

SISTEME CU IA

RNA PENTRU CLASIFICARE SI RECUNOASTERE

Contextul general

Clasificarea - procesul de învățare a asemănarilor și deosebirilor dintre instanțele unor obiecte dintr-o populație de obiecte neidentice.

Contextul general

DOUA etape

Clasificarea - sistemul învață caracteristicile generale ale claselor pe baza caracteristicilor specifice instanțelor .

Recunoașterea - sistemul suprapune caracteristicile unei instanțe peste caracteristicile claselor cunoscute și identifică clasa căreia aparține.

Contextul general

DOUA tipuri de invatare

Invatare supravegheata – se aplica atunci cand pentru fiecare model de intrare (instanta a obiectelor clasificate) se cunoaste apriori clasa la care se asociaza.

Invatare nesupravegheata – se aplica atunci cand nu se cunoaste in prealabil clasa la care se asociaza fiecare model de intrare.

Masuri de similaritate / disimilaritate

ABORDARE

Se porneste de la un set de vectori de intrare $X = \{x(1), x(2), \dots, x(M)\}$, cu structura $x(m) = (x1(m), x2(m), \dots, xN(m))$.

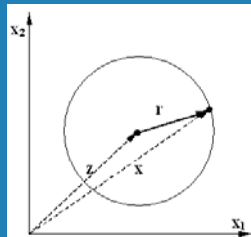
Se doreste separe a K clase – notate $X1, X2, \dots, XK$ – fiecare caracterizată de un prototip z^k ($k = 1, \dots, K$).

In final, pe baza celor K prototipuri z^k se asociaza fiecare vector $x(m)$ la una din cele K clase.

Masuri de similaritate / disimilaritate

Distanta Euclidiană – norma 2

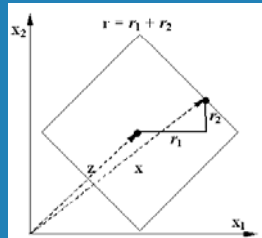
$$\|x - z^k\|_2 = \sqrt{\sum_{n=1}^N (x_n - z_n^k)^2}$$



Masuri de similaritate / disimilaritate

Distanța după norma 1

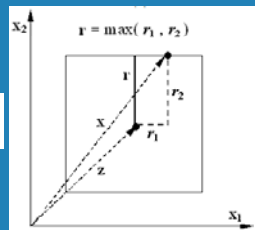
$$\|x - z^k\|_1 = \sum_{n=1}^N |x_n - z_n^k|$$



Masuri de similaritate / disimilaritate

Distanța după norma ∞

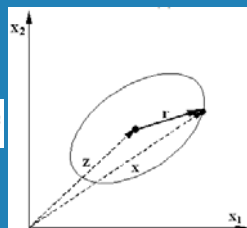
$$\|x - z^k\|_\infty = \max_{n=1, \dots, N} (|x_n - z_n^k|)$$



Masuri de similaritate / disimilaritate

Distanța după ponderată (Mahalanobis)

$$\|x - \mu\|_M = [(x - \mu)^T \cdot C^{-1} \cdot (x - \mu)]^{1/2}$$



Masuri de similaritate / disimilaritate

Distanța după norma Hamming

$$\|x - z\|_H = \sum_{n=1}^N |x_n - z_n|$$

În cazul reprezentării binare, **distanța Hamming** măsoară numărul de componente diferite ale vectorilor **x** și **z**.

Pentru: **x = (1101)** și **z = (0110)** se obține:

$$\|x - z\|_H = |1 - 0| + |1 - 1| + |0 - 1| + |1 - 0| = 3$$

Algoritmi / Modele de clasificare

1. Algoritmii celor mai apropiați **K** vecini
2. Algoritmii **C**-medii
3. Algoritmii ISODATA
4. Rețele Kohonen (cu autoorganizare)
 - 4.1. Rețele VQ – **Vector Quantization**
 - 4.2. Rețele LVQ – **Learning Vector Quantization**
 - 4.3. Rețele SOFM – **Self Organizing Feature Maps**

Algoritmii celor mai apropiați **K** vecini

Principiu

Se realizează selectarea dintre vectorii deja clasificați a primilor **K** (valoarea lui **K** este precizată a priori), cei mai apropiați de vectorul curent, iar acesta din urmă este asociat clasei dominante asociată celor **K** vectori de referință.

Algoritmul celor mai apropiați K vecini

Initializari

Se pornește de la setul de învățare format din vectorii $\{x^{(1)}, x^{(2)}, \dots, x^{(M)}\}$ și de la numărul de clase C în care se face clasificarea. Atunci când unul din acești vectori $x^{(m)}$ este asociat unei clase c , se folosește notația:

$$\text{Clasa}(m) = c$$

După reordonarea aleatorie a vectorilor $x^{(m)}$ din setul de antrenare se consideră arbitrar că primii C vectori sunt incluși fiecare în una din cele C clase.

Algoritmul celor mai apropiați K vecini

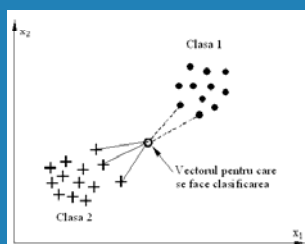
Clasificarea propriu-zisă

În continuare, pentru fiecare din restul vectorilor $x^{(C+1)}, x^{(C+2)}, \dots, x^{(M)}$ se parcurg următorii pași:

- se calculează distanțele față de vectorii deja clasificați;
- se ordonează crescător aceste distanțe și se consideră vectorii corespunzători primelor K distanțe reordonate;
- pentru cei K vectori se determină clasa dominantă, iar vectorul curent se asociază acestei clase.

Algoritmul celor mai apropiați K vecini

Selectarea celor mai apropiați K vecini (cazul K = 5)



Algoritmul celor mai apropiați K vecini

Algoritmul

1. Date de intrare: modelele de antrenare, sub forma vectorilor de intrare $x^{(m)}$, $m = 1, \dots, M$; parametrul k (numărul de vecini); numărul de clase C .
2. Reordonarea aleatorie a vectorilor de antrenare $\{x^{(m)}\}$ $m = 1, \dots, M$.
3. Asocierea primelor C modele la cele C clase:
for $c = 1$ to C do $Clasa(c) = c$
4. Asocierea claselor pentru restul modelelor de antrenare:
for $m = C + 1$ to M do
// Calculează distanțele față de vectorii deja clasificați
for $q = 1$ to $m - 1$ do
 $d(q) = \|x^{(m)} - x^{(q)}\|$

Algoritmul celor mai apropiați K vecini

Algoritmul - continuare

```
// Ordonează crescător distanțele  $d(q)$  prin reordonarea  
// indecșilor  $q$   
  
for  $i = 1$  to  $m - 1$  do  $lx(i) = i$   
for  $i = 1$  to  $m - 2$  do  
for  $j = i + 1$  to  $m - 1$  do  
if  $d(lx(j)) < d(lx(i))$  then  $lx(j) \leftrightarrow lx(i)$ 
```

Algoritmul celor mai apropiați K vecini

Algoritmul - continuare

```
// Calculează numărul de vectori asociați fiecărei clase  
// pentru primele  $K$  distanțe  $d(q)$ , corespunzător  
// indecșilor reordonați  
 $vmax = 0$ ;  
for  $c = 1$  to  $C$  do  
v = 0  
for  $j = 1$  to  $k$  do  
if  $Clasa(lx(j)) = c$  then  $v = v + 1$   
if  $v > vmax$  then  
Clasa(m) = c  
vmax = v  
  
END.
```

Algoritmul celor mai apropiați K vecini

Observatii

- (a) necesită precizarea a priori a numărului de clase ce urmează a fi separate;
- (b) Rezultatele clasificării sunt influențate de ordinea de prezentare a vectorilor;
- (c) chiar principiul celor mai apropiați vecini compară vectorul curent cu vectori care, în majoritatea cazurilor, se găsesc către periferia zonelor
- (d) nu stabilește vectori caracteristici / prototipuri.

Algoritmul C - medii

Principiu

Algoritmul c-medii introduce pentru prima dată noțiunea de **prototip** sau **vector-centru** sau **vector de codare**, care descrie în mod unitar o clasă. Astfel, pentru o clasă c , prototipul $\mathbf{z}^{(c)}$ se calculează ca medie aritmetică a celor $n(c)$ vectori ce au fost asociați clasei respective:

$$\mathbf{z}^{(c)} = \frac{1}{n(c)} \sum_{m=1}^M \mathbf{x}^{(m)}$$

$\mathbf{z}^{(c)} = \mathbf{x}^{(m)} \quad \forall m$, cu proprietatea **Clasa (m) = c**

Algoritmul C - medii

Initializare

Se pornește de la setul de antrenare $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(M)}\}$ și de la numărul de clase $C < M$ care trebuie separate. După reordonarea aleatorie a setului de antrenare, se atribuie arbitrar primele C modele de antrenare celor C clase, vectorii $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(C)}$ devenind prototipurile $\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(C)}$.

Algoritmul C - medii

Clasificarea propriu-zisa

Fiecare din restul de $M - C$ modele de antrenare este asociat unei clase, pe baza distanțelor minime față de cele C prototipuri.

Se recalculeaza prototipurile claselor și se fac noi asocieri.

Procesul se reia, până când – în două iterații succesive – prototipurile claselor nu se schimbă sau se schimbă într-o măsură ne semnificativă.

Algoritmul C - medii

Cuantificarea preciziei

Evaluarea dimensiunilor zonei acoperite de fiecare clasă – abaterea medie pătratică dintre prototipul clasei și vectorii din setul de antrenare asociați acesteia:

$$\sigma_c^2 = \frac{1}{n(c)} \sum_{m=1}^M \|\mathbf{x}^{(m)} - \mathbf{z}^{(c)}\|^2 \quad \forall m, \text{ cu proprietatea } \text{Clasa}(m) = c$$

Abaterea patratice totala:

$$\sigma_T^2 = \sum_{c=1}^C \sigma_c^2$$

reprezintă o măsură a preciziei de clasificare.

Algoritmul C - medii

Algoritmul

1. Date de intrare: modelele de antrenare, sub forma vectorilor de intrare $\mathbf{x}^{(m)}$, $m = 1, \dots, M$; numărul de clase C .
2. Reordonarea aleatorie a vectorilor de antrenare $\{\mathbf{x}^{(m)}\}$ $m = 1, \dots, M$.
3. Asocierea primelor C modele la cele C clase:

for $c = 1$ to C do

$\mathbf{z}^{(c)} = \mathbf{x}^{(c)}$; $n(c) = 0$

Algoritmul C - medii

Algoritmul - continuare

4. Asocierea fiecărui model de antrenare la clasa cu prototipul cel mai apropiat:

```
for m = 1 to M do
  Dmin = 106 // Inițializarea distanței minime
  for c = 1 to C do
    if ||x(m) - z(c)|| < Dmin then
      Dmin = ||x(m) - z(c)||
      Cmin = c
  Clasa(m) = Cmin
  n(Cmin) = n(Cmin) + 1
```

Algoritmul C - medii

Algoritmul - continuare

5. Calculul noilor prototipuri pentru cele C clase:

```
for c = 1 to C do
  w(c) = 0
  for m = 1 to M do
    if Clasa(m) = c then
      w(c) = w(c) + x(m)
  w(c) = w(c) / n(c) // Prototipuri temporare
```

Algoritmul C - medii

Algoritmul - continuare

6. Calculul abaterilor pătratice:

```
σT2 = 0
for c = 1 to C do
  σc2 = 0
  for m = 1 to M do
    if Clasa(m) = c then
      σc2 = σc2 + ||x(m) - w(c)||2
  σc2 = σc2 / n(c)
σT2 = σT2 + σc2
```

Algoritmul C - medii

Algoritmul - continuare

7. Criteriul de oprire:

dacă prototipurile s – au modificat nesemnificativ – adică

$$||z^{(c)} - w^{(c)}|| < \varepsilon \text{ pentru toate centrele } c = 1, \dots, C$$

algoritmul se încheie. În caz contrar, se memorează noile prototipuri:
for $c = 1$ to C do $z^{(c)} = w^{(c)}$
și se revine la pasul 4.

Algoritmul C - medii

Observatii

- (a) necesită precizarea a priori a numărului de clase ce urmează a fi separate;
- (b) nu este garantată convergența;
- (c) rezultate mai bune decât algoritmul celor mai apropiați K vecini
- (d) Da rezultate bune, în special în cazul claselor net separabile.
