

# CHOSEN DATA MINING PROBLEMS E.G. DECISION TREES

dr Piotr Wąsiewicz

1. From the training set shown below in the table with the help of top-down decision tree induction algorithm create a decision tree (use the entropy cost). The attribute age should be discretized using two thresholds 30 and 65 years. The attribute risk will be the label class.

$x$	age	car	risk
1	18	mini	big
2	35	mini	small
3	50	racer	big
4	66	van	big
5	18	racer	big
6	35	van	small
7	60	mini	small
8	70	racer	big
9	25	van	small

SOLUTION:

The attribute **age** has three values after discretization:

$w_1$ : **age** < 30,  $w_2$ : **age**  $\geq$  30  $\wedge$  **age** < 65,  $w_3$ : **age**  $\geq$  65.

First the information of the whole set  $I(P)$  and attributes is calculated using the entropy impurity measures.

$$I(P) = -\frac{|P^{small}|}{|P|} \log_2\left(\frac{|P^{small}|}{|P|}\right) - \frac{|P^{big}|}{|P|} \log_2\left(\frac{|P^{big}|}{|P|}\right) = -\frac{4}{9} \log_2\left(\frac{4}{9}\right) - \frac{5}{9} \log_2\left(\frac{5}{9}\right) = 0.991,$$

$$E_{age,w_1}(P) = -\frac{|P_{age,w_1}^{small}|}{|P_{age,w_1}|} \log_2\left(\frac{|P_{age,w_1}^{small}|}{|P_{age,w_1}|}\right) - \frac{|P_{age,w_1}^{big}|}{|P_{age,w_1}|} \log_2\left(\frac{|P_{age,w_1}^{big}|}{|P_{age,w_1}|}\right) = -\frac{1}{3} \log_2\left(\frac{1}{3}\right) - \frac{2}{3} \log_2\left(\frac{2}{3}\right) = 0.918,$$

$$E_{age,w_2}(P) = -\frac{|P_{age,w_2}^{small}|}{|P_{age,w_2}|} \log_2\left(\frac{|P_{age,w_2}^{small}|}{|P_{age,w_2}|}\right) - \frac{|P_{age,w_2}^{big}|}{|P_{age,w_2}|} \log_2\left(\frac{|P_{age,w_2}^{big}|}{|P_{age,w_2}|}\right) = -\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) = 0.811,$$

$$E_{age,w_3}(P) = -\frac{|P_{age,w_3}^{small}|}{|P_{age,w_3}|} \log_2\left(\frac{|P_{age,w_3}^{small}|}{|P_{age,w_3}|}\right) - \frac{|P_{age,w_3}^{big}|}{|P_{age,w_3}|} \log_2\left(\frac{|P_{age,w_3}^{big}|}{|P_{age,w_3}|}\right) = -\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{2}{2}\right) = 0,$$

$$E_{car,mini}(P) = -\frac{|P_{car,mini}^{small}|}{|P_{car,mini}|} \log_2\left(\frac{|P_{car,mini}^{small}|}{|P_{car,mini}|}\right) - \frac{|P_{car,mini}^{big}|}{|P_{car,mini}|} \log_2\left(\frac{|P_{car,mini}^{big}|}{|P_{car,mini}|}\right) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.918,$$

$$E_{car,van}(P) = -\frac{|P_{car,van}^{small}|}{|P_{car,van}|} \log_2\left(\frac{|P_{car,van}^{small}|}{|P_{car,van}|}\right) - \frac{|P_{car,van}^{big}|}{|P_{car,van}|} \log_2\left(\frac{|P_{car,van}^{big}|}{|P_{car,van}|}\right) = -\frac{2}{3} \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2\left(\frac{1}{3}\right) = 0.918,$$

$$E_{car,racer}(P) = -\frac{|P_{car,racer}^{small}|}{|P_{car,racer}|} \log_2\left(\frac{|P_{car,racer}^{small}|}{|P_{car,racer}|}\right) - \frac{|P_{car,racer}^{big}|}{|P_{car,racer}|} \log_2\left(\frac{|P_{car,racer}^{big}|}{|P_{car,racer}|}\right) = -\frac{0}{3} \log_2\left(\frac{0}{3}\right) - \frac{3}{3} \log_2\left(\frac{3}{3}\right) = 0,$$

Next the weighted entropies are obtained:

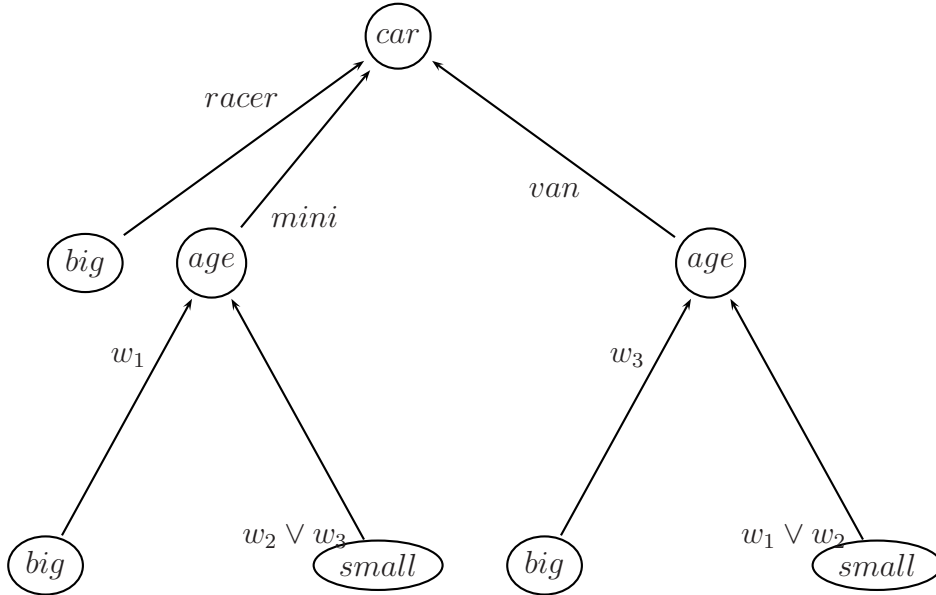
$$\begin{aligned}
E_{\text{age}}(P) &= \frac{|P_{\text{age},w_1}|}{|P|} E_{\text{age},w_1}(P) + \frac{|P_{\text{age},w_2}|}{|P|} E_{\text{age},w_2}(P) + \frac{|P_{\text{age},w_3}|}{|P|} E_{\text{age},w_3}(P) = \frac{3}{9}(0.918) + \\
&\frac{4}{9}(0.811) + \frac{2}{9}0 = 0,666, \\
E_{\text{car}}(P) &= \frac{|P_{\text{car},\text{mini}}|}{|P|} E_{\text{car},\text{mini}}(P) + \frac{|P_{\text{car},\text{van}}|}{|P|} E_{\text{car},\text{van}}(P) + \frac{|P_{\text{car},\text{racer}}|}{|P|} E_{\text{car},\text{racer}}(P) = \frac{3}{9}(0.918) + \\
&\frac{3}{9}(0.918) + \frac{3}{9}0 = 0,612,
\end{aligned}$$

And all attributes information measures:

$$\begin{aligned}
IV_{\text{age}}(P) &= -\frac{|P_{\text{age},w_1}|}{|P|} \log_2\left(\frac{|P_{\text{age},w_1}|}{|P|}\right) - \frac{|P_{\text{age},w_2}|}{|P|} \log_2\left(\frac{|P_{\text{age},w_2}|}{|P|}\right) - \frac{|P_{\text{age},w_3}|}{|P|} \log_2\left(\frac{|P_{\text{age},w_3}|}{|P|}\right) = \\
&-\frac{3}{9} \log_2\left(\frac{3}{9}\right) - \frac{4}{9} \log_2\left(\frac{4}{9}\right) - \frac{2}{9} \log_2\left(\frac{2}{9}\right) = 0,528 + 0,519 + 0,482 = 1,53, \\
IV_{\text{car}}(P) &= -\frac{|P_{\text{car},\text{mini}}|}{|P|} \log_2\left(\frac{|P_{\text{car},\text{mini}}|}{|P|}\right) - \\
&\frac{|P_{\text{car},\text{van}}|}{|P|} \log_2\left(\frac{|P_{\text{car},\text{van}}|}{|P|}\right) - \frac{|P_{\text{car},\text{racer}}|}{|P|} \log_2\left(\frac{|P_{\text{car},\text{racer}}|}{|P|}\right) = \\
&-\frac{3}{9} \log_2\left(\frac{3}{9}\right) - \frac{3}{9} \log_2\left(\frac{3}{9}\right) - \frac{3}{9} \log_2\left(\frac{3}{9}\right) = 0,528 + 0,528 + 0,528 = 1,584,
\end{aligned}$$

At last information increase coefficients:

$$\begin{aligned}
\vartheta_{\text{age}}(P) &= \frac{I(P) - E_{\text{age}}(P)}{IV_{\text{age}}(P)} = \frac{0,991 - 0,666}{1,53} = 0,212 \\
\vartheta_{\text{car}}(P) &= \frac{I(P) - E_{\text{car}}(P)}{IV_{\text{car}}(P)} = \frac{0,991 - 0,612}{1,584} = 0,239
\end{aligned}$$



The attribute **car** has the biggest coefficient and wins to be the first node.  
For its value **racer** every example with this value has the label **big** of the target risk.