

Test AD si Regresie

RD (beat) >>> Dorinel

Regularizarea regresiiilor liniare și logistice presupune

Select one:

- a. Selectarea acelor attribute ale datelor de intrare care sunt mai informative
- b. Adăugarea unui termen în funcția de cost care penalizează ponderile cu valori mari**
- c. Creșterea dimensiunii setului de antrenare prin adăugarea de noi exemple
- d. Transformarea spațiului de intrare într-unul mai complex

Care dintre următoarele reprezintă criterii de oprire pentru construcția arborelui de decizie cu algoritmul ID3. Alegeți opțiunile corecte dintre:

1. Toate exemplele din subset fac parte din aceeași clasă
2. Câștigul informațional ajunge să fie negativ
3. Nu mai există attribute valide cu care să creăm un nod test
4. Adâncimea maximă este atinsă
5. Câștigul informațional depășește un anumit prag

☐ Select one:

☐ a.1

☐ b.4 și 5

☒ c.1, 3 și 4

☐ d.1 și 4

☐ e.1, 2 și 3

Care dintre următoarele reprezintă criterii de oprire pentru construcția arborelui de decizie cu algoritmul ID3. Alegeți opțiunile corecte dintre:

1. Toate exemplele din subset fac parte din aceeași clasă
2. Câștigul informațional ajunge să fie negativ
3. Nu mai există attribute valide cu care să creăm un nod test
4. Adâncimea maximă este atinsă
5. Câștigul informațional depășește un anumit prag

Select one:

a.

1, 2 și 3

b.

1 și 4

c.

4 și 5

d.

1, 3 și 4

e.

1

care e aici?

Arborii de decizie construiți cu algoritmul ID3 prezintă robustețe crescută la „outliers”?
Care dintre răspunsuri este cel mai plauzibil?

Select one:

a. Nu, pentru că vom obține valori negative pentru câștigul informațional în cazul lor.

b. Da, pentru că unele attribute nu sunt luate în considerare.++

c. Nu, pentru că ele (outliers) influențează sub-arborii generați ca fii ai unui nod test.

d.Da, întrucât acest lucru este garantat prin calculul entropiei. zic eu ++

<https://datascience.stackexchange.com/questions/37394/are-decision-trees-robust-to-outliers>

Ce tip de regresie ar trebui folosită pentru o problemă în care dorim să împărțim un set de date în 2 categorii? (presupunem că dorim să antrenăm un algoritm pentru care avem setul de date de antrenare)

Select one:

a.Regresie liniară deoarece cele 2 categorii (din setul de antrenare) pot fi etichetate cu 0 și 1. În etapa de antrenare algoritmul va avea de găsit o dreaptă care trece prin cele 2 puncte (0 și 1)

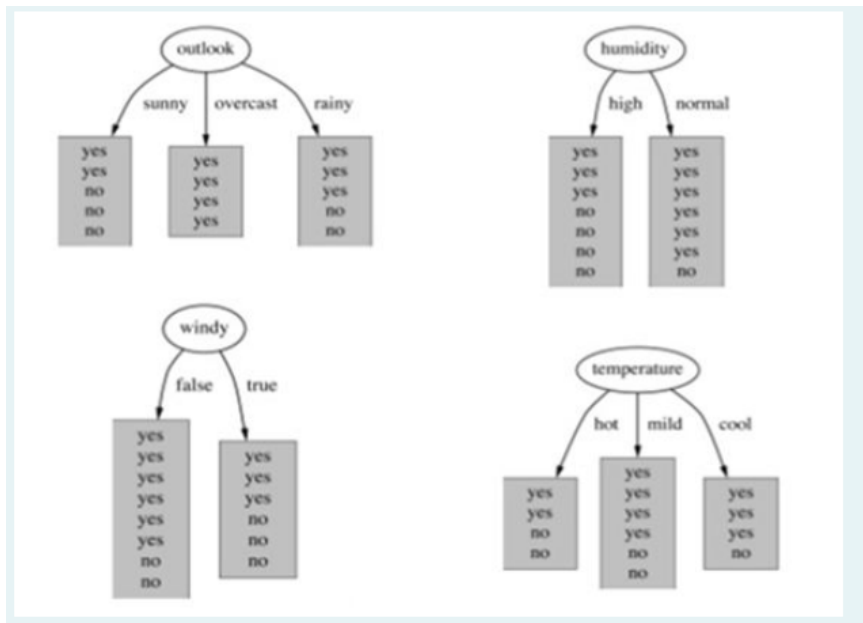
b.Regresie liniară deoarece datele pot fi separate cu ajutorul unei drepte

c.Regresie logistică deoarece nu știm cum să codificăm "categoriile" în nume reale necesare regresiei liniare.

d.Regresia logistică deoarece folosește funcția logistică pentru a modela probabilitatea ca un punct din setul de date să se găsească într-o anumită clasă.

Se antrenează un arbore de decizie având "maximum information gain" (câștigul maxim informațional, bazat pe diferența de entropie) drept criteriu de selecție a atributelor.

Considerând acest criteriu și următoarea situație de antrenare (cf. imaginii), care va fi atributul selectat pentru a "împărți" nodul curent:



Select one:

- a.Outlook (+?) ++
- b.Temperature (?)
- c.Humidity
- d.Windy

Dorim să prezicem prețul unui apartament din zona Pipera. Pentru a face acest lucru ce putem folosi dintre:

Select one: **ziceti repede care e**

- a.Regresie liniară **sau arbori de decizie** ++++
- b.Arbori de deciziecare di
- c.Regresie logistică
- d.Regresie liniară
- e.Regresie liniară sau logistică

-----over-----

In cazul algoritmului C4.5 daca avem valori necunoscute pentru un atribut in anumite
exemple de invatare trebuie sa

Select one:

- a. atribuim acelui atribut valoarea cu cea mai mare frecventa din exemplele de invatare
- b. atribuim acelui atribut valoarea cu cea mai mica frecventa din exemplele de invatare

Care dintre urmatoarele sunt bune practici pentru reducerea overfitting-ului?

- (a) Utilizarea unei functii de cost cu 2 componente, care include un regulator pentru a penaliza complexitatea modelului
- (b) Utilizarea unui optimizator bun pentru a reduce erorile pe datele de antrenare
- (c) Construirea unei structuri de subseturi de modele imbricate, antrenarea pe fiecare subset pornind de la cel mai mic, si oprirea cand eroarea de cross-validare incepe sa creasca
- (d) Eliminarea aleatoare a 50% din datele de antrenare

Select one:

- a. (a) si (c)
- b. (a) si (b) si (c)
- c. (b) si (c)
- d. (a) si (b)
- e. (c)

macar a si c cred ++

Ce presupune "lama lui Occam" în contextul arborilor de decizie?

Select one:

- a. Dacă sunt mai mulți arbori de decizie corecți, se preferă cel mai simplu**
- b. Un arbore corect de decizie trebuie să cuprindă toate atributele din setul de date
- c. Nu există mai mulți arbori corecți pentru aceeași problemă
- d. Nu se poate aplica principiul în contextul arborilor de decizie

În cazul algoritmului C4.5 dacă avem valori necunoscute pentru un atribut în anumite exemple de învățare trebuie să

Select one:

- a. atribuim acelui atribut valoarea cu cea mai mică frecvență din exemplele de învățare
- b. atribuim acelui atribut valoarea cu cea mai mare frecvență din exemplele de învățare**

Într-un arbore de decizie pentru clasificarea unor exemple care conțin atribute cu valori continue, are sens să repetăm pe o anumită ramură a arborelui (cale de la rădăcina la o frunză) un același atribut.

Select one:

True

False

Avem următoarea matrice de confuzie pentru o problemă de clasificare binară (pe coloane sunt valorile reale, iar pe linii sunt valorile prezise).

Care afirmatie este adevarata?

	Positive	Negative
Positive	23	1
Negative	12	556

Select one:

a.

Accuracy=568/592

b.

Accuracy=557/556

c.

Recall=23/24

d.

Precision=23/24

Un set de date utilizat pentru învățare supervizată are o valoare mare a entropiei informaționale dacă:

Select one:

a.

Numărul de exemple din fiecare clasă este relativ similar

b.

Foarte multe dintre exemple aparțin unui număr mic de clase, în timp ce restul de clase au un număr mic de exemple. cred ca asta? asta pare

c.

Toate exemplele fac parte din aceeași clasă

Pipera duplicat

Ce presupune "lama lui Occam" în contextul arborilor de decizie?

Select one:

a.

Dacă sunt mai mulți arbori de decizie corecți, se preferă cel mai simplu

b.

Nu există mai mulți arbori corecți pentru aceeași problemă

c.

Nu se poate aplica principiul în contextul arborilor de decizie

d.

Un arbore corect de decizie trebuie să cuprindă toate atributele din setul de date

Avem următoarea matrice de confuzie pentru o problema de clasificare binară (pe coloane sunt valorile reale, iar pe linii sunt valorile prezise).

Care afirmație este adevărată?

	Positive	Negative
Positive	23	1
Negative	12	556

Select one:

a.

Accuracy=568/592

b.

Accuracy=557/556

c.

Recall=23/24

d.

Precision=23/24
