# Decision Trees and Random Forests

A decision tree is a hierarchical (tree-like) classifier containing two kinds of nodes viz., decision nodes (internal nodes) and leaf nodes.

The internal nodes constitute tests on the features of the dataset and the leaf nodes represent the outcomes of the decisions (interms of class labels)

Defn: A decision tree is a graphical representation of all possible solutions to a decision based on certain conditions.

Process

① The root node contains entire population (dataset)

② Based on decisions, the root nodes is split into two or more smaller groups.

③ If uniformity (homogeneity) of a group is observed, then the group is made into a leaf node with the common class as its label. Otherwise, the splitting of the dataset continues.

④ Repeat step 4 until all groups are divided ultimately into leaf nodes.

# Example:

| Outlook | Temperature | Humidity | Windy | Play |
|---------|-------------|----------|-------|------|
| sunny | hot | high | false | no |
| sunny | hot | high | true | no |
| overcast | hot | high | false | yes |
| rainy | mild | high | false | yes |
| rainy | cool | normal | false | yes |
| rainy | cool | normal | true | no |
| overcast | cool | normal | true | yes |
| sunny | mild | high | false | no |
| sunny | cool | normal | false | yes |
| rainy | mild | normal | false | yes |
| sunny | mild | normal | true | yes |
| overcast | mild | high | true | yes |
| overcast | hot | normal | false | yes |
| rainy | mild | high | true | no |

Task : Predict whether you can play or not given the
    day's attributes

? - Which attribute to pick first?
    outlook / temperature / humidity / windy?

A - Determine the ~~best~~ attribute that best classifies
    the training data

    -? how to choose the best attribute
        (or)
    -? how does the tree decision where to split

# Different attribute selection measures

- Information Gain
- Gain Ratio
- Gini Index
- chi-square test

## Information Gain

- decrease in entropy after a dataset is split based on an attribute.

so, constructing a decision tree is all about finding the attribute that returns the highest information gain.

Entropy – metric that measures impurity in a given dataset

$$Entropy(s) = - \sum_i P(i) \log_2 P(i)$$

case 1

| $n$ | yes | 5 |
|---|---|---|
| | no | 5 |

Entropy $= -P(yes) \log_2 P(yes)$
$\qquad - P(no) \log_2 P(no)$

$= -0.5 \log_2(0.5) - 0.5 \log_2(0.5)$

$= 0.5 + 0.5 = 1$

$$\boxed{\begin{array}{l} \log_2(0.5) = \log_2(\frac{1}{2}) \\ = \log_2(2^{-1}) = -1 \end{array}}$$

case 2

| $n$ | yes | 10 |
|---|---|---|
| | no | 0 |

Entropy $= -P(yes) \log_2 P(yes)$
$\qquad - P(no) \log_2 P(no)$

$= -1 \log_2(1) - 0 \log_2 0$

$= -1 \times 0 - 0 \times 1 = 0$

# Information Gain

- measures reduction in entropy
- decides which attribute must be selected as decision node.

Information Gain = Entropy (s) - [(weighted Avg) × (Entropy (each feature))]

For the provided dataset (s)

Entropy (s) = $-P(Yes) \log_2 P(Yes) - P(No) \log_2 P(No)$

$$= -\left(\frac{9}{14}\right) \log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right)$$

$$= 0.41 + 0.53 = 0.94$$

which node to select as root node

outlook / temperature / humidity / windy?

## Outlook?



Outlook

sunny

| yes |
| yes |
| no |
| no |
| no |

overcast

| yes |
| yes |
| yes |
| yes |

rainy

| yes |
| yes |
| yes |
| no |
| no |

E (outlook = sunny) = $-\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.971$

E (outlook = overcast) = $-1 \log_2 1 \rightarrow 0 \log_2 0 = 0$
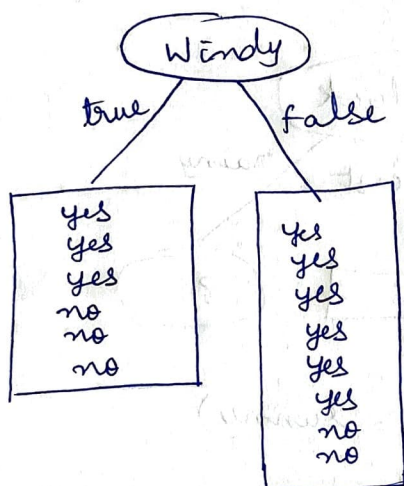
E (outlook = rainy) = $-\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.971$

Information from outlook

$$I(outlook) = \frac{5}{14} \times 0.971 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.971 = 0.693$$

⇒ Gain (outlook) = E(s) − I(outlook)

$$= 0.94 - 0.693$$

$$= 0.247$$

Windy?



E (windy = true) = 1

E (windy = false) = 0.811

Information from windy

$$I(windy) = \underbrace{\frac{8}{14} \times 0.811}_{false} + \underbrace{\frac{6}{14} \times 1}_{true} = 0.892$$

Gain (windy) = E(s) − I(windy)

$$= 0.94 - 0.892$$

$$= 0.048$$

illy Temperature

Info I(Temperature) = 0.911

Gain (Temperature) = E(s) − I(Temperature)

$$= 0.94 - 0.911$$

## Humidity

$I(\text{Humidity}) = 0.788$

$\text{Gain}(\text{Humidity}) = \overset{0.94 - 0.788}{\cancel{0.788 - 0}}$

$\qquad\qquad = 0.152$

## Information Gain

outlook $\rightarrow 0.247$ ✓

Temperature $\rightarrow 0.029$
Humidity $\rightarrow 0.152$
Windy $\rightarrow 0.048$



data at
__Left subtree__ (Outlook = sunny)

| Temperature | Humidity | Windy | Play |
|---|---|---|---|
| hot | high | false | no |
| hot | high | true | no |
| mild | high | false | no |
| cool | normal | false | yes |
| mild | normal | true | yes |

repeat the process

Data at right subtree
(Outlook = rainy)

| Temperature | Humidity | Windy | Play |
|---|---|---|---|
| mild | high | false | yes |
| cool | normal | false | yes |
| cool | normal | true | no |
| mild | normal | false | yes |
| mild | high | true | no |

repeat the process

## complete decision tree



Advantages
- simple to build
- easy to understand the solution (interpretable)

Disadvantages
- susceptible to overfitting
  - soln? pruning

Pruning
- reducing the complexity
- improving the generality

Pruning — Pre pruning (decide whether or not to split a particular node during model building)

└ Post pruning (build the decision tree, then prune the branches to avoid overfitting)