

Milan Zanussi

CSCI 4350/5350

Homework 7

Due: Thu. Nov. 8, 11:00 PM

1. (1 point) How many bits are required to encode an overall class membership of 12.5%? [hint: $I(x) = -\log_2(P(x))$]

$$12.5\% = \frac{1}{8}. \text{ So } I(x) = I\left(\frac{1}{8}\right) = -\log_2\left(\frac{1}{8}\right) = \log_2(2^3) = \boxed{3}$$

2. (1 point) What is the Shannon entropy of an overall class membership of 12.5%? [hint: $H(x) = P(x) * I(x)$]

$$I(x) = 3, P(x) = \frac{1}{8}. \text{ So } H(x) = P(x)I(x) = \boxed{\frac{3}{8}}$$

3. (2 points) Given two classes, + and -, what is the total information for an overall class membership of 12.5% and 87.5%, respectively? [hint: $I(x,y) = H(x)+H(y)$]

$$P(+)=\frac{1}{8}, P(-)=\frac{7}{8}, I(+)=3, I(-)=-\log_2\left(\frac{7}{8}\right)=\log_2\left(\frac{8}{7}\right)=3-\log_2(7). \text{ So } I(+,-)=\frac{1}{8}(3)+\frac{7}{8}(3-\log_2(7))=\boxed{3-\frac{7}{8}\log_2(7)}$$

Use the data set below to answer all remaining questions:

size	shape	color	class
^^^^	^^^^^	^^^^^	^^^^^
big	triangle	red	+
big	square	green	-
big	square	red	-
small	circle	green	-
big	triangle	green	+
small	square	red	+
small	square	green	+
small	circle	red	-

$$\begin{aligned} P(+) &= 0.5 = \frac{1}{2}, P(-) = \frac{1}{2} \\ P(\text{big}) &= \frac{1}{2}, P(\text{small}) = \frac{1}{2} \\ P(+|\text{big}) &= \frac{1}{2}, P(-|\text{big}) = \frac{1}{2} \\ P(-|\text{small}) &= \frac{1}{2}, P(+|\text{small}) = \frac{1}{2} \end{aligned}$$

4. (2 points) What is $I(+,-)$?

$$I(+,-) = H(+)+H(-) = P(+)I(+)+P(-)I(-) = \frac{1}{2}(1) + \frac{1}{2}(1) = \boxed{1}$$

5. (2 points) What is $P(\text{big})$?

$$P(\text{big}) = \frac{1}{2}$$

6. (2 points) What is the $P(+|\text{big})$?

$$P(+|\text{big}) = \frac{1}{2}$$

7. (2 points) What is $E(\text{size})$?

[hint: $E(\text{size}) = -P(\text{big}) * [P(+|\text{big}) * \log_2(P(+|\text{big})) + P(-|\text{big}) * \log_2(P(-|\text{big}))] - P(\text{small}) * [P(+|\text{small}) * \log_2(P(+|\text{small})) + P(-|\text{small}) * \log_2(P(-|\text{small}))]$]

$$\begin{aligned} E(\text{size}) &= -\frac{1}{2} \left[\frac{1}{2} (\log_2(\frac{1}{2})) + \frac{1}{2} (\log_2(\frac{1}{2})) \right] - \frac{1}{2} \left[\frac{1}{2} (\log_2(\frac{1}{2})) + \frac{1}{2} (\log_2(\frac{1}{2})) \right] \\ &= -\frac{1}{2} \log_2(\frac{1}{2}) - \frac{1}{2} \log_2(\frac{1}{2}) = -\log_2(\frac{1}{2}) = \boxed{1} \end{aligned}$$

8. (2 points) What is $\text{gain}(\text{size})$? [hint: $\text{gain}(\text{size}) = I(+, -) - E(\text{size})$]

$$\text{gain}(\text{size}) = I(+, -) - E(\text{size}) = 1 - 1 = 0$$

9. (2 points) Calculate $\text{gain}(\text{shape})$ and $\text{gain}(\text{color})$.

$$\text{gain}(\text{shape}) = I(+, -) - E(\text{shape}) = \frac{1}{2} \quad \text{gain}(\text{color}) = I(+, -) - E(\text{color}) = 0$$

10. (1 point) Which is the best attribute to select for making the first split? (max gain)

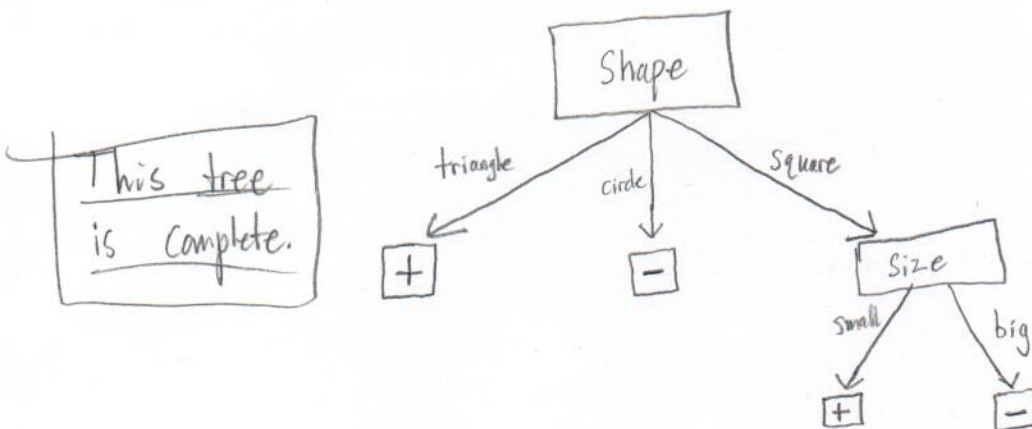
Shape

11. (3 points) Draw a complete ID3 decision tree for the data (show your work for how you decided to make each split in the tree)

$$\begin{aligned} &P(\text{triangle}) = \frac{1}{4} \quad P(+|\text{triangle}) = 1 \quad P(-|\text{triangle}) = 0 \\ &P(\text{circle}) = \frac{1}{4} \quad P(+|\text{circle}) = 0 \quad P(-|\text{circle}) = 1 \\ &P(\text{square}) = \frac{1}{2} \quad P(+|\text{square}) = \frac{1}{2} \quad P(-|\text{square}) = \frac{1}{2} \\ &P(\text{red}) = \frac{1}{2} \quad P(+|\text{red}) = \frac{1}{2} \quad P(-|\text{red}) = \frac{1}{2} \\ &P(\text{green}) = \frac{1}{2} \quad P(+|\text{green}) = \frac{1}{2} \quad P(-|\text{green}) = \frac{1}{2} \end{aligned}$$

$$E(\text{shape}) = -\frac{1}{4} [1 \log_2(1) + 0 \log_2(0)] - \frac{1}{4} [0 \log_2(0) + 1 \log_2(1)] - \frac{1}{2} [\frac{1}{2} \log_2(\frac{1}{2}) + \frac{1}{2} \log_2(\frac{1}{2})] = -\frac{1}{2} [\log_2(\frac{1}{2})] = \frac{1}{2} \log_2(2) = \frac{1}{2}$$

$$E(\text{color}) = 1 \quad (\text{since all values are the same as the size values})$$



Restricting Sample to Squares:

big	square	green	-
big	square	red	-
small	square	red	+
small	square	green	+

$$P(+)=\frac{1}{2} \quad P(-)=\frac{1}{2} \quad I(+, -) = 1$$

$$E(\text{size}|\text{shape}) = -\frac{1}{2} [1 \log_2(1) + 0 \log_2(0)] - \frac{1}{2} [0 \log_2(0) + 1 \log_2(1)] = 0$$

$$E(\text{color}|\text{shape}) = -\frac{1}{2} [\frac{1}{2} \log_2(\frac{1}{2}) + \frac{1}{2} \log_2(\frac{1}{2})] - \frac{1}{2} [\frac{1}{2} \log_2(\frac{1}{2}) + \frac{1}{2} \log_2(\frac{1}{2})] = 1$$

$$\begin{aligned} &P(\text{big}) = \frac{1}{2} \quad P(+|\text{big}) = 0 \quad P(-|\text{big}) = 1 \\ &P(\text{small}) = \frac{1}{2} \quad P(+|\text{small}) = 1 \quad P(-|\text{small}) = 0 \\ &P(\text{red}) = \frac{1}{2} \quad P(+|\text{red}) = \frac{1}{2} \quad P(-|\text{red}) = \frac{1}{2} \\ &P(\text{green}) = \frac{1}{2} \quad P(+|\text{green}) = \frac{1}{2} \quad P(-|\text{green}) = \frac{1}{2} \end{aligned}$$

$$\text{gain}(\text{size}|\text{shape}) = I(+, -|\text{shape}) - E(\text{size}|\text{shape}) = 1$$

$$\text{gain}(\text{color}|\text{shape}) = I(+, -|\text{shape}) - E(\text{color}|\text{shape}) = 0$$

Size is chosen to branch from square case of shape