

## Learning Decision Trees

What is the best model in the form of a decision tree?

Top-Down Induction of Decision Tree Algorithm on learning set L:

IF all samples from L belong to the same class C:

Make a terminal node (i.e. a leaf) and label it C;

ELSE:

Pick the most informative attribute A.

Split L with respect to the values of A.

Recursively run this algorithm on all subsets from the previous step.

How can we measure the information score of attribute A?

### Information Gain

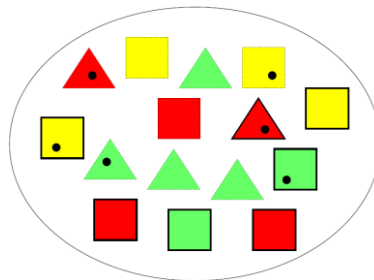
$$IG(Y|X) = H(Y) - H(Y|X)$$

$$H(Y|X) = \sum_v p(X=v) * H(Y | X=v)$$

$$H(Y | X=v) = - \sum_r p(Y=r | X=v) * \log_2(Y=r | X=v)$$

---

## Shapes



Given is a set of 15 shapes, described by 3 features: Color, Dot, and Edge.

The task is to learn the concept of a shape based on these features.

The class (or outcome) is a binary variable with 2 values: ▲, ■

		color			
shape		red	yellow	green	
	▲	2	0	4	6
	■	3	4	2	9
		5	4	6	15

$$H(\text{shape}) = -6/15 \log_2(6/15) - 9/15 \log_2(9/15) = 0.971$$

$$IG(\text{color}) = H(\text{shape}) - Ires(\text{color}) = 0.971 - 0.690 = 0.281$$

$$\begin{aligned} Ires(\text{color}) &= 5/15 H(\text{red}) + 4/15 H(\text{yellow}) + 6/15 H(\text{green}) = \\ &= 1/3 * 0.971 + 4/15 * 0 + 6/15 * 0.918 = 0.690 \end{aligned}$$

$$H(\text{red}) = -2/5 \log_2(2/5) - 3/5 \log_2(3/5) = 0.971$$

$$H(\text{yellow}) = 0 \log_2(0) - 4/4 \log_2(4/4) = 0$$

$$H(\text{green}) = -4/6 \log_2(4/6) - 2/6 \log_2(2/6) = 0.918$$

$$H(\text{color}) = -5/15 \log_2(5/15) - 4/15 \log_2(4/15) - 6/15 \log_2(6/15) = 1.565$$

$$RIG(\text{color}) = IG(\text{color})/H(\text{color}) = 0.281/1.565 = 0.179$$

		Dot		
Shape		y	n	
	▲	3	3	6
	■	3	6	9
		6	9	15

$$H(\text{Shape}) = -6/15 \log_2(6/15) - 9/15 \log_2(9/15) = 0.971$$

$$IG(\text{Dot}) = H(\text{Shape}) - Ires(\text{Dot}) = 0.971 - 0.951 = 0.02$$

$$\begin{aligned} Ires(\text{Dot}) &= 6/15 H(y) + 9/15 H(n) = \\ &= 2/5 * 1 + 3/5 * 0.918 = 0.951 \end{aligned}$$

$$H(y) = -2 * 1/2 \log_2(1/2) = 1$$

$$H(n) = -1/3 \log_2(1/3) - 2/3 \log_2(2/3) = 0.918$$

$$H(\text{Dot}) = -6/15 \log_2(6/15) - 9/15 \log_2(9/15) = 0.971$$

$$RIG(\text{Dot}) = IG(\text{Dot})/H(\text{Dot}) = 0.02/0.971 = 0.021$$

		Edge		
Shape		y	n	
	▲	1	5	6
	■	6	3	9
		7	8	15

$$H(\text{Shape}) = -6/15 \log_2(6/15) - 9/15 \log_2(9/15) = 0.971$$

$$IG(\text{Edge}) = H(\text{Shape}) - Ires(\text{Edge}) = 0.971 - 0.784 = 0.187$$

$$\begin{aligned} Ires(\text{Edge}) &= 7/15 H(y) + 8/15 H(n) = \\ &= 7/15 * 0.591 + 8/15 * 0.954 = 0.784 \end{aligned}$$

$$H(y) = -1/7 \log_2(1/7) - 6/7 \log_2(6/7) = 0.591$$

$$H(n) = -5/8 \log_2(5/8) - 3/8 \log_2(3/8) = 0.954$$

$$H(\text{Edge}) = -7/15 \log_2(7/15) - 8/15 \log_2(8/15) = 0.996$$

$$RIG(\text{Edge}) = IG(\text{Edge})/H(\text{Edge}) = 0.187/0.996 = 0.188$$