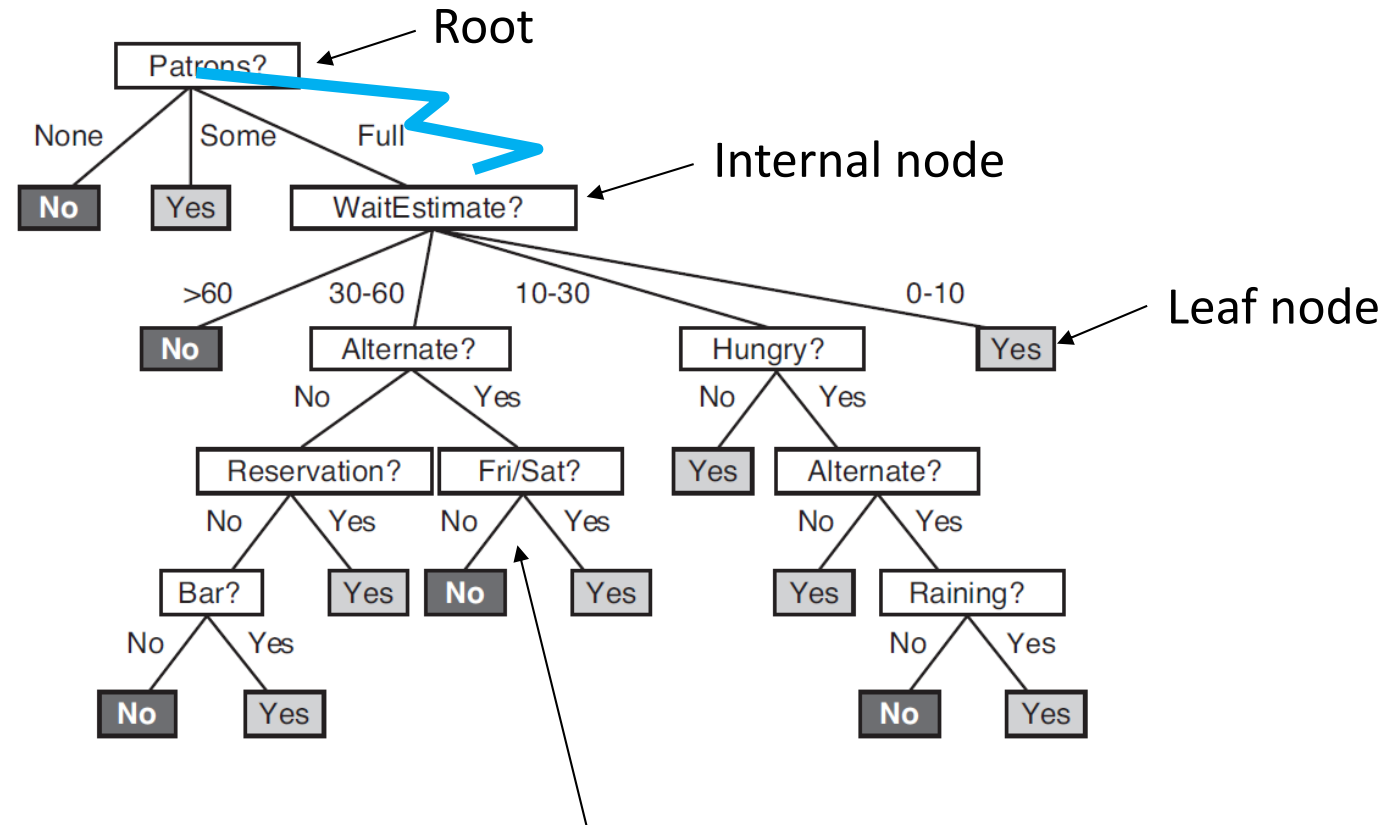


Chapter 9 – Tree

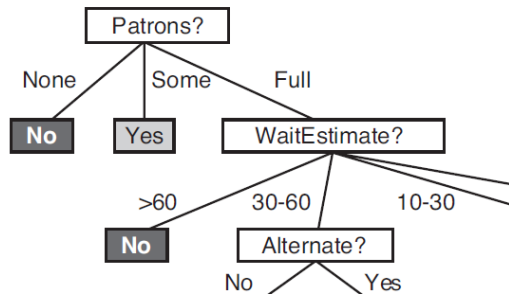
Instructor: Dr. Hongfu Liu
Email: hongfuliu@brandeis.edu

A Decision Tree Example

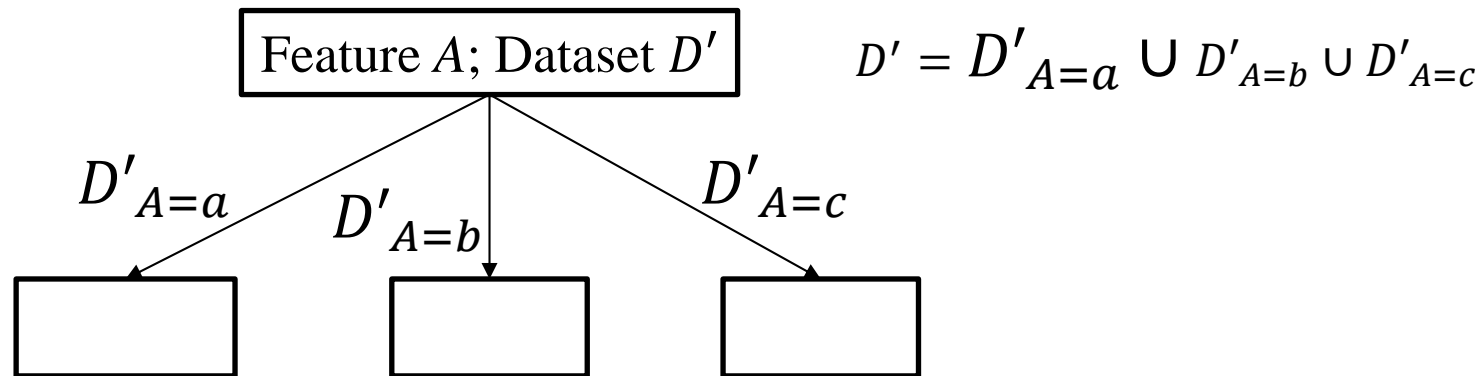


$(\text{Patrons} = \text{Full}) \wedge (\text{WaitEstimate} = 30-60) \wedge (\text{Alternate} = \text{Yes}) \wedge (\text{Fri/Sat} = \text{No})$

Basic Idea of Decision-Tree Learning Algorithms



1. Recursively builds a tree, starts with a root node and the full dataset D .
2. At each node v and the associated dataset $D' \subseteq D$
 - Select a feature A .
 - Uses A to make local decision at v . If decide v as a leaf node, stop.
 - Divide D' into subsets accordingly to A , and create one child node of v for each subset.
 - Repeat 2 for each child node



An Example Dataset

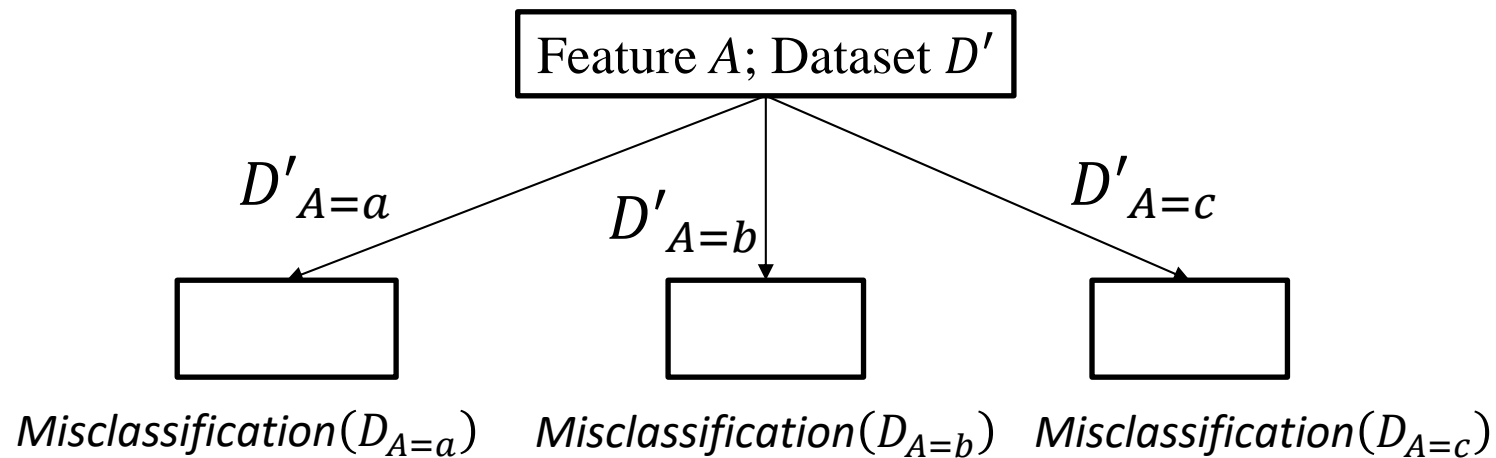
Example	Attributes										Wait
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	
1	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Some</i>	<i>\$\$\$</i>	<i>No</i>	<i>Yes</i>	<i>French</i>	<i>0-10</i>	<i>Yes</i>
2	<i>Yes</i>	<i>No</i>	<i>No</i>	<i>Yes</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Thai</i>	<i>30-60</i>	<i>No</i>
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>	<i>Yes</i>
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>	<i>Yes</i>
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>	<i>No</i>
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>	<i>Yes</i>
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>	<i>No</i>
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>	<i>Yes</i>
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>	<i>No</i>
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>	<i>No</i>
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>	<i>No</i>
12	<i>Yes</i>	<i>Yes</i>	<i>Yes</i>	<i>No</i>	<i>Full</i>	<i>\$</i>	<i>No</i>	<i>No</i>	<i>Burger</i>	<i>30-60</i>	<i>Yes</i>

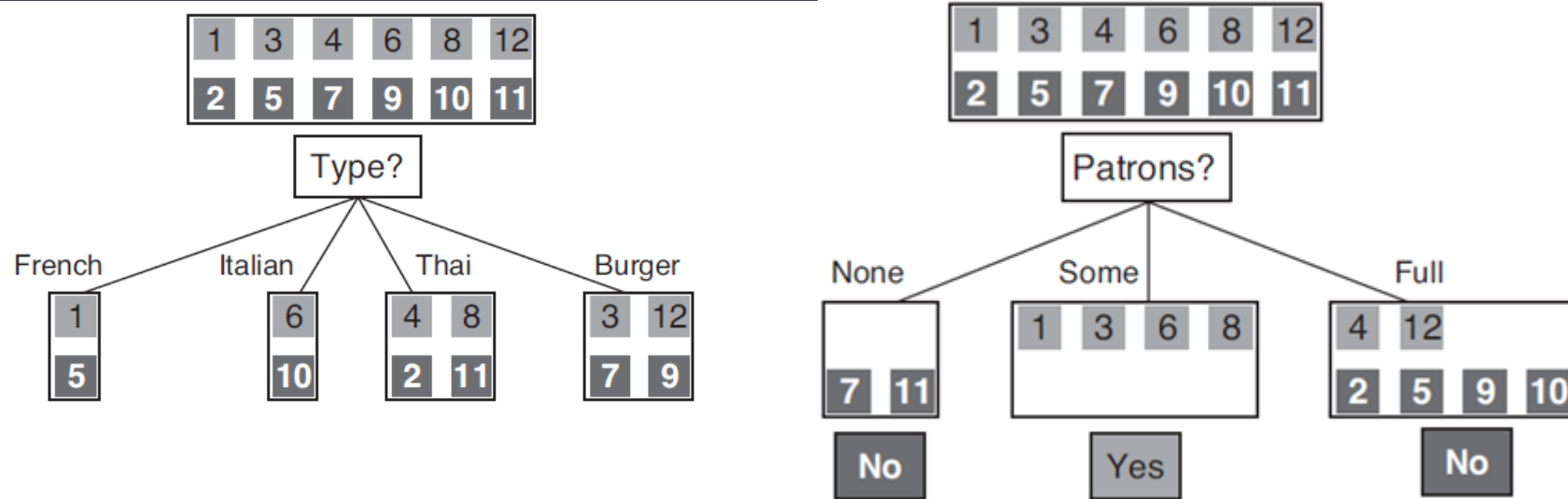
Select the Best Feature I

Misclassification error

$$\sum_{a \in \text{Values}(A)} \frac{|D_{A=a}|}{|D|} \text{Misclassification}(D_{A=a})$$

$\text{Misclassification}(D_{A=a})$ can be as simple as the error of majority voting





$$\text{Misclassification}([\text{ }]) = 0.5$$

$$\text{Misclassification}(\text{Type}) = \frac{2}{12} \times 0.5 + \frac{2}{12} \times 0.5 + \frac{4}{12} \times 0.5 + \frac{4}{12} \times 0.5 = 0.5$$

$$\text{Misclassification}(\text{Patrons}) = \frac{2}{12} \times 0 + \frac{4}{12} \times 0 + \frac{6}{12} \times \frac{2}{6} = \frac{1}{6}$$

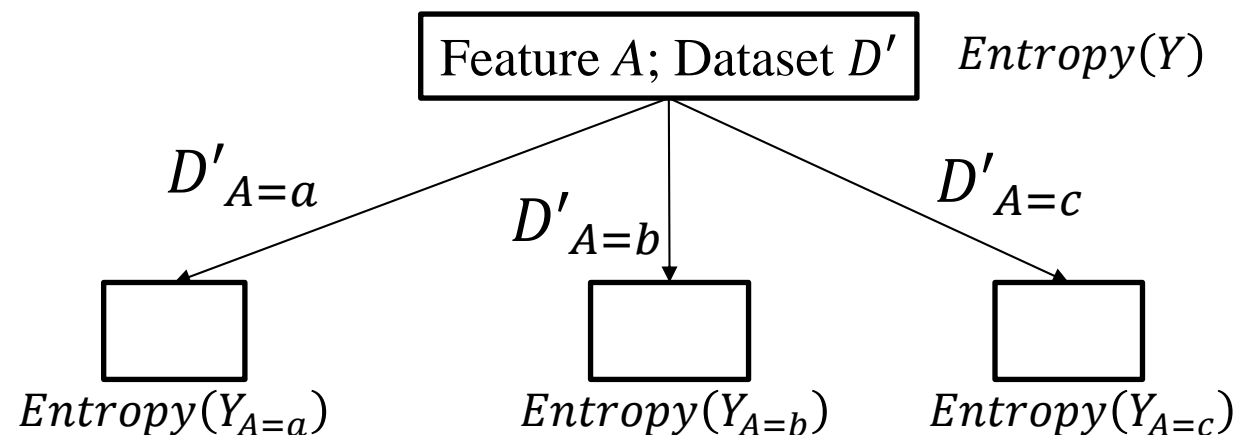
Select the Best Feature II

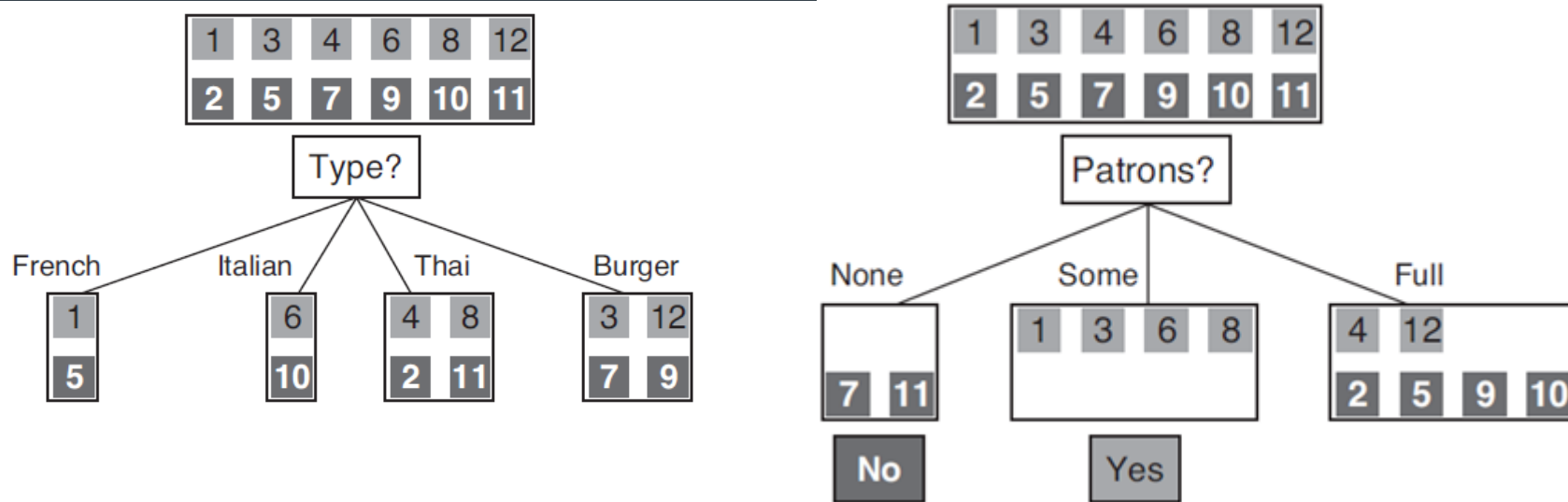
- Entropy – characterizes the amount of uncertainty.

$$Entropy(Y) = \sum_y P(y) \log_2 \frac{1}{P(Y)} = P(+)\log_2 \frac{1}{P(+)} + P(-)\log_2 \frac{1}{P(-)}$$

- Information Gain

$$Gain(Y, A) = Entropy(Y) - \sum_{a \in Values(A)} \frac{|D_{A=a}|}{|D|} Entropy(Y_{A=a})$$





$$Gain([]) = 1 - \frac{12}{12} H\left(\frac{1}{2}\right) = 0$$

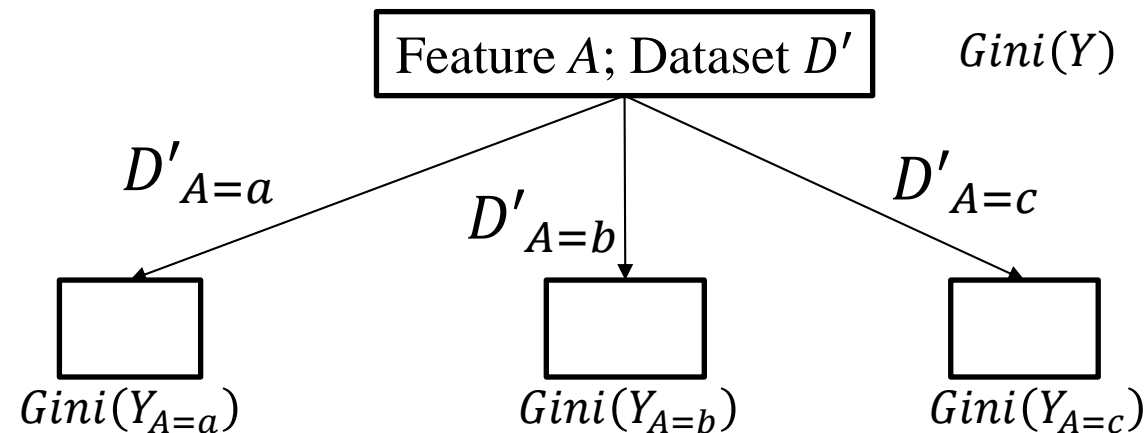
$$Gain(Type) = 1 - \left[\frac{2}{12} H\left(\frac{1}{2}\right) + \frac{2}{12} H\left(\frac{1}{2}\right) + \frac{4}{12} H\left(\frac{1}{2}\right) + \frac{4}{12} H\left(\frac{1}{2}\right) \right] = 0$$

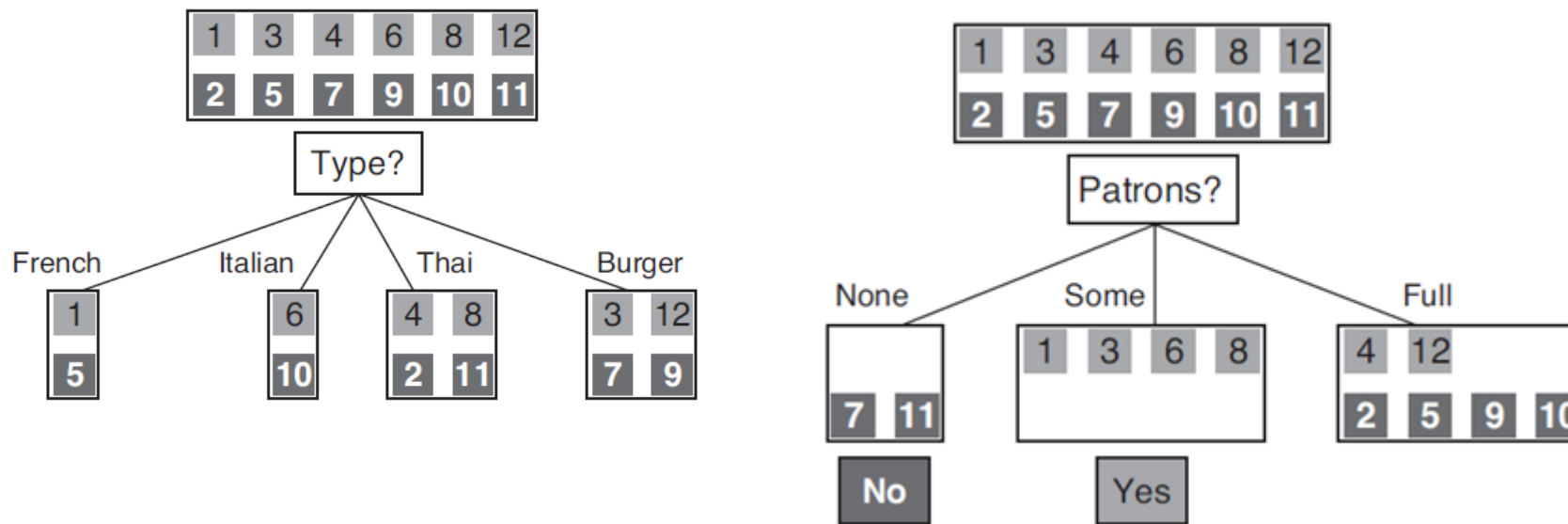
$$Gain(Patrons) = 1 - \left[\frac{2}{12} H\left(\frac{0}{2}\right) + \frac{4}{12} H\left(\frac{4}{4}\right) + \frac{6}{12} H\left(\frac{2}{6}\right) \right] = 0.54$$

Select the Best Feature III

Impurity $Gini(Y) = 1 - \sum_y P(y)^2$

Splitting based on attribute A $\sum_{a \in \text{Values}(A)} \frac{|D_{A=a}|}{|D|} Gini(Y_{A=a})$





$$Gini([]) = 1 - 0.5^2 - 0.5^2 = 0.5$$

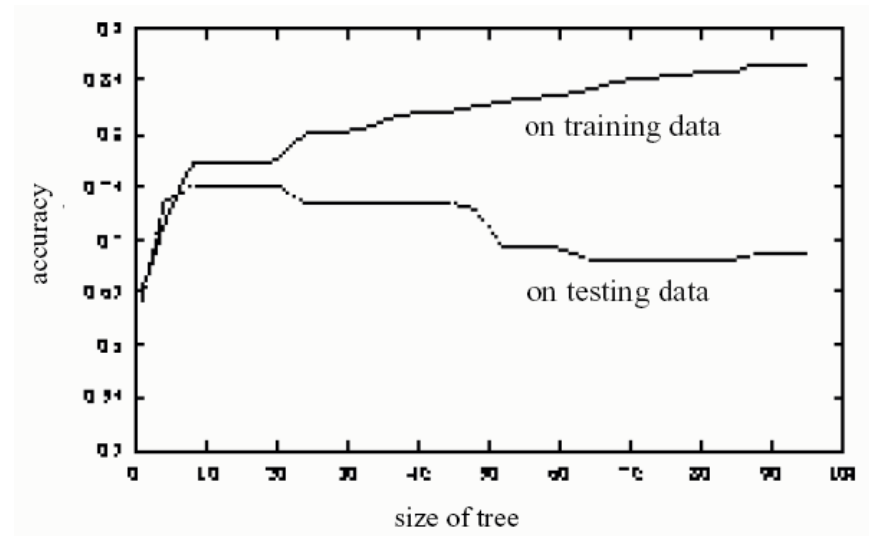
$$Gini(Type) = \frac{2}{12} \times 0.5 + \frac{2}{12} \times 0.5 + \frac{4}{12} \times 0.5 + \frac{4}{12} \times 0.5 = 0.5$$

$$Gini(Patrons) = \frac{2}{12} \times 0 + \frac{4}{12} \times 0 + \frac{6}{12} \times \left(1 - \frac{1}{9} - \frac{4}{9}\right) = \frac{2}{9}$$

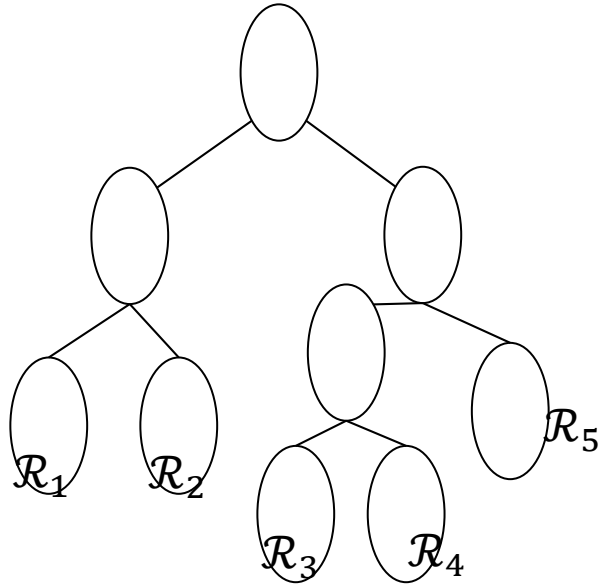
Overfitting Problem

Usually, the training data set is relatively small with respect to the complexity of the model

Large tree, some leaves only represent a tiny set of observations. (not robust to noise)



Pruning by Examining Decision Regions



Each leaf represents a decision region \mathcal{R}_t

The optimal prediction for \mathcal{R}_t is

$$h_t = \frac{1}{|\mathcal{R}_t|} \sum_{\vec{x}_n \in \mathcal{R}_t} y_n$$

The error of \mathcal{R}_t is

$$Q_t = \sum_{\vec{x}_n \in \mathcal{R}_t} \text{error}(h_t, y_n)$$

Objective function:

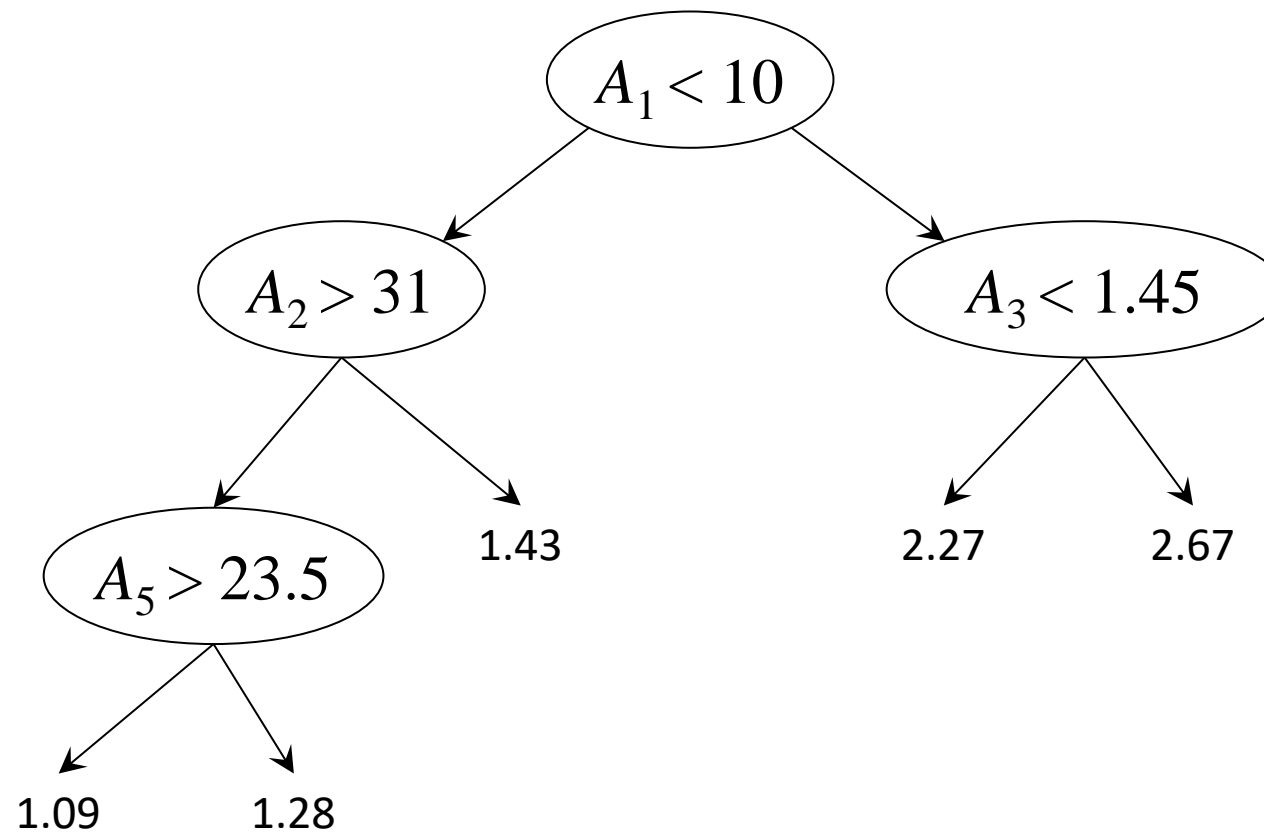
$$C = \sum_{t=1}^L Q_t + \lambda L$$

L is the number of regions

Note: after pruning a node, the corresponding decision region is merged with another region whose Q_t needs to be updated.

- Continuous attributes
- Continuous output (Regression)
 - Leaves will be the regression values rather than the predicted classification values
 - How to choose a feature for an internal node?

For example, a feature that results in the biggest reduction in regression error.



- Missing attributes
 - Impute: fill in the missing values, e.g., using the mean of that attribute over the non-missing observations.
 - Make “missing” a category.
 - Explore the correlations between attributes.

- Tree is optimal at each split – it may not be globally optimal.
- Building tree as gradient search.

