

Machine Learning: Decision Trees

by Elena Battini Sönmez
İstanbul Bilgi University

- based on: '*Artificial Intelligence, a Modern Approach*', by Stuart Russell and Peter Norvig, Prentice Hall.

by Elena Battini Sönmez
İstanbul Bilgi University

Decision Trees:

1. Def: Decision Trees: vector of attribute values \rightarrow decision
2. A decision tree reaches its decision by performing a sequence of tests
3. Each internal node of the tree corresponds to a test on the values of one of the input attributes.
4. All branches from a node are labeled with the possible values of the attribute
5. Each leaf node is a decision
6. Boolean classification: positive and negative samples

Example: wait at the restaurant?

Let us build a decision tree to decide whether the agent has to wait for a table at a restaurant:

Attributes:

- *Alternate*: whether there is a suitable alternative restaurant nearby
- *Bar*: whether the restaurant has a comfortable bar area to wait in
- *Fri/Sat*: true on Friday and Saturday
- *Hungry*: whether we are hungry
- *Patrons*: how many people are in the restaurant (values: None, Some, Full)
- *Price*: the restaurant price range (\$, \$\$, \$\$\$)
- *Raining*: whether it is raining outside
- *Reservation*: whether we made a reservation
- *Type*: the kind of restaurant (French, Italian, Tai, Burger)
- *Wait Estimate*: the wait estimated by the host (0-10 minutes, 10-30, 30-60, >60)

Example 1: 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>>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>>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>	$y_{12} = \text{Yes}$

Figure 18.3 Examples for the restaurant domain.

Example: Decision Tree (1st solution)

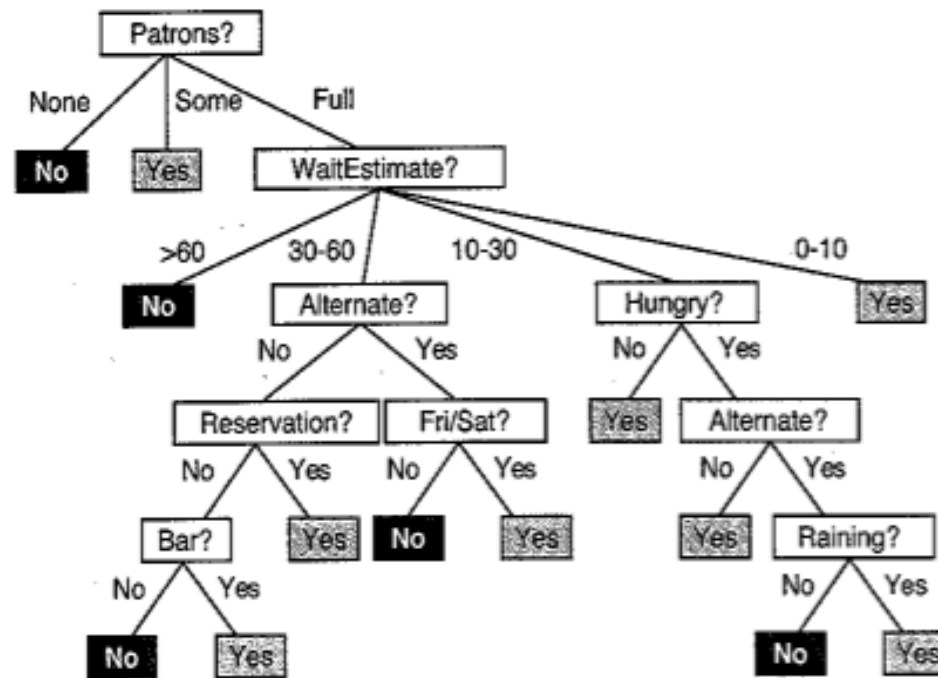


Figure 18.2 A decision tree for deciding whether to wait for a table.

Example: Important attributes

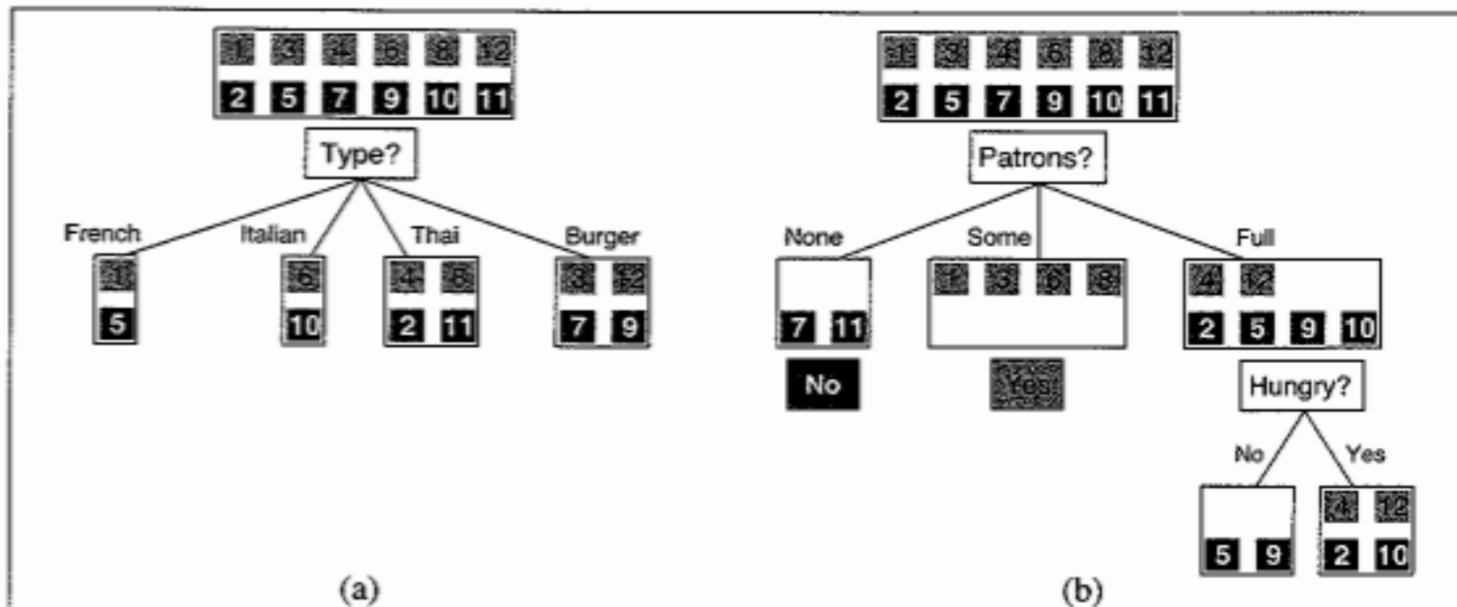


Figure 18.4 Splitting the examples by testing on attributes. At each node we show the positive (light boxes) and negative (dark boxes) examples remaining. (a) Splitting on *Type* brings us no nearer to distinguishing between positive and negative examples. (b) Splitting on *Patrons* does a good job of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

Measure of Impurity: GINI

- **Gini index** for a given node t :

$$GINI(t) = 1 - \sum_j p(j \setminus t)^2$$

where $p(.)$ is the relative frequency of class j at node t

- Maximum value of Gini index: non-homogeneous node, high degree of impurity, when records are equally distributed among all classes
- Minimum value of Gini index (0): homogeneous node, low degree of impurity, when all records belong to one class

Examples:

C1	0
C2	6

$$GINI(t) = 1 - \sum_j p(j|t)^2$$

$$P(C1)=0/6=0 \quad P(C2)=6/6=1$$
$$Gini=1-P(C1)^2-P(C2)^2=1-0-1=0$$

C1	1
C2	5

$$P(C1)=1/6 \quad P(C2)=5/6$$
$$Gini=1-P(C1)^2-P(C2)^2=1-(1/6)^2 - (5/6)^2=0.278$$

C1	2
C2	4

$$P(C1)=2/6 \quad P(C2)=4/6$$
$$Gini=1-P(C1)^2-P(C2)^2=1-(2/6)^2 - (4/6)^2=0.444$$

Splitting based on Gini index:

When a node t is split into k partitions (children), the quality of split is computed as:

$$GINI_{split} = \sum_{i=1}^k \frac{n_i}{n} GINI(i)$$

where: n_i = number of records at child i ,
 n = number of records at node t

The best splitting attribute is the one with minimum $GINI_{split}$

Example of calculation:

- Which is the better attribute for the root node of the decision tree of the given example? Let us restrict our interest into the “Type” and “patrons” attributes?
- Which is the $\text{Gini}_{\text{split}}(\text{type})$?
- Which is the $\text{Gini}_{\text{split}}(\text{patrons})$?

Alternative Splitting Criteria based on Info:

Entropy at a given node t :

$$Entropy(t) = - \sum_j p(j|t) \log p(j|t)$$

where $p(.)$ is the relative frequency of class j at node t

- It measures the homogeneity of a node
- Maximum value: when records are equally distributed among all classes implying least information
- Minimum value: when all records belong to one class, implying most information
- Entropy based computations are similar to the GINI index computations

by Elena Battini Sonmez
Istanbul Bilgi University

Splitting based on info:

Information Gain:

$$Gain_{split} = Entropy(parent_node) - \left(\sum_{i=1}^k \frac{n_i}{n} Entropy(i) \right)$$

where the parent node is split into k partitions, and n_i is number of records in partition i

- Choose the attribute that maximizes the GAIN
- Disadvantage: Tends to prefer splits that result in large number of partitions, each being small but pure

Exercises:

- Considering the training set of example1, use the Gini index to build the best decision tree
- Considering the training set of the example1, use the entropy to build the best decision tree