

- ① First we have to find the best split for the data using entropy based information gain

The dataset S consist of M instances. 6 of them are positive and 5 negative. This leads us to:

$$I(S) = -p_{\text{pos}} \log_2(p_{\text{pos}}) - p_{\text{neg}} \log_2(p_{\text{neg}})$$

$$= -\frac{6}{11} \log_2 \left(\frac{6}{11} \right) - \frac{5}{11} \log_2 \left(\frac{5}{11} \right) \approx 0.934$$

For each of the attributes we get the following gain

Outlook: Sunny [2+, 2-]

Overcast [1+, 2-]

Rain [3+, 1-]

$$\text{Gain } I(S) - \frac{4}{11} I(S_{\text{sunny}}) - \frac{3}{11} I(S_{\text{overcast}}) - \frac{4}{11} I(S_{\text{rain}})$$

$$= 0,0849$$

Wind: Weak [3+, 2-]

Strong [3+, 3-]

$$\text{Gain } I(S) - \frac{5}{11} I(S_{\text{weak}}) - \frac{6}{11} I(S_{\text{strong}}) \approx 0,0072$$

→ Temperature: T° [1, 6] 8° [1, 5] 10° [1, 4] 11° [1, 0]

out examples: 20° [1+, 0] | 21° [0+, 1] 25° [1+, 1] | 31° [0+, 1]

32° [0+, 1] 34° [0+, 1]

$$\text{Gain } I(S) - \frac{1}{11} I(S_{1^{\circ}}) - \frac{1}{11} I(S_{8^{\circ}}) - \frac{1}{11} I(S_{10^{\circ}}) - \frac{1}{11} I(S_{11^{\circ}})$$

$$- I(S_{20^{\circ}}) - \frac{1}{11} I(S_{21^{\circ}}) - \frac{2}{11} I(S_{25^{\circ}}) - \frac{1}{11} I(S_{31^{\circ}})$$

$$- \frac{1}{11} I(S_{32^{\circ}}) - \frac{1}{11} I(S_{34^{\circ}}) \approx 0,812$$

// everything 0, except $I(S_{25^{\circ}}) = 1$

The temperature gives the best information gain. This category will be our first node. And lets next the two red lines above mark adjacent examples

with different classes, values already sorted, see above (should be sorted now, already due by accident)

Mean values for the two class-changes:

$$a) \frac{20+21}{2}$$

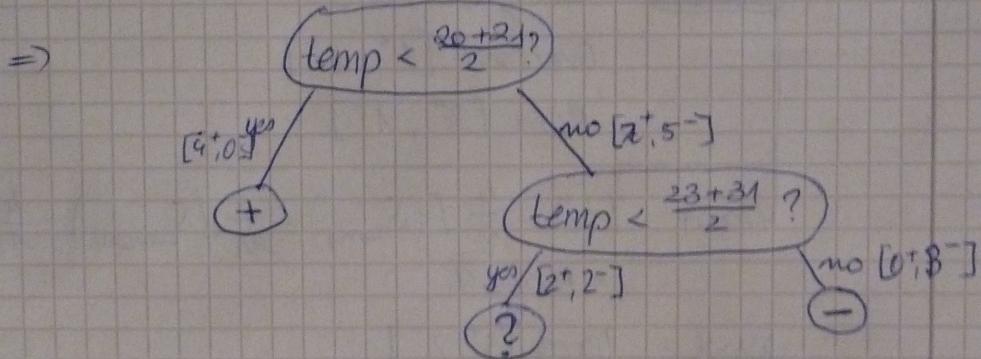
$$b) \frac{23+31}{2}$$

Calculate the information gain for both mean tests, namely $\text{temp} < \frac{20+21}{2}$ and $\text{temp} < \frac{23+31}{2}$.

$$I(S_{\text{temp} < \frac{20+21}{2}}) = -\frac{5}{11} \log\left(\frac{5}{11}\right) - \frac{6}{11} \log\left(\frac{6}{11}\right) \approx 0,994$$

$$I(S_{\text{temp} < \frac{23+31}{2}}) = -\frac{8}{11} \log\left(\frac{8}{11}\right) - \frac{3}{11} \log\left(\frac{3}{11}\right) \approx 0,8454$$

test a) gives the higher information gain.



The children hot ($\text{temp} > 27^\circ\text{C}$) and cool ($\text{temp} < \frac{20+21}{2}$) are classified perfectly, no further operations here. For the descendant with temperatures between "hot" and "cool" another split has to be done.

There are four instances with temperature belonging to this group. The instances are:

	Outlook	Wind	Injoy spad
S^* {	Rain	weak	+
	Sunny	strong	+
	Overcast	strong	-
	Rain	strong	-

let's call this set S^* .



Now we must calculate the gain for these attributes.

$I(S^*) = 1$ since there are ^{disjointed} equally divided (+) and (-) instances (2 of each)

Outlook : Rain $[1^+, 1^-]$

Sunny $[1^+, 0^-]$

Overcast $[0^+, 1^-]$

$$\text{Gain: } I(S^*) - \frac{2}{4} \underbrace{I(S_{\text{Rain}}^*)}_1 - \frac{1}{4} \underbrace{I(S_{\text{Sunny}}^*)}_0 - \frac{1}{4} \underbrace{I(S_{\text{Overcast}}^*)}_0$$

$$= 1 - \frac{1}{2} = \frac{1}{2} = 0,5 //$$

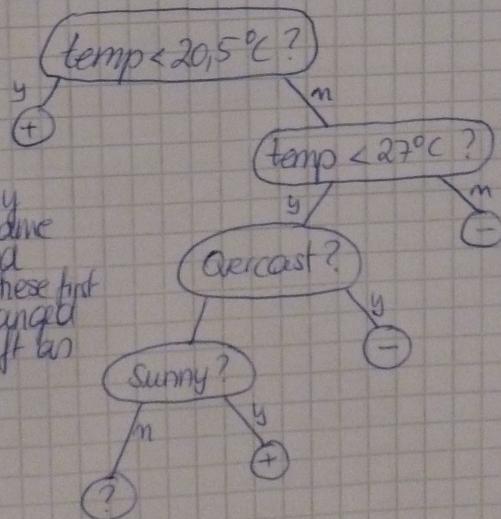
Wind: Weak $[1^+, 0^-]$

Strong $[1^+, 2^-]$

$$\text{Gain: } I(S^*) - \underbrace{\frac{1}{4} I(S_{\text{Weak}}^*)}_0 - \frac{3}{4} I(S_{\text{Strong}}^*) = 1 - \frac{3}{4} \cdot \cancel{I(S^*)} \approx 0,3113,$$

Outlook gives the largest gain and is the next attribute for the subtree

Overcast and Sunny obviously give the same information gain and Rainy is left, so these first values can be exchanged and Rainy is left as last test mode.



The descendants "Overcast" and "Sunny" are perfectly ~~split~~ classified. This leaves the attribute wind. This classifies the rest perfectly hence the following tree:

Result:

