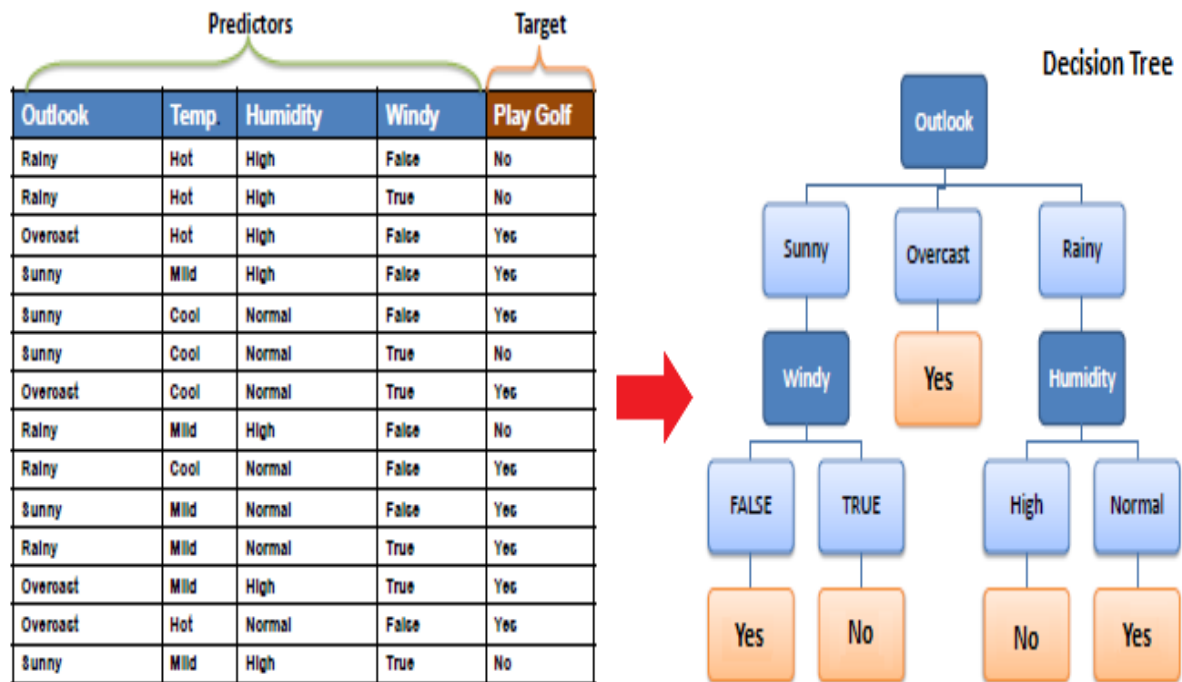Decision Tree - Classification

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy). Leaf node (e.g., Play) represents a classification or decision. The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.
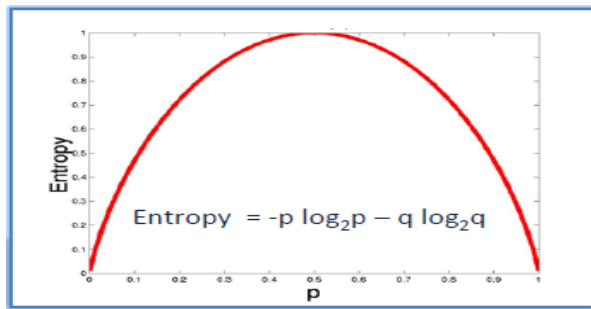


Algorithm

The core algorithm for building decision trees called ID3 by J. R. Quinlan which employs a top-down, greedy search through the space of possible branches with no backtracking. ID3 uses Entropy and Information Gain to construct a decision tree. In ZeroR model there is no predictor, in OneR model we try to find the single best predictor, naive Bayesian includes all predictors using Bayes' rule and the independence assumptions between predictors but decision tree includes all predictors with the dependence assumptions between predictors.

## Entropy

A decision tree is built top-down from a root node and involves partitioning the data into subsets that contain instances with similar values (homogenous). ID3 algorithm uses entropy to calculate the homogeneity of a sample. If the sample is completely homogeneous the entropy is zero and if the

sample is an equally divided it has entropy of one.



Entropy = -0.5 $\log_2 0.5$ – 0.5 $\log_2 0.5$ = 1

To build a decision tree, we need to calculate two types of entropy using frequency tables as follows:

a) Entropy using the frequency table of one attribute:

$$E(S) = \sum_{i=1}^{c} - p_i \log_2 p_i$$

| Play Golf | |
|---|---|
| Yes | No |
| 9 | 5 |

Entropy(PlayGolf) = Entropy (5,9)
= Entropy (0.36, 0.64)
= - (0.36 $\log_2$ 0.36) - (0.64 $\log_2$ 0.64)
= 0.94

b) Entropy using the frequency table of two attributes:

$$E(T,X) = \sum_{c \in X} P(c)E(c)$$

| | | Play Golf | | |
|---|---|---|---|---|
| | | Yes | No | |
| | Sunny | 3 | 2 | 5 |
| Outlook | Overcast | 4 | 0 | 4 |
| | Rainy | 2 | 3 | 5 |
| | | | | 14 |

E(PlayGolf, Outlook) = P(Sunny)*E(3,2) + P(Overcast)*E(4,0) + P(Rainy)*E(2,3)

= (5/14)*0.971 + (4/14)*0.0 + (5/14)*0.971

= 0.693

| Play Golf | | | | | |
|---|---|---|---|---|---|
| | | **Yes** | **No** | **Total** | |
| Humidity | **High** | 3 | 4 | 7 | Entropy (4,4) |
| | **Normal** | 6 | 1 | 7 | Entropy (1,6) |
| | | **Total Probability** | | **14** | |

Entropy with Respect to **High**

$$Entropy (4,4) = - 4*m.log(4,2)-4*m.log(4,2)$$

Entropy with Respect to **Normal**

$$Entropy (1,6) = -1*m.log(1,2) – 6*m.log(6,2)$$

P(Play Golf, Humidity) = P(High)*Entropy(4,4)   +  P(Normal)*Entropy(1,6)
                = P(7/14)*Entropy(4,4)   +  P(7/14)*Entropy(1,6)
                =

| Play Golf | | | | | |
|---|---|---|---|---|---|
| | | **Yes** | **No** | **Total** | |
| Windy | **False** | 6 | 2 | 8 | Entropy (2,6) |
| | **True** | 3 | 3 | 6 | Entropy (3,3) |
| | | **Total Probability** | | **14** | |

| Play Golf | | | | | |
|---|---|---|---|---|---|
| | | **Yes** | **No** | **Total** | |
| TEMP | **HOT** | 2 | 2 | 4 | Entropy (2,2) |
| | **MILD** | 4 | 2 | 6 | Entropy (4,2) |
| | **COOL** | 3 | 1 | 4 | Entropy (1,3) |

| | | |
|---|---|---|
| **Total Probability** | **14** | |

Entropy with Respect to **Mild**

$$Entropy (4,2) = - 4*m.log(4,2)-2*m.log(2,2)$$

Entropy with Respect to **Cool**

$$Entropy (1,3) = -1*m.log(1,2) – 3*m.log(3,2)$$

P(Play Golf, TEMP) = P(Hot)*Entropy(2,2)   + P(MILD)*Entropy(4,2)+P(Cool)*Entropy(1,3)
=P(4/14)* Entropy(2,2)   + P(6/14)*Entropy(4,2)+P(4/14)*Entropy(1,3)
= 0.2857 * 1 + 0.4285 * 0.918+0.2857
=0.96

# Information Gain

The information gain is based on the decrease in entropy after a dataset is split on an attribute. Constructing a decision tree is all about finding attribute that returns the highest information gain (i.e., the most homogeneous branches).

Step 1: Calculate entropy of the target.

$$\textbf{Entropy(PlayGolf)} = Entropy\ (5,9)$$
$$= Entropy\ (0.36, 0.64)$$
$$= - (0.36\ log_2\ 0.36) - (0.64\ log_2\ 0.64)$$
$$= 0.94$$

Step 2: The dataset is then split on the different attributes. The entropy for each branch is calculated. Then it is added proportionally, to get total entropy for the split. The resulting entropy is subtracted from the entropy before the split. The result is the Information Gain, or decrease in entropy.

| Outlook | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Outlook | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| Gain = 0.247 | | | |

| Temp. | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Temp. | Hot | 2 | 2 |
| | Mild | 4 | 2 |
| | Cool | 3 | 1 |
| Gain = 0.029 | | | |

| Humidity | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Humidity | High | 3 | 4 |
| | Normal | 6 | 1 |
| Gain = 0.152 | | | |

| Windy | | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Windy | False | 6 | 2 |
| | True | 3 | 3 |
| Gain = 0.048 | | | |

$$Gain(T,X) = Entropy(T) - Entropy(T,X)$$

G(PlayGolf, Outlook) = E(PlayGolf) – E(PlayGolf, Outlook)
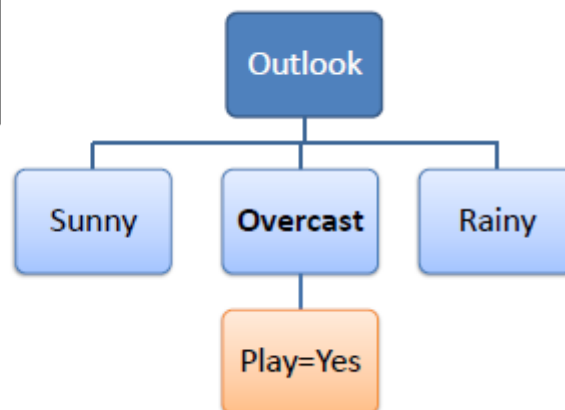
= 0.940 – 0.693 = 0.247

Step 3: Choose attribute with the largest information gain as the decision node, divide the dataset by its branches and repeat the same process on every branch.

| Outlook | ★ | Play Golf | |
|---|---|---|---|
| | | Yes | No |
| Outlook | Sunny | 3 | 2 |
| | Overcast | 4 | 0 |
| | Rainy | 2 | 3 |
| Gain = 0.247 | | | |

| Outlook | Temp | Humidity | Windy | Play Golf |
|---|---|---|---|---|
| Sunny | Mild | High | FALSE | Yes |
| Sunny | Cool | Normal | FALSE | Yes |
| Sunny | Cool | Normal | TRUE | No |
| Sunny | Mild | Normal | FALSE | Yes |
| Sunny | Mild | High | TRUE | No |

| Overcast | Hot | High | FALSE | Yes |
|---|---|---|---|---|
| Overcast | Cool | Normal | TRUE | Yes |
| Overcast | Mild | High | TRUE | Yes |
| Overcast | Hot | Normal | FALSE | Yes |

| Rainy | Hot | High | FALSE | No |
|---|---|---|---|---|
| Rainy | Hot | High | TRUE | No |
| Rainy | Mild | High | FALSE | No |
| Rainy | Cool | Normal | FALSE | Yes |
| Rainy | Mild | Normal | TRUE | Yes |

Step 4a: A branch with entropy of 0 is a leaf node.

| Temp | Humidity | Windy | Play Golf |
|---|---|---|---|
| Hot | High | FALSE | Yes |
| Cool | Normal | TRUE | Yes |
| Mild | High | TRUE | Yes |
| Hot | Normal | FALSE | Yes |

Step 4b: A branch with entropy more than 0 needs further splitting.

| Temp | Humidity | Windy | Play Golf |
|------|----------|-------|-----------|
| Mild | High | FALSE | Yes |
| Cool | Normal | FALSE | Yes |
| Mild | Normal | FALSE | Yes |
| Cool | Normal | TRUE | No |
| Mild | High | TRUE | No |



Step 5: The ID3 algorithm is run recursively on the non-leaf branches, until all data is classified.

Decision Tree to Decision Rules

A decision tree can easily be transformed to a set of rules by mapping from the root node to the leaf nodes one by one.
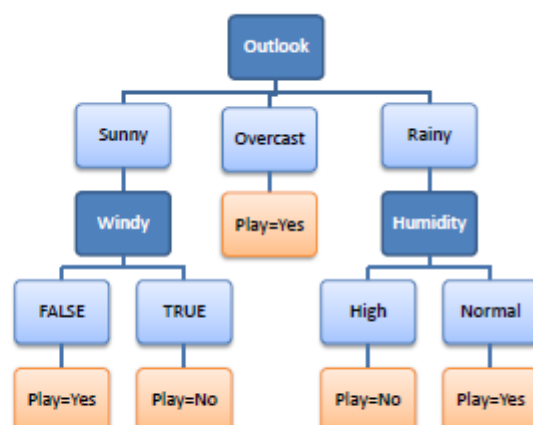
R₁: IF (Outlook=Sunny) AND (Windy=FALSE) THEN Play=Yes

R₂: IF (Outlook=Sunny) AND (Windy=TRUE) THEN Play=No

R₃: IF (Outlook=Overcast) THEN Play=Yes

R₄: IF (Outlook=Rainy) AND (Humidity=High) THEN Play=No

R₅: IF (Outlook=Rain) AND (Humidity=Normal) THEN Play=Yes



## Decision Trees - Issues

Working with continuous attributes (binning)

Avoiding overfitting

Super Attributes (attributes with many unique values)

Working with missing values