

b. soluții:

b.1. proceduri bazate pe gradient: $a_n = a_{n-1} + \delta_n \nabla_{\mathbf{a}} Q(y_n, f(x_n, a_{n-1}))$, unde a_n = vectorul de probe la pasul n , δ_n = rata de învățare, $\nabla_{\mathbf{a}}$ = gradientul (în ce direcție și cu ce intensitate trebuie modificat parametrul), $Q(y_n, f(x_n, a_{n-1}))$ = o funcție, care măsoară eroarea (sau pierdere) între valoarea actuală a_{n-1} și predicția modelului $f(x_n, a_{n-1})$.

b.2. ERM cu condiție convergenței uniforme;

c. Demonstrații:

- c.1. Teoria aproximațiilor stochastice prin procedurile bazate pe gradient
- c.2. Teoria convergenței uniforme pentru principiul ERM.

Perceptronul - primește toate inputurile, cu ponderile lor și, printr-o funcție ce utilizează un prag, returnează una din 2 valori (peste/sub prag). Ultimul input are, de obicei, valoare și corespunde la deplasare (bias) pentru perceptron. Notăm $y = \text{sgn}(x^T w)$, funcția ce returnează +1 sau -1, pornind de la input și vectorul comun ponderilor w . Rezultatul (+1) e compari cu datele din setul de antrenament, numite generic, $d: \{x^i, d^i\}$. Scopul perceptronului este găsirea vectorului ponderilor optim. General, putem considera că perceptronul împarte un spațiu în 2 părți.

Regula de învățare a perceptronilor (Rosenblatt, 1962): $w^{k+1} = w^k + \eta (d^k - y^k) x^k$, unde η = constantă denumită rata de învățare;

d^k, y^k din $\{-1, +1\}$. Există soluție (3 una \rightarrow 3 o infinitate) doar dacă $k=1, 2, \dots, n$, antrenamentul este separabil liniar (subplanul doar cu valori +1 și cel doar cu -1).

Se alege $w^1 = 0 \Rightarrow$ limita de corecție pe care o poate face perceptronul în procesul de învățare k_0 : $k_0 = \frac{\alpha^2 \|w^*\|^2}{\beta^2} = \frac{\beta^2}{\alpha^2} \|w^*\|^2$ sau $k_0 = \frac{\max_i \|x^i\|^2}{[\min_i (x^i)^T w^*]^2}$

- unde w^* e vectorul de greutate ideal
- x^i sunt exemplele de antrenament
- k_0 = numărul maxim de corecții, pe care le face în procesul de învățare
- derivăm și obținem formule alăturată, unde: $\beta = \max_i \|x^i\|$ - maximul normelor (lungimilor) exemplarelor de antrenament

Regula de învățare Widrow-Hoff (n-LMS): cea mai utilizată implementare a principiului perturbării minime

$\delta = \min_i x_i^T w^*$ - margin: distanța min a punctelor față de hiperplanul ideal definit de w^* .

actualizează ponderile cât mai puțin în direcție erorii \Rightarrow învățare stohastică și lină, nu heuristică, ce duce ponderile la w^* . Acest vector minimizează media pătratelor erorilor dintre rețeauă și valorile corecte:

$J(w) = \frac{1}{2} \sum (d^i - y^i)^2$. Formula: $w^{k+1} = w^k - \mu \nabla J(w) = w^k + \mu \sum_{i=1}^m (d^i - y^i) x^i$

- w = ponderile, μ = rata de învățare, d = etichetă, y = rezultat din rețea, $d^i = y^i$ = eroare, m = pasul. Putem folosi J pentru a regla gradientul descendent pe partea cea mai abruptă a graficului $J(w)$ în punctul w^k , făcând rezultatul Δw^k proporțional cu gradientul lui J în w^k . Obținem regula cunoscută ca gradientului descendent, cel mai abrupt, sau a gradientului descendent. $J(w)$ definește o suprafață hiperparabolică convexă cu un singur minimum: w^* (cel global).

(COURS)

$$w^{k+1} = w^k - p \nabla J(w) \Big|_{w=w^k} = w^k - p \left[\frac{\partial J}{\partial w_1} \frac{\partial J}{\partial w_2} \dots \frac{\partial J}{\partial w_{n+1}} \right]^T \Big|_{w=w^k}$$

> p = rata de învățare (step-size) - cît de mare e pasul la $\nabla J(w)$

este vectorul din rețea de direcție în care $J(w)$ scade cel mai rapid → soluția optimă ei este la capăt. Aceasta este regula batch LMS (Least Mean Squares) - se actualizează după toate exemplele, o dată pe pas.

- Regula μ -LMS sau LMS (se actualizează după fiecare exemplu/pas): $w' = 0$ sau arbitrar,

$$w^{k+1} = w^k + \mu (d^k - y^k) x^k, \text{ unde } 0 < \mu < 2 / (\max_i \|x^i\|^2). \text{ Se poate demonstra c\^a:}$$

$$w^* = X^+ d, \text{ unde } X = [x^1 x^2 \dots x^m] \text{ matricea de intrare, } d = [d^1 d^2 \dots d^m]^T \text{ vectorul de etichete,}$$

dar $X^+ = (X X^T)^{-1} X$ - e inversa generalizată → putem obține direct (fără iterații) soluția w^* dacă e o problemă liniar-separabilă și avem toate datele,

dar: cost computațional mare pentru mărime + nu se poate aplica dacă datele vin în flux.

Regula batch e mai eficientă, dar poate fi aplicată însoțită de o diferență de timp din ce în ce mai neglijabilă, cu cît se folosesc achiziții de epoci consecutive mai mici (ex. 20) pt. învățare.

Regula de învățare Delta - Extensie a regulii LMS, care permite învățarea și a funcției de activare este neliniară, dar derivabilă. Aplicăm gradientul descendent → $\nabla J(w) =$

$$= -(d - y) f'(\text{net}) x, \text{ unde } (d - y) = \text{eror între valorile dorite și valorile actuale};$$

$f'(\text{net})$ = derivata funcției de activare față de suma ponderată de intrări w ; x = vectorul de intrare. → regula Delta: $w^{k+1} = w^k + p [d^k - f(\text{net}^k)] f'(\text{net}^k) x^k = w^k + p \delta^k x^k$

unde $\text{net}^k = (x^k)^T w^k$ și $f' = \frac{df}{d\text{net}}$, dar $(x^k)^T$ e vectorul de intrare de la pasul k , transpus.
 Dacă $f(\text{net}) = \tanh(\beta \text{net}) \Rightarrow f'(\text{net}) = \beta [1 - f^2(\text{net})]$. Dacă $f(\text{net}) = \frac{1}{1 + e^{-\beta \text{net}}} \Rightarrow$

$$\Rightarrow f'(\text{net}) = \beta f(\text{net}) [1 - f(\text{net})].$$

- Funcție de activare este o funcție matematică aplicată individual pe fiecare percepție (neuron), care primește ca argument suma ponderată a valorilor de la toate percepțiile din straturile anterioare.

- Pentru a preveni stagnarea când $f'(\text{net}) \rightarrow 0$, se adaugă un bias pozitiv ϵ :

$$w^{k+1} = w^k + p [d^k - f(\text{net}^k)] [f'(\text{net}^k) + \epsilon] x^k$$

Capabilități computaționale ale percepțiilor: Dacă avem un set de antrenament $T =$

$$= \{(x_i, t_i), i=1:n\} \text{ și seturile de inputuri } S^+ = \{x: (\exists i) t_i = 1, x = x_i\} \text{ și}$$

$$S^- = \{x: (\exists i) t_i = -1, x = x_i\} \text{ - seturi cu eticheta 1, respectiv -1.}$$

Def. Separabilitate linieară - spunem că S^+ și S^- sunt liniar separabile, dacă (\exists)

$$S^+ C H^+, S^- C H^- \text{ (H-urile sunt semiplanuri/semispafii) sau } (\forall x \in S^+) w^T x - b \geq 1$$

$(\forall x \in S^-) w^T x - b < 0$, dar b = pragul (echivalent cu biasul - în rețele neuronale)

(1) Un set de antrenament poate fi omotet fără eror de un perceptron, dacă acesta este liniar separabil.

Notăm nr. maxim de dimensiuni care poate fi introdus de un perceptron asupra unui set S de n elemente din R^d . $P(n, d) = \max_{S: \|S\| = n} L(S)$

Introducem sist de ecuații a diferinței finite: $\Delta(n, d) = \Delta(n-1, d) + \Delta(n-1, d-1)$ și
 corect $\Delta(n, d) = \text{nr. maxim de etichetări } (\pm 1) \text{ pe un set de } n \text{ puncte, este}$
 calculat decă adăugăm un pct nou (coste n) sau coste optinal (d) și satisface
 condiție: $\Delta(1, d) = \Delta(1, 1) = 2$
 $\Delta(n, 1) = \Delta(n, 1) = 2(n+1) - 2 = 2n$

1. Al un pct pe care-l pui etichetăm în 2 feluri: ± 1 .
2. Pentru optinal, cu 1 dimensiune (o axă), percepțional după axa în intervale prin
 n puncte \Rightarrow se pot genera doar 2n separări diferite (nu 2^n), între fiecare 2 puncte
 separate: $+1$ la stg și -1 la dreapta sau invers.

Teorema limitei superioare (pt. dihotomie / Teorema lui Cover): Ptn. (± 1) set de n puncte din
 optinal \mathbb{R}^d , numărul maxim de etichetări ± 1 (di. dihotomie), pe care le poate realiza
 un perceptron este: $\Delta(n, d) = \begin{cases} 2 \sum_{i=0}^d \binom{n-1}{i}, & \text{if } n > d+1 \\ 2^n, & \text{if } n \leq d+1 \end{cases}$, unde $\binom{n-1}{i}$ sunt combinații de
 $n-1$ luate câte $i =$

Deci, un perceptron poate înveta corect (± 1) set, care are $\Delta(n, d) \leq 2^n$.
 Deci, probabilitatea ca un perceptron Rosenblatt să fie capabil să distingă Telenor e:

$$P(n, d) \geq \Delta(n, d) / 2^n = \left(\frac{1}{2} \right)^{n-1} \sum_{k=0}^d \binom{n-1}{k}$$

inputuri binare: primesc ca input o secvență binară de lungime d (ex. 10...010...) și
 returnează 0 sau 1 ptn. fiecare. Numărul total de combinații (inputuri) posibile $\pm 2^d$.
 Numărul total de funcții booleene posibile cu d inputuri: $2(2^d)$ - pentru că fiecare combinație
 poate returna 2 valori: ± 1 . Numărul maxim de funcții booleene implementate de un perceptron

$$B(d) \leq \Delta(2^d, d) < 2^{d^2} / (d-1)!, \text{ unde } B(d) = \text{nr. exact de funcții booleene}$$

$\Delta(2^d, d) = \text{limita superioară teoretică (dar nu exactă) ptn. care dihotomie se pot realiza}$
 în \mathbb{R}^d cu 2^d puncte de un perceptron. $2(d^2)/(d-1)! \approx 0$ limită mare, dar care scade
 Regula de învățare adaptivă Ho-Kashyap (rapid ca proprietate din $2(2^d)$).

3 reguli (prinse de la LMS): $ATK \bar{1}$ (probleme separabile liniar), $ATK \bar{u}$ (probleme non-separabile liniar), $ATK \bar{u}$ (probleme separabile
 liniar), $ATK \bar{u}$ (probleme non-separabile liniar).

Considerăm o problemă de clasificare în 2 clase c_1 și c_2 a unui seturi de n
 numere, notate fiecare, generate, cu x . Avem m prechii (set de numere etichetate) ptn.
 antrenament $\{x^i, d^i\}$, $i=1, 2, \dots, m$, $x^i \in \mathbb{R}^{n+1}$, iar ultimul x^i este un bias constant
 cu valoarea 1. Considerăm sgn -funcția de activare a percepționalului. Atunci un
 perceptron poate fi învățat să clasifice corect prechile de antrenament, decă sunt

satisfăcute următoarele m inegalități, ptn. $d=1, \dots, m$:

$$w = \text{un singur vector de ponderi, care se aplică fiecărui } x^i: (x^i)^T w \begin{cases} > 0, \text{ if } d^i = +1 \\ < 0, \text{ if } d^i = -1 \end{cases}$$

Decă notăm vectorii $z^i = \begin{cases} +x^i, & \text{if } d^i = +1 \\ -x^i, & \text{if } d^i = -1 \end{cases}$ și avem pseudoinverse $Z =$
 (numită valoare) obținem $Z^T W = b$. Facem aceste (frecare b^i să fie peste 0
 prag de performanță a clasificării, ce trebuie egdat sau depășit de fiecare dată
 la alea arbitrari e o alegere prostă, el trebuie să se fie dinamic găsind.

În algoritmul Ho-Kashyap se încearcă așezarea la o valoare a lui b strict
 pozitiv de la un b mic. Apoi, se încearcă minimizarea funcției de cost:

$$J(w, b) = \frac{1}{2} \| Z^T w - b \|^2$$

Norma euclidiană a vectorului: $\|w\|^2 = w_1^2 + \dots + w_n^2$

Deci $w = Z^T b$, unde $Z^T = (ZZ^T)^{-1} Z$, ptr. $m > n+1$. În continuare, recalculăm unușorul $b^{k+1} > 0$ cu formule de gradient descent: $b^{k+1} = b^k + \frac{1}{2} [\Sigma + |\Sigma|^K]$, unde $\Sigma^K = Z^T W^K - b^K$. Apoi, recalculăm w . Așadar, ptr. $\Sigma = 0$ sau foarte aproape de pozitiv (indice separabilitate liniară) sau ptr. $\Sigma \leq 0$ (indice lipsa separabilității liniare). Calcularea directă a Z^T poate fi evitată cu un algoritm gradient-descent: procedură adaptivă mod batch Ho-Kashyap. Pornim de la $J(w, b) = \frac{1}{2} \|Z^T w - b\|^2$. Apoi, gradientul lui J în funcție de w și b este dat de: $\nabla_b J(w, b) | w^K, b^K = -(Z^T w^K - b^K)$ respectiv $\nabla_w J(w, b) | w^K, b^{k+1} = -Z(Z^T w^K - b^{k+1})$. Deci folosim în locul ecuației de gradient descent, $-0,5(\Sigma + |\Sigma^K|)$, cu $\Sigma^K = Z^T w^K - b^K$. atuncea: $b^{k+1} = b^K + \frac{\rho_1}{2} (|\Sigma^K| + \Sigma)$ cu $\Sigma^K = Z^T w^K - b^K$
 $w^{k+1} = w^K - \rho_2 Z(Z^T w^K - b^{k+1}) = w^K + \frac{\rho_1 \rho_2}{2} Z[|\Sigma^K| + \Sigma(1 - \frac{2}{\rho_1})]$ unde:

1) $\rho_2 = \text{constante de învățare strict pozitivă}$
 Dacă $\rho_1 = 0$ și $b^1 = 1$, ecuația se reduce la regula de învățare μ -LMS. Iar dacă $0 < \rho_1 < 2$
 și $0 < \rho_2 < \frac{2}{\lambda_{\max}}$, unde λ_{\max} e cea mai mare valoare proprie a matricii definite ZZ^T .

0 procedură adaptivă Ho-Kashyap se obține de la funcția costului de instantă:

$$J(w, b) = \frac{1}{2} [(z^i)^T w - b^i]^2, \text{ care se generalizează regulile de incrementare:}$$

$$b_i^{k+1} = b_i^K + \frac{\rho_1}{2} (|\Sigma_i^K| + \Sigma_i^K), \text{ cu } \Sigma_i^K = (z^i)^T w^K - b_i^K$$

$$w^{k+1} = w^K - \rho_2 Z^i [(z^i)^T w^K - b_i^{k+1}] = w^K + \frac{\rho_1 \rho_2}{2} [|\Sigma_i^K| + \Sigma_i^K(1 - \frac{2}{\rho_1})] z^i$$

Alternativ, scriem: $\Delta b_i = \rho_1 \Sigma_i^K$ și $\Delta w = \rho_2 (\rho_1 - 1) \Sigma_i^K z^i$, dacă $\Sigma_i^K < 0$

$\Delta b_i \geq 0$ și $\Delta w = -\rho_2 \Sigma_i^K z^i$, dacă $\Sigma_i^K \geq 0$, unde Δb și Δw = diferențele dintre valorile actualizate și curente ale lui b și w . Aceasta este regula de învățare AHK_I.

Regula de învățare AHK_I pornește de la AHK_I dar permite verificarea Δb și în special valorilor negative, cel timp nu se obține o mărime negativă:

$$\Delta b_i = \rho_1 \Sigma_i^K \text{ și } \Delta w = \rho_2 (\rho_1 - 1) \Sigma_i^K z^i \text{ dacă } b_i^K + \rho_1 \Sigma_i^K > 0$$

$$\Delta b_i \geq 0 \text{ și } \Delta w = -\rho_2 \Sigma_i^K z^i, \text{ dacă } b_i^K + \rho_1 \Sigma_i^K \leq 0, \text{ unde:}$$

$$0 < \rho_2 \leq \frac{2}{\max \|z^i\|^2} \text{ și } 0 < \rho_1 < 2$$

Regula de învățare AHK_{II} se poate aplica și problemelor ne-separabile liniar. Atci Δw e setat la 0 în regula AHK_{II} ptr. $\Delta b_i \geq 0$

Norma Manhattan

Versiune batch: $J(w) = \frac{1}{n} \sum_{i=1}^n |d^i - y^i|$ sau versiune de instantă:

$$J(w) = \frac{1}{n} |d^i - y^i|, \text{ Forme generală a funcției gradient:}$$

$$\nabla J(w) = -\text{sgn}(d - y) |d - y|^{n-1} f'(\text{net}) x$$

Regula de învățare prin coabare Pornește de la funcția costului de calculare a gradientului descent: $J(x) = -\sum_{i=1}^m y^i d^i$. Calculăm w : $w^1 = 0$, $w^{k+1} = w^k + \rho d^k x^k$
 unde $y^i = (x^i)^T w$

Deci setăm $p=1$, atunci $\Rightarrow w^x = \sum_{i=1}^m d_i x^i = Xd$.

Regula de înmulțire prin corelație (nu înmulțim doar prin corelație)

Formula de la funcția costului pentru gradient:

$J(w) = -\frac{1}{n} \sum_{i=1}^n (y_i - \langle y \rangle)(d_i - \langle d \rangle)$ unde $\langle y \rangle$ și $\langle d \rangle$ sunt medii calculate, pentru toate perechile de antrenament, ptr. output și respectiv, etichete.

O altă funcție este noaptea erorii relative de entropie instantanee:

$J(w) = -\frac{1}{2} \left[(1+d) \ln \left(\frac{1+d}{1+y} \right) + (1-d) \ln \left(\frac{1-d}{1-y} \right) \right]$ unde $d \in (-1, 1)$. Deci $y = \tanh(\beta \text{net})$, atunci $\nabla J(w) = -\beta(d-y)x$.

Considerăm următoarea funcție costului generală ptr. gradientul descendent:

$J(w) = \sum_{i=1}^m g(z^T w)$ unde $z \in X$, deci $x \in \text{clasa } 1$, $g(x) = -x$, deci $x \in \text{clasa } 2$.

Pe $S = z^T w$, atunci $J(w)$ este bine formulat (nu ne grăbim garantat în w^*), deci (condițiile lui Karstner și Boker):

1. $g(s)$ este diferentabil;
2. $\text{ptr. } (v) s, -dg(s)/ds \geq 0$
3. $\exists \epsilon > 0$ și $-dg(s)/ds \geq \epsilon$, $\text{ptr. } (v) s \leq 0$
4. $g(s)$ este mărginit inferior