

## UČENJE DREVES

Iščemo najboljši model v obliki drevesa. Prostor iskanja so načeloma vsa možna drevesa; velika kombinatorična zahtevnost. (U ... učna množica)

V praksi gradimo drevo T s požrešnim algoritmom:

- če vsi primeri iz U pripadajo istemu razredu R, naredi list in ga označi z R;
- sicer: izberi najbolj informativen atribut A in razdeli U na podmnožice glede na vrednosti atributa A. Rekurzivno izvedi algoritem na vsaki od pravkar generiranih podmnožic U.

Kako merimo informativnost atributa A?

### Informacijski prispevek (*Information Gain*)

iz teorije informacije:

$IG(Y|X)$  (v strojnem učenju: Y je razred, X je atribut)

Poslati moram sporočilo Y; koliko bitov v povprečju lahko prihranim, če obe strani poznata vrednost X?

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

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

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

Oblike likov:

robot učimo koncepta "oblika". Robot ima senzorje za zaznavanje barve, pike in roba (atributi). Razred je oblika: trikotnik, kvadrat

		barva			
oblika		rdeča	rumena	zelena	
	trikotnik	2	0	4	6
	kvadrat	3	4	2	9
		5	4	6	15

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

$$IG(\text{barva}) = H(\text{oblika}) - I_{\text{res}}(\text{barva}) = 0.971 - 0.690 = 0.281$$

$$I_{\text{res}}(\text{barva}) = 5/15 H(\text{rdeča}) + 4/15 H(\text{rumena}) + 6/15 H(\text{zelena}) = \\ = 1/3 * 0.971 + 4/15 * 0 + 6/15 * 0.918 = 0.690$$

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

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

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

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

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

		pika		
oblika		da	ne	
	trikotnik	3	3	6
	kvadrat	3	6	9
		6	9	15

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

$$IG(\text{pika}) = H(\text{oblika}) - I_{\text{res}}(\text{pika}) = 0.971 - 0.951 = 0.02$$

$$\begin{aligned} I_{\text{res}}(\text{pika}) &= 6/15 H(\text{da}) + 9/15 H(\text{ne}) = \\ &= 2/5 * 1 + 3/5 * 0.918 = 0.951 \end{aligned}$$

$$H(\text{da}) = -2 * 1/2 \log_2(1/2) = 1$$

$$H(\text{ne}) = -1/3 \log_2(1/3) - 2/3 \log_2(2/3) = 0.918$$

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

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

		rob		
oblika		da	ne	
	trikotnik	1	5	6
	kvadrat	6	3	9
		7	8	15

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

$$IG(\text{rob}) = H(\text{oblika}) - I_{\text{res}}(\text{rob}) = 0.971 - 0.784 = 0.187$$

$$\begin{aligned} I_{\text{res}}(\text{rob}) &= 7/15 H(\text{da}) + 8/15 H(\text{ne}) = \\ &= 7/15 * 0.591 + 8/15 * 0.954 = 0.784 \end{aligned}$$

$$H(\text{da}) = -1/7 \log_2(1/7) - 6/7 \log_2(6/7) = 0.591$$

$$H(\text{ne}) = -5/8 \log_2(5/8) - 3/8 \log_2(3/8) = 0.954$$

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

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