sub>	no	burger	friday	"let's go"
x <sub>8</sub>	yes	thai	sunday	"let's stay"
x <sub>9</sub>	no	burger	friday	"not sure"
x <sub>10</sub>	no	italian	friday	"not sure"
x <sub>11</sub>	no	thai	friday	"not sure"
x <sub>12</sub>	yes	burger	sunday	"not sure"

# Decision tree learning : Importance?

**function** DECISION-TREE-LEARNING(*examples*, *attributes*, *parent\_examples*) **returns**  
a tree

**if** *examples* is empty **then return** PLURALITY-VALUE(*parent\_examples*)

**else if** all *examples* have the same classification **then return** the classification

**else if** *attributes* is empty **then return** PLURALITY-VALUE(*examples*)

**else**

$attr \leftarrow \operatorname{argmax}_{a \in \text{attributes}} \text{IMPORTANCE}(a, \text{examples})$

*tree*  $\leftarrow$  a new decision tree with root test *attr*

**for each** value  $v_i$  of *attr* **do**

$exs \leftarrow \{e : e \in \text{examples} \text{ and } e.attr = v_i\}$

*subtree*  $\leftarrow$  DECISION-TREE-LEARNING(*exs*, *attributes* - *attr*, *examples*)

add a branch to *tree* with label (*attr* =  $v_i$ ) and subtree *subtree*

**return** *tree*

---

**Figure 18.5** The decision-tree learning algorithm. The function IMPORTANCE is described in Section 18.3.4. The function PLURALITY-VALUE selects the most common output value among a set of examples, breaking ties randomly.



# Outline

1 Decision trees

2 NLTK

NLTK decision tree demo

# Take-away today

- Introduction to decision trees
- (Almost) fully understand one complex machine learning model
- Basis for hands-on training and applying this model in next practical exercise
- Exam: questions testing basic understanding of decision trees, but no formulas