

Alberi decisionali

- Si consideri il dataset

Istanza	a1	a2	classe
1	T	T	+
2	T	T	+
3	T	F	-
4	F	F	+
5	F	T	-
6	F	T	-

- Si costruisca a mano l'albero seguendo l'algoritmo di c4.5 nell'ipotesi che il minimo numero di esempi in almeno due sottoinsiemi sia pari a 2

1

Risposta

- Nodo radice: $\text{info}(T)=1$
- Test su A1: $\text{info}_{A1}(T)=3/6*(-2/3*\log_2 2/3-1/3*\log_2 1/3)+3/6*(-1/3*\log_2 1/3-2/3*\log_2 2/3)=0.5*0.92+0.5*0.92=0.92$
 $\text{info}(A1)=1$
- $\text{gain}(A1)=(1-0.92)/1=0.08$
- Test su A2: $\text{info}_{A2}(T)=4/6*(-2/4*\log_2 2/4-2/4*\log_2 2/4)+2/6*(-1/2*\log_2 1/2-1/2*\log_2 1/2)=0.66*1+0.33*1=1$
 $\text{info}(A2)=-4/6*\log_2 4/6-2/6*\log_2 2/6=0.92$
- $\text{gain}(A2)=(1-1)/0.92=0$
- Viene preferito A1

2

Risposta

• $T_{A1=T}=\{$

1 T T +

2 T T +

3 T F -

$\}$

$T_{A1=F}=\{$

4 F F +

5 F T -

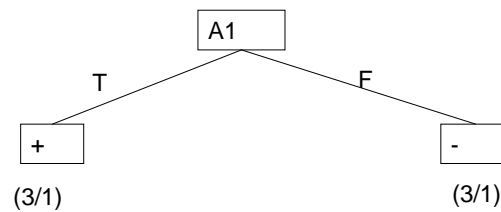
6 F T -

$\}$

- c4.5 si ferma in quanto $T_{A1=T}$ e $T_{A1=F}$ non possono essere suddivisi in modo che almeno due sottoinsiemi abbiano almeno due elementi

3

Risposta



4

Alberi decisionali

- Il maestro Yoda e' preoccupato dal numero di apprendisti Jedi che hanno deciso di darsi al Lato Oscuro, quindi ha deciso di apprendere un albero di decisione su alcuni dati storici per identificare i casi problematici in futuro.
- La tabella T descrive 12 iniziati specificando se sono passati al Lato Oscuro sulla base dell'eta' in cui il loro apprendistato Jedi e' cominciato, se hanno completato il loro apprendistato, la loro disposizione generale e la loro specie.

5

Tabella T

Eta' di inizio dell'apprendistato	Apprendistato completato	Disposizione	Specie	Lato Oscuro
5	1	Felice	Umana	0
9	1	Felice	Gungan	0
6	0	Felice	Wookie	0
6	1	Triste	Mon Calamari	0
7	0	Triste	Umana	0
8	1	Arrabbiata	Umana	0
5	1	Arrabbiata	Ewok	0
9	0	Felice	Ewok	1
8	0	Triste	Umana	1
8	0	Triste	Umana	1
6	0	Arrabbiata	Wookie	1
7	0	Arrabbiata	Mon Calamari	1

6

Domande

- Qual'è l'entropia della tabella T rispetto all'attributo Lato Oscuro?

$$info(S) = -\sum_{j=1}^k \frac{freq(C_j, S)}{|S|} \times \log_2 \left(\frac{freq(C_j, S)}{|S|} \right)$$

- Dire ad occhio (senza calcolare la funzione euristica) quale attributo verrebbe scelto come radice dell'albero dall'algoritmo di apprendimento di alberi di decisione?
- Qual'è il guadagno di informazione dell'attributo scelto nella risposta precedente?

7

Risposte

- $info(T) = -5/12 \log_2(5/12) - 7/12 \log_2(7/12) = 0.980$
- Apprendistato completato
- $gain(App) = info(T) - (5/12 * (-0.5 \log_2(0.5) - 5/5 \log_2(5/5))) + 7/12 * (-5/7 \log_2(5/7) - 2/7 \log_2(2/7)) = 0.476$

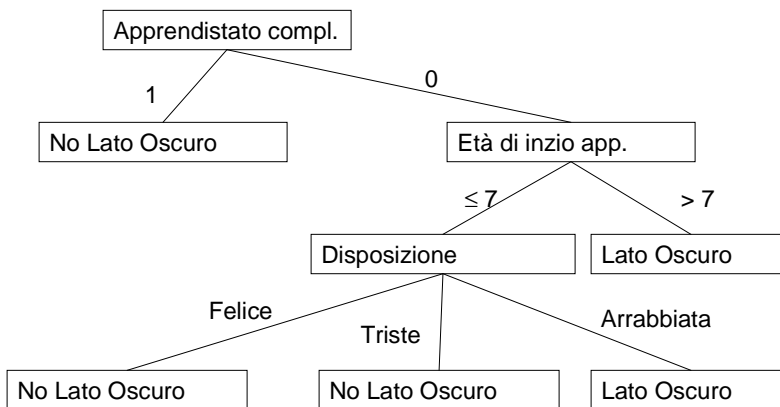
8

Domanda

- Disegnare l'albero decisionale che sarebbe appreso da questi dati (suggerimento: l'albero ha al massimo 3 divisioni. Si costruisca l'albero guardando solo la composizione degli insiemi, senza calcolare il guadagno di informazione)

9

Risposta



10

Domanda

- Si consideri la seguente tabella

Es.	attributi										Dec
	alt.	bar	V/S	fame	noC	prez	piov	pren	tipo	att	
x1	si	no	no	si	alc	£££	no	si	F	0-10	Si
x2	si	no	no	si	pieno	£	no	no	Thai	30-60	No
x3	no	si	no	no	alc	£	no	no	hamb	0-10	Si
x4	si	no	si	si	pieno	£	no	no	Thai	10-30	Si
x5	si	no	si	no	pieno	£££	no	si	F	>60	No
x6	no	si	no	si	alc	££	si	si	I	0-10	Si
x7	no	si	no	no	ness	£	si	no	hamb	0-10	No
x8	no	no	no	si	alc	££	si	si	Thai	0-10	Si
x9	no	si	si	no	pieno	£	si	no	hamb	>60	No
x10	si	si	si	si	pieno	£££	no	si	I	10-30	No
x11	no	no	no	no	ness	£	no	no	Thai	0-10	No
x12	si	si	si	si	pieno	£	no	no	hamb	30-60	Si

11

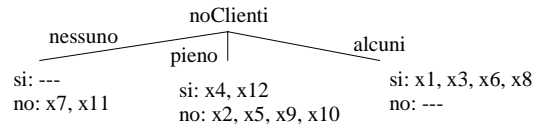
Domanda

- Si costruisca un albero di decisione a partire dalla tabella.
 - Al primo passo si scelga usando il criterio del guadagno quale tra gli attributi no Clienti e tipo e' piu' conveniente usare.
 - Dopo questa scelta si prosegua selezionando l'attributo fame, poi l'attributo tra noClienti e tipo non ancora usato per la generazione dell'albero e poi l'attributo ven/sab.
 - NOTA: si prosegua fino a trovare foglie omogenee senza usare il criterio di terminazione di C4.5

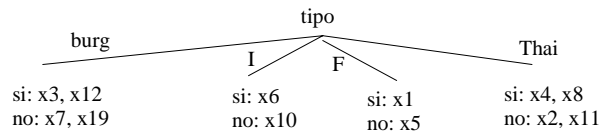
12

Soluzione

- attributo **noClienti**: per due valori discrimina completamente ("nessuno" e "alcuni")



- attributo **tipo**: discrimina male per tutti i valori

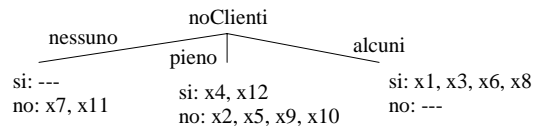


- Tra i due noClienti è la scelta migliore
- in generale tra tutti è quello con entropia più bassa

13

Soluzione

- Esempio
 - attributo **noClienti**: per due valori discrimina completamente ("nessuno" e "alcuni")

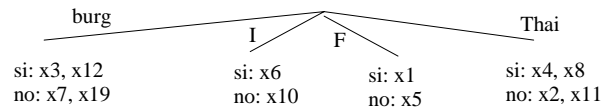


- $\text{info}(T) = -1/2 \log_2(1/2) - 1/2 \log_2(1/2) = 1$
- $\text{info}(T_{\text{clienti}}) = 2/12 \text{info}(T_{\text{alcuni}}) + 4/12 \text{info}(T_{\text{ness}}) + 6/12 \text{info}(T_{\text{pieno}})$
 - $\text{info}(T_{\text{alcuni}}) = 0$
 - $\text{info}(T_{\text{ness}}) = 0$
 - $\text{info}(T_{\text{pieno}}) = -2/6 \log_2(2/6) - 4/6 \log_2(4/6)$
- $\text{gain}(T_{\text{clienti}}) = 1 - \text{info}(T_{\text{clienti}}) \approx 0.541$

14

Soluzione

- Esempio
 - attributo **tipo**: discrimina male per tutti i valori

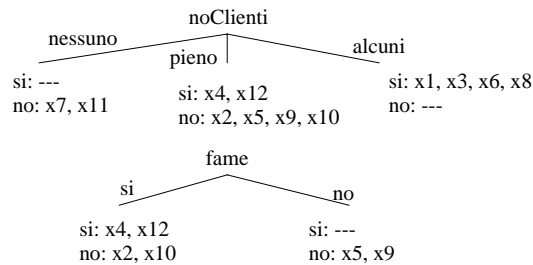


- $\text{info}(T) = -1/2 \log_2(1/2) - 1/2 \log_2(1/2) = 1$
- $\text{info}(T_{\text{tipo}}) = 2/12 \text{info}(T_I) + 2/12 \text{info}(T_F) + 4/12 \text{info}(T_{\text{burg}}) + 4/12 \text{info}(T_{\text{Thai}})$
 - $\text{info}(T_I) = 1$
 - $\text{info}(T_F) = 1$
 - $\text{info}(T_{\text{burg}}) = 1$
 - $\text{info}(T_{\text{Thai}}) = 1$
- $\text{gain}(T_{\text{tipo}}) = 1 - 1 = 0$

15

Soluzione

- L'algoritmo procede ricorsivamente considerando il valore "pieno" di noClienti e considerando gli esempi per quel valore
 - si analizzano gli altri attributi
 - si seleziona quello che discrimina meglio
 nel caso "fame": per uno dei due valori si ha classificazione completa



- Analogamente si procede ricorsivamente su ramo "si"

16

Albero risultante

