

ML-101:

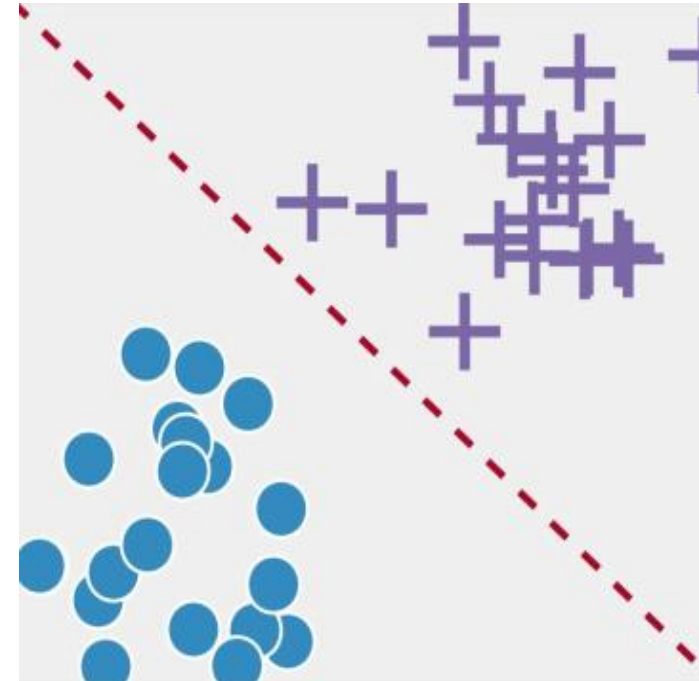
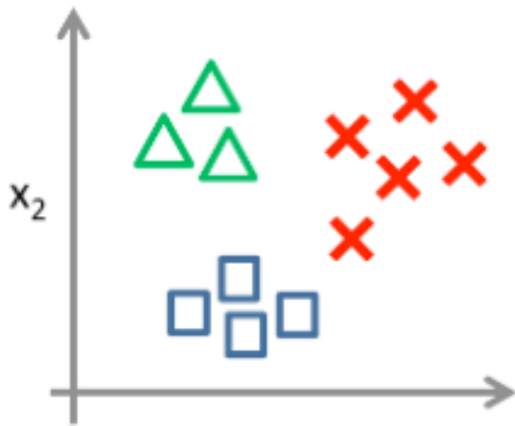
Classification, Random Forest Algorithm

BY SARTHAK CONSUL

Classification

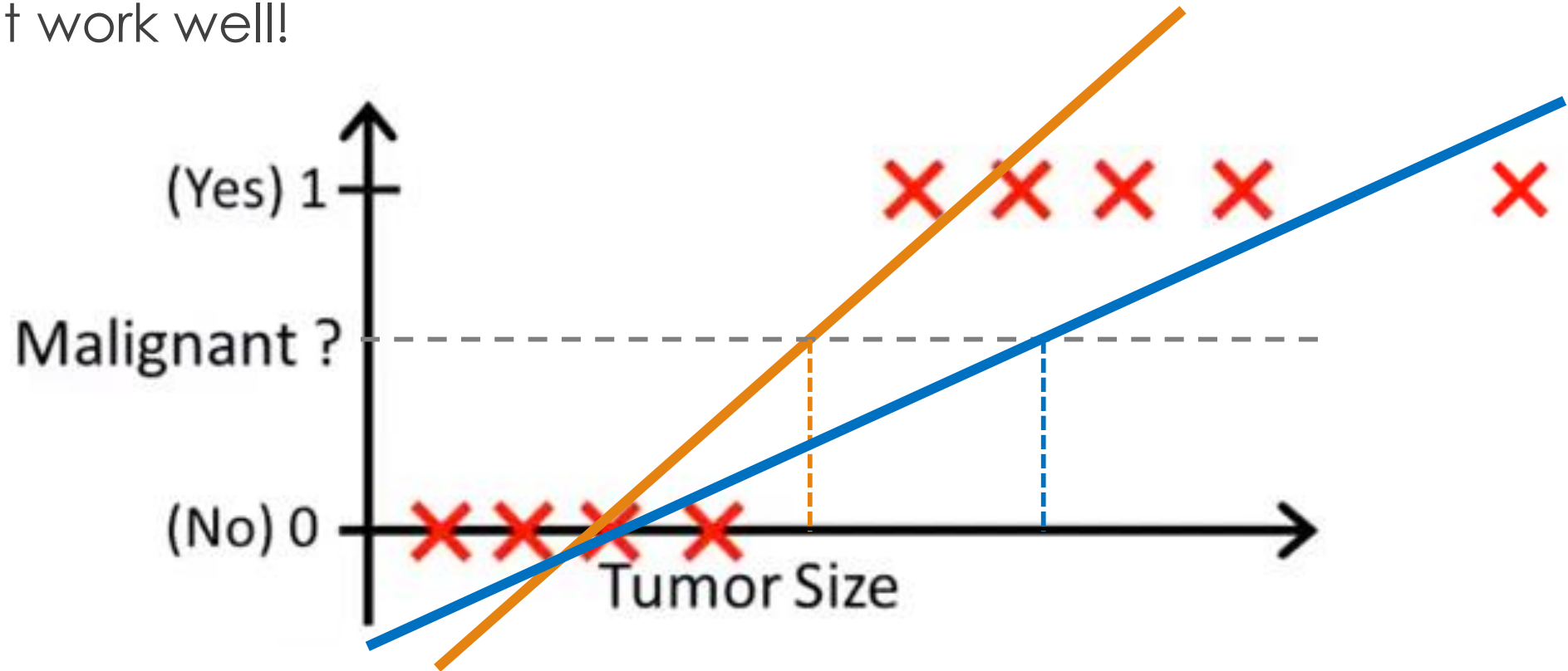
- ❖ $Y=\{0,1\}$ – Binary
- ❖ $Y=\{0,1,\dots,k\}$ – Multiclass

Eg. Spam filters, Image Recognition, Tumour type



Attempt 1: Linear Regression

❖ Won't work well!

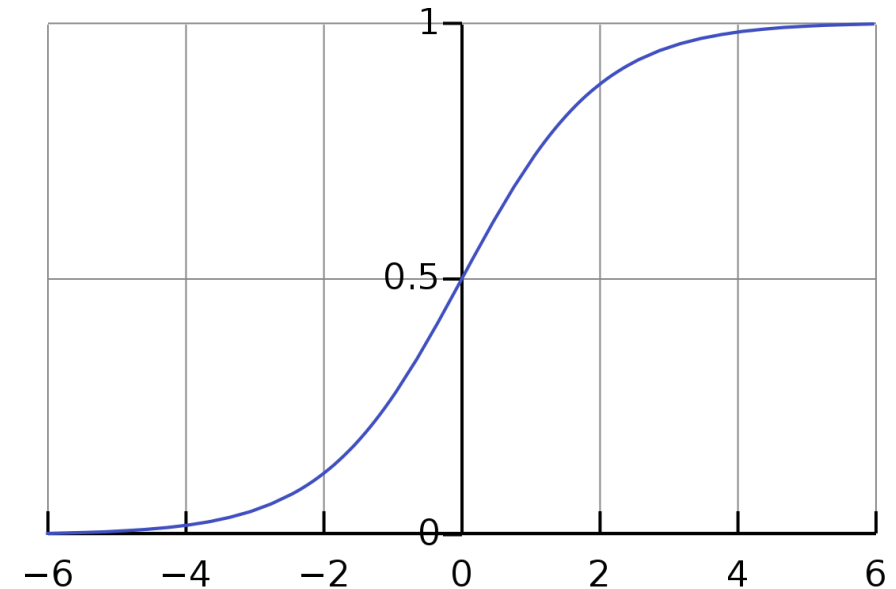


Source: Coursera - Machine Learning by Andrew Ng

A Better Way: Logistic Regression

❖ The sigmoid Function: $g(z) = \frac{1}{1 + e^{-z}}$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$



- ❖ Classification on the basis of $h_{\theta}(x)$ compared to 0.5
- ❖ Multiclass extension: One-vs-All

Cost Function for Logistic Regression

- ❖ Loss should be of the form that penalizes wrong sign of prediction and actual value i.e.

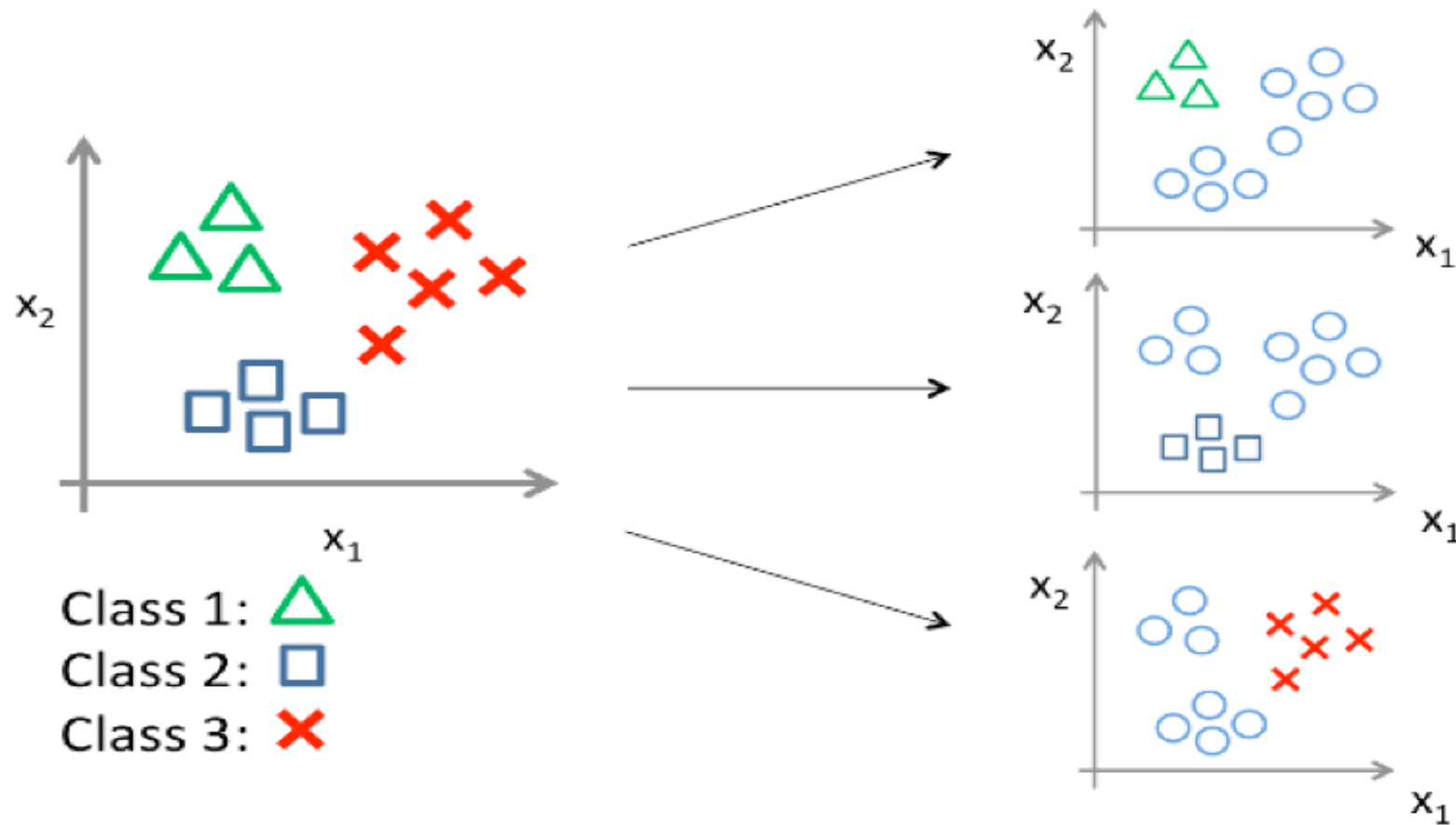
$$L(y\hat{y}) = (1 - \text{sign}(y\hat{y}))/2$$

- ❖ MSE won't work → non convex
- ❖ Instead objective function (to minimize) is:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m -y^{(i)} \log h_{\theta}(x^{(i)}) - (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

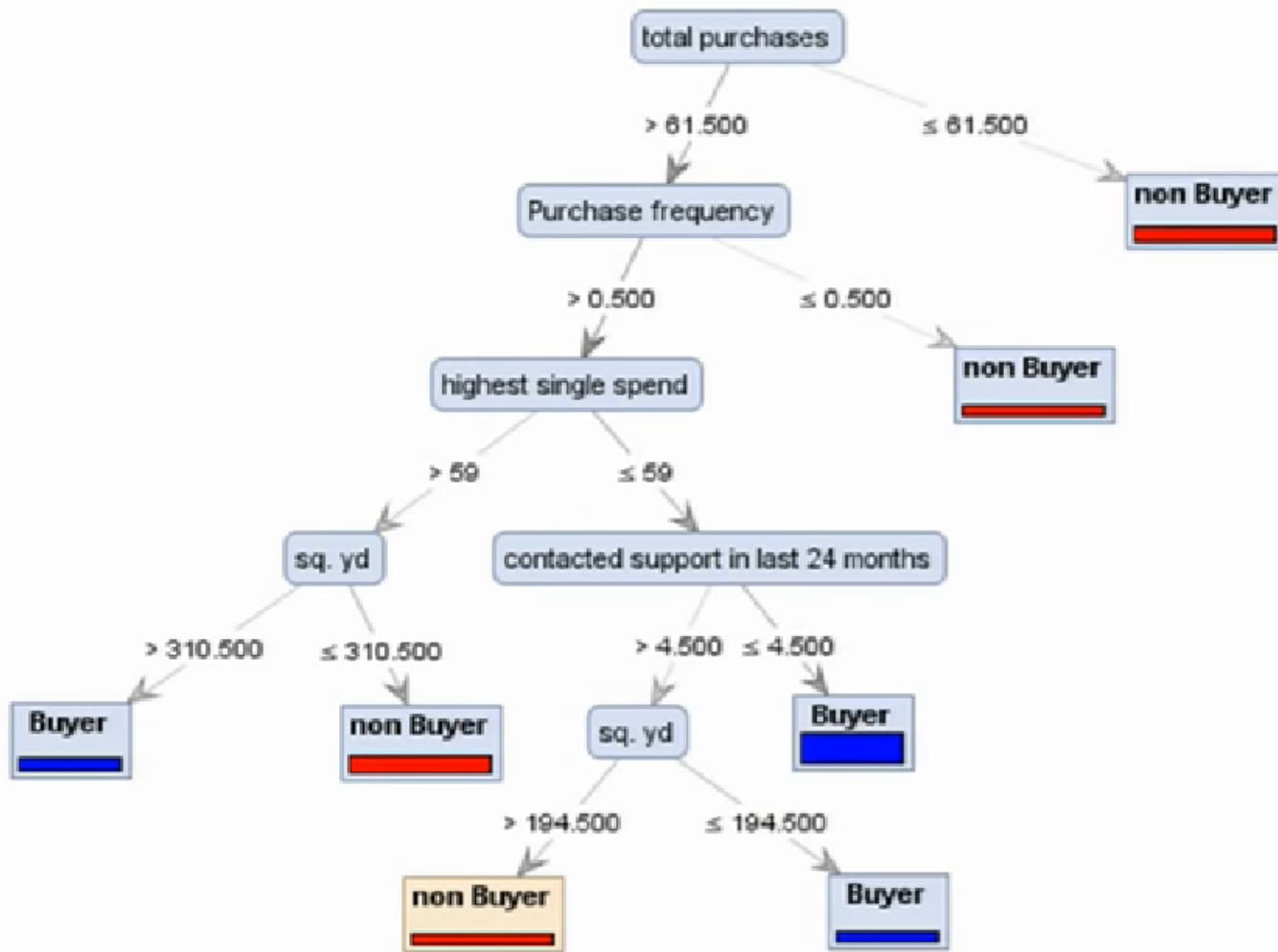
- ❖ Gradient Descent $\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$

Multiclass Classification: One-vs-All

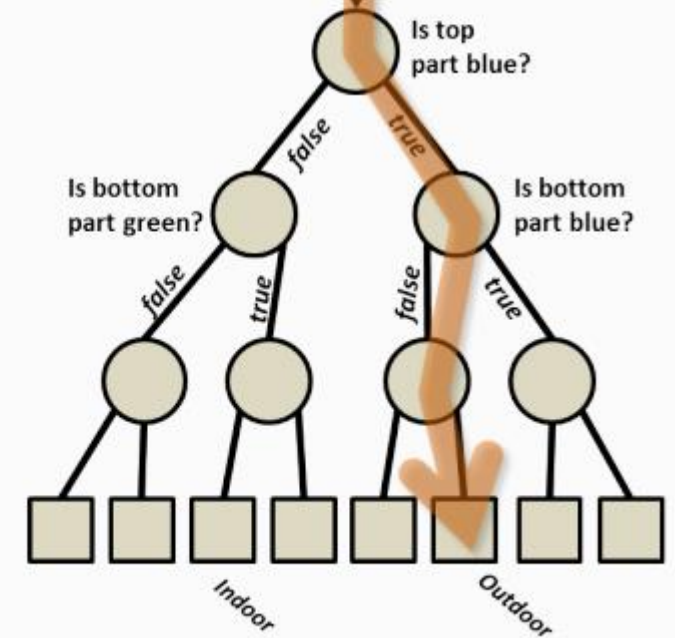


Decision Tree

- ❖ Used a lot 2-3 decades ago
- ❖ Fell out of fashion as they tend to not generalize well
- ❖ Trees have large variance, averaging out many trees reduces the variance
- ❖ Modified to make very powerful algorithms (eg. Random Forest/ Decision Forest)



A decision tree



Source: MSR Tutorial on decision forests by Criminisi et al, 2011

Constructing a Decision Tree

Example	Attributes										Target
	<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>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0-10</i>	<i>T</i>
X_2	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30-60</i>	<i>F</i>
X_3	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>T</i>
X_4	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10-30</i>	<i>T</i>
X_5	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>>60</i>	<i>F</i>
X_6	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0-10</i>	<i>T</i>
X_7	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>F</i>
X_8	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0-10</i>	<i>T</i>
X_9	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>>60</i>	<i>F</i>
X_{10}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10-30</i>	<i>F</i>
X_{11}	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0-10</i>	<i>F</i>
X_{12}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30-60</i>	<i>T</i>

Source: Artificial Intelligence: A Modern Approach by Stuart Russel and Peter Norvig

Entropy and Information Gain

- ❖ For a training set containing p positive and n negative examples,

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

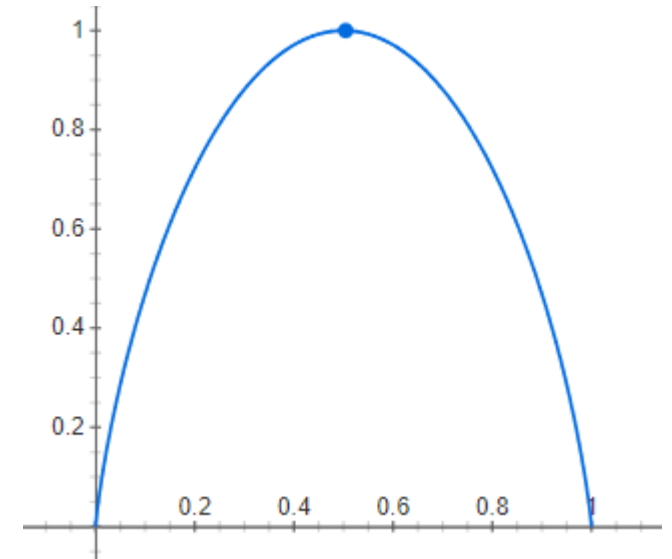
- ❖ Attribute A divides the training set into K subsets
 - ❖ Expected Entropy remaining after A , EH

$$EH(A) = \sum_{i=1}^K \frac{p_i + n_i}{p+n} H\left(\frac{p_i}{p_i + n_i}\right)$$

- ❖ Information Gain (reduction in entropy), I

$$I(A) = H\left(\frac{p}{p+n}, \frac{n}{p+n}\right) - EH(A)$$

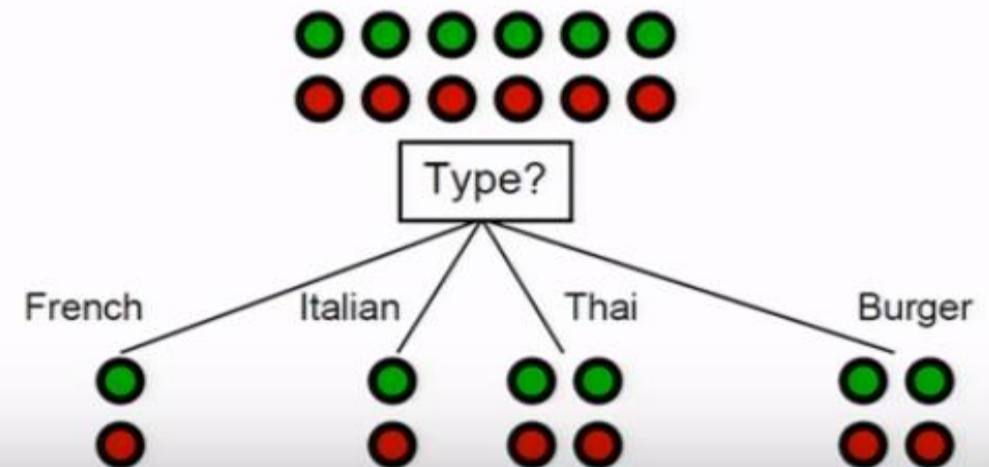
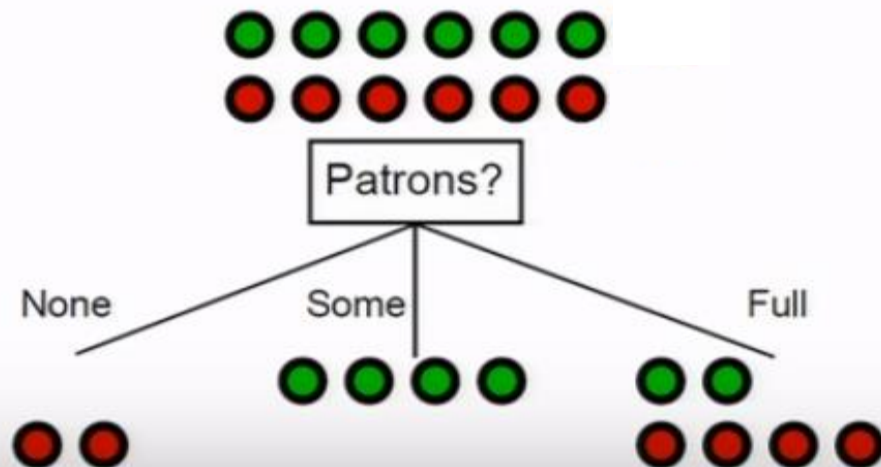
- ❖ Advanced Idea: Gaussians to decide



Choosing the attributes

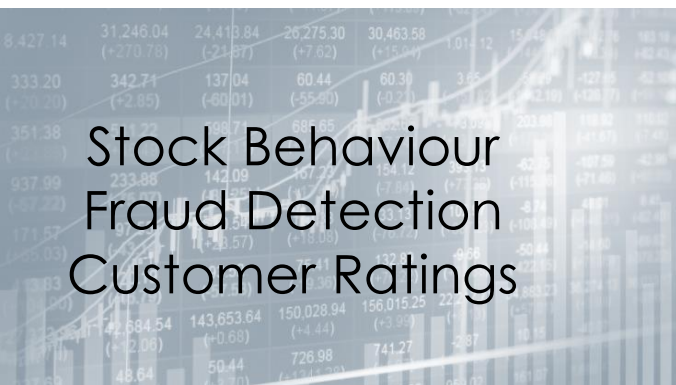
$$I(Patrons) = 1 - \left[\frac{2}{12} H(0,1) + \frac{4}{12} H(1,0) + \frac{6}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right] = .0541 \text{ bits}$$

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



Random Forest Algorithm

- ❖ One of the most popular powerful supervised learning algorithm
- ❖ Can perform both regression and classification tasks
- ❖ It builds multiple decision trees and merges them together to get a more accurate and stable prediction
- ❖ Uses:



Growing a Forest from Trees

- ❖ $d = \text{\# of Features}$, $n = \text{\# of data points}$
- ❖ Pick any 2 features at **RANDOM**, split according to it \rightarrow RANDOM TREE
 - Reduces computation time
 - Randomness is introduced
- ❖ Take trees and average them to make a FOREST

Growing a Forest from Trees: the Algo

(Bagging) Bootstrap Aggregating

For $b=1$ to T :

- a. Draw a *bootstrap sample*, Z^* of size N from training data
- b. Grow a random-forest tree, T_b from Z^*
(via recursion – till min node size is reached)

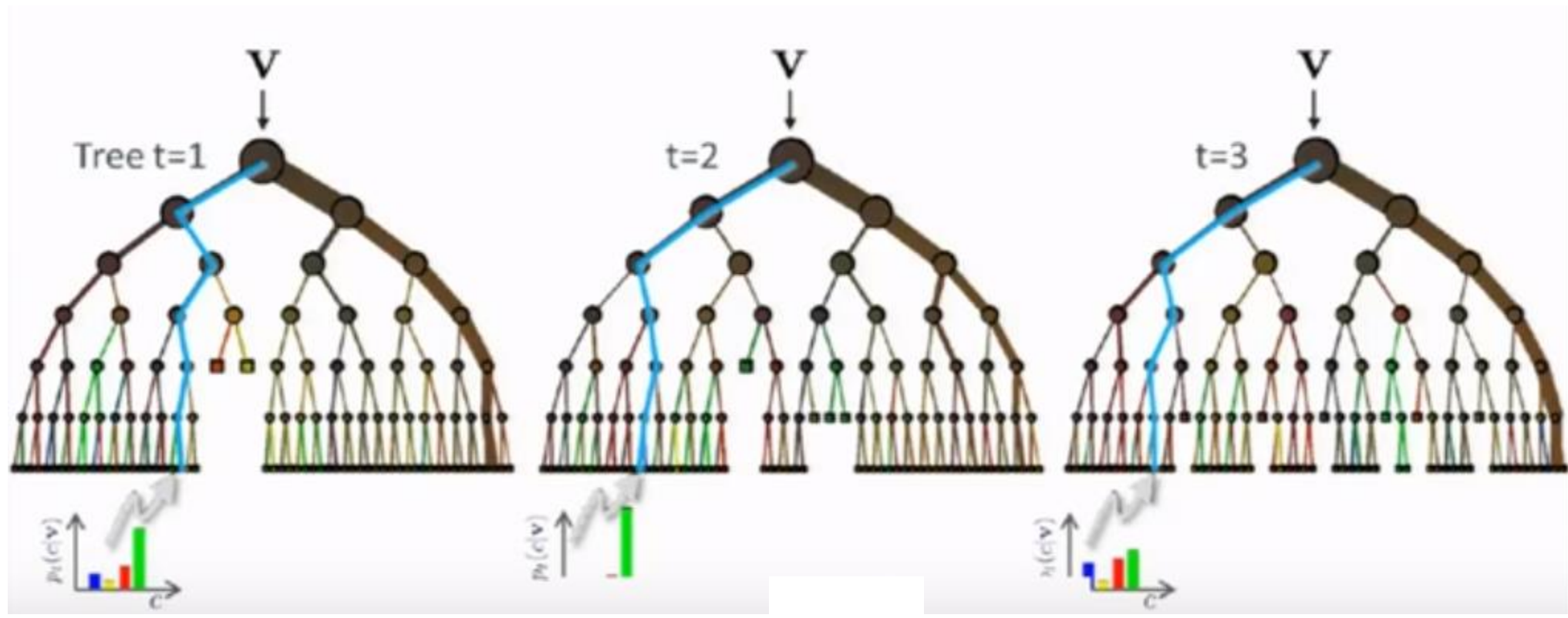
Ensemble of trees $\{T_b\}$

Word of the day: *Bootstrap Sample*
Sampling uniformly at random, with replacement

During testing,

- i. each point is passed through all trees
- ii. Average out the result (AM or GM), Max votes, etc.

Source: The Elements of Statistical Learning by Friedman, Tibshirani, and Hastie



$$p(c|\mathbf{v}) = \frac{1}{T} \sum_{t=1}^T p_t(c|v)$$

Source: Criminisi et al, 2011