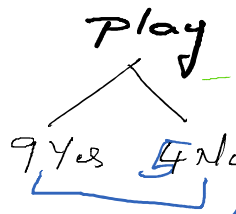


Weather	Humidity	Wind	play
sunny	high	weak	no
sunny	high	strong	no
cloudy	high	weak	yes
rainy	high	weak	yes
rainy	normal	weak	yes
rainy	normal	strong	no
cloudy	normal	strong	yes
sunny	high	weak	no
sunny	normal	weak	yes
rainy	normal	weak	yes
sunny	normal	strong	yes
cloudy	high	strong	yes
cloudy	normal	weak	yes
rainy	high	strong	no



$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

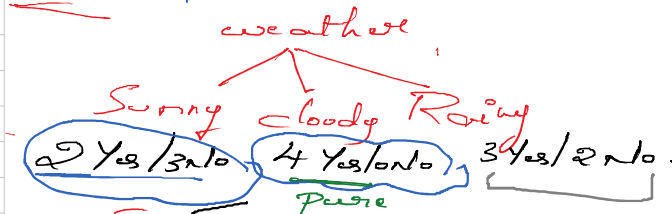
$$H(S) = -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{No}) \log_2 P(\text{No})$$

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

$$P(\text{Yes}) = \frac{9}{14}$$

$$P(\text{No}) = \frac{5}{14}$$

$$H(S) = 0.97$$



$$H(\text{Sunny}) = -\left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) - \left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) = 0.97$$

$$H(\text{Cloudy}) = -\left(\frac{4}{4}\right) \log_2 \left(\frac{4}{4}\right) - 0 = 0$$

$$H(\text{Rainy}) = -\left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right) = 0.97$$

$$\text{Gain} = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

$$= 0.97 - \left(\frac{5}{14}\right) H(\text{Sunny}) - \left(\frac{4}{14}\right) H(\text{Cloudy}) - \left(\frac{5}{14}\right) H(\text{Rainy})$$

$$= 0.97 - \left(\frac{5}{14}\right)(0.97) - \left(\frac{4}{14}\right)(0) - \left(\frac{5}{14}\right)(0.97)$$

$$\text{Gain} = 0.247$$

weather
Gain

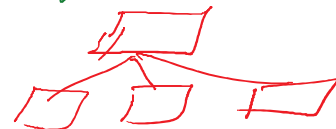
$$\text{Gain}(\text{weather}) = 0.247$$

$$\text{Gain}(\text{wind}) = 0.12$$

$$\text{Gain}(\text{Humidity}) = 0.07$$

out of all these, weather feature is giving maximum information gain.

weather is the root node



Limitations with Entropy

Computation time taken will be very high because of log function.

Gini Index / Gini Impurity 1. 1. Entropy

Gini Index / Gini Impurity

Since computation time is high in Entropy we don't use entropy instead we use Gini.

Gini

- * It measures impurity of the split.
- * It is used to decide the root node

Steps

- Find Gini of each category in column
- Find weighted Gini of each column.
- Repeat above steps for all the columns
- Select the column which has less weighted Gini.

$$Gini = 1 - \sum P_i^2 = 1 - P_+^2 - P_-^2$$

P_i → probability of classes

$$\text{weighted Gini} = \sum \frac{S_v}{S} Gini(S)$$

Note → Select the variable which has less weighted Gini as a root node.

weighted Gini (weather) = 0.43
 weighted Gini (windy) = 0.8
 weighted Gini (Humidity) = 0.5

root node

wind

Weather	Humidity	Wind	play
sunny	high	weak	no
sunny	high	strong	no
cloudy	high	weak	yes
rainy	high	weak	yes
rainy	normal	weak	yes
rainy	normal	strong	no
cloudy	normal	strong	yes
sunny	high	weak	no
sunny	normal	weak	yes
rainy	normal	weak	yes
sunny	normal	strong	yes
cloudy	high	strong	yes
cloudy	normal	weak	yes
rainy	high	strong	no

weak strong
 6 Yes, 2 No 3 Yes, 3 No

$$Gini = 1 - [P(Yes)]^2 - [P(No)]^2$$

$$Gini(weak) = 1 - \left(\frac{6}{8}\right)^2 - \left(\frac{2}{8}\right)^2 = 0.375$$

$$Gini(strong) = 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 = 0.5$$

$$\text{weighted Gini} = \sum \frac{|S_v|}{|S|} Gini(S_v)$$

S → Sample before split

$$\text{weighted Gini} = \left(\frac{S_v}{S}\right) Gini(weak) + \left(\frac{S_v}{S}\right) Gini(strong)$$

$S \rightarrow$ Sample before split
 $S_v \rightarrow$ Sample after split
 $\text{weighted } Gini = \left(\frac{S_v}{S}\right) Gini(\text{weak}) + \left(\frac{S_v}{S}\right) Gini(\text{strong})$
 $\text{weighted } Gini = \left(\frac{8}{14}\right)(0.375) + \frac{6}{14}(0.5) = 0.43$

$S = 14$
 $\text{wind}[14]$
 weak 8
 $S_v = 8$
 strong 6
 $S_v = 6$

Age
 20 - Yes
 25 - Yes
 30 - No
 35 - Yes
 40 - No
 45 - No

Age
 < 30 2Y, 0No
 ≥ 30 1Y, 3No

Temp
 * Scaling and handling outliers are not required for Decision Tree algorithms.

Overfitting

- Training score is high and testing score is less.
- Model works extremely well on train data and score will be high, but during testing it will make wrong prediction eventually testing score will be less.
- Low bias and high variance → overfitting

Bias → Training error
 Variance → Testing error

Underfitting

- * Training score will be less and testing score is less.
- * High bias, high variance

Score Variance.

* high bias, high variance

How do you say that model is Generalized

→ Training score and testing score should be high.

→ Low bias, low variance

→ high bias, low variance → Training score can be low and testing score should be high.

How to overcome overfitting?

→ Hyperparameters