

Învățare automată

— Licență, anul III, 2018-2019, re-examinare, parțial I —

Nume student:

Grupa:

1. (Câștigul de informație — determinarea celui mai „bun” atribut; arbori de decizie consistenti cu un set de date)

Eticheta instanței	X_1	X_2	X_3	X_4
1	T	T	T	F
1	T	T	T	F
1	F	T	T	F
1	F	T	F	F
0	T	T	F	F
0	T	T	F	F
0	F	T	F	F
0	F	T	F	F

Fie atributele binare X_1, X_2, X_3, X_4 (ale căror valori pot fi doar T sau F), precum și două tipuri de etichete, 0 și 1. Veți considera cele 8 instanțe din tabelul alăturat.

a. Vrem să învățăm un arbore de decizie din acest set de exemple. În vederea selectării celui mai bun candidat pentru nodul rădăcină, calculați câștigul de informație pentru fiecare atribut X_i , cu $i = 1, \dots, 4$. În prealabil, veți desena compasul de decizie corespunzător fiecărui atribut. (La calcul câștigului de informație puteți folosi aproximația următoare: $\log_2 3 = 1.585$.) Ce atribut veți selecta?

Observație: Veți scrie definiția câștigului de informație. În cazul în care pentru calculul câștigurilor de informație veți folosi o altă formulă decât definiția, vă cerem să demonstrați formula respectivă.

b. Există oare un arbore de decizie care poate clasifica în mod perfect instanțele date? În caz afirmativ, desenați acel arbore de decizie. În caz contrar, dați o explicație simplă.

Handwritten calculations for information gain:

For X_1 : $H(1) = 1$, $H(0) = 1$, $H(1|X_1) = 0$, $H(0|X_1) = 0$, $IG(X_1) = 0$.

For X_2 : $H(1) = 1$, $H(0) = 1$, $H(1|X_2) = 0$, $H(0|X_2) = 0$, $IG(X_2) = 0$.

For X_3 : $H(1) = 1$, $H(0) = 1$, $H(1|X_3) = 2 \cdot \frac{1}{2} \cdot H(\frac{3}{4}) = H(\frac{3}{4}) = \frac{1}{2} \cdot 2 + \frac{3}{4} \cdot \log_2 \frac{4}{3} = 1.585$, $IG(X_3) = 1 - 1.585 = -0.585$.

For X_4 : $H(1) = 1$, $H(0) = 1$, $H(1|X_4) = 0$, $H(0|X_4) = 0$, $IG(X_4) = 0$.

Conclusion: $IG(X_3) = 1 - 0.585 = 0.415$ is the highest, so X_3 is selected.

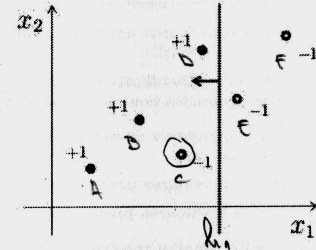
2.

(AdaBoost: întrebări în legătură cu aplicarea algoritmului pe un set de date din \mathbb{R}^2)

Folosind algoritmul AdaBoost, vrem să obținem un ansamblu de compasi de decizie (engl., decision stumps) h_t , de forma

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

În figura alăturată sunt desenate câteva puncte (instanțe) etichetate în planul bidimensional, precum și primul compas de decizie care a fost ales de către algoritmul AdaBoost. Un compas de decizie oarecare produce valori binare ± 1 , ținând cont doar de un anumit prag (engl., the split point). Săgeata mică din figură, care este perpendiculară pe dreapta care reprezintă compasul de decizie indică zona de decizie pentru care compasul de decizie va produce valoarea +1.



a. Încercuți toate acele instanțe din figură pentru care ponderea / probabilitatea [atribuită de către AdaBoost] va crește ca urmare a incorporării [în ipoteza combinată H] primului compas de decizie. Justificați răspunsul în mod riguros.

b. Desenați pe aceeași figură un compas de decizie care va putea fi selectat la următoarea iterație a algoritmului AdaBoost. Veți trasa atât dreapta care reprezintă compasul de decizie cât și o săgeată care să indice zona sa de decizie pozitivă. Justificați în mod riguros.

c. Va fi oare coeficientul / votul α_2 , care este asociat celui de-al doilea compas de decizie mai mare decât α_1 , coeficientul [din ansamblul H] pentru primul compas de decizie? Cu alte cuvinte, vom avea oare $\alpha_2 > \alpha_1$? Justificați în mod riguros.

Handwritten calculations for AdaBoost:

a. $\frac{1}{2} \ln \frac{1+\epsilon_1}{1-\epsilon_1} = \frac{1}{2} \ln \frac{1+\frac{1}{6}}{1-\frac{1}{6}} = \frac{1}{2} \ln \frac{7}{5} = \frac{1}{2} \ln 1.4 = 0.183$.

b. A decision boundary is drawn, and a small arrow indicates the region where the decision is +1.

c. $\alpha_2 > \alpha_1$ because $\epsilon_2 < \epsilon_1$. $\alpha_2 = \frac{1}{2} \ln \frac{1+\epsilon_2}{1-\epsilon_2} = \frac{1}{2} \ln \frac{1+\frac{1}{10}}{1-\frac{1}{10}} = \frac{1}{2} \ln \frac{11}{9} = \frac{1}{2} \ln 1.222 = 0.105$. Since $\alpha_2 > \alpha_1$, the answer is yes.