

3. Instruire prin minimizarea erorilor pătratice

Algoritmul perceptronului și metoda Gallant reprezintă proceduri iterative de instruire care pot fi aplicate doar în cazul unor ieșiri binare ale neuronilor. Ele sunt metode de *corectare a erorilor* deoarece modifică vectorul pondere numai în cazul unei clasificări incorecte. Acești algoritmi converg numai pentru clase liniar separabile. Pentru clase neseperabile, algoritmi devin oscilanți la obținerea soluției, nefiind posibilă determinarea unui vector de separare aproximativ. În acest capitol vom expune o altă metodă de instruire a rețelelor neuronale, care se va introduce mai întâi pentru o rețea cu un singur neuron.

3.1. METODĂ DE INSTRUIRE GLOBALĂ (ÎNTR-UN SINGUR PAS) A REȚELELOR NEURONALE

Considerăm o mulțime de instruire formată din vectorii normalizați în raport cu semnul $\mathbf{z}^1, \mathbf{z}^2, \dots, \mathbf{z}^k, \dots, \mathbf{z}^p$. Așa cum s-a arătat prin instruire se urmărește determinarea unui vector pondere \mathbf{w} , astfel încât

$$\mathbf{w}^T \mathbf{z}^k > 0 \quad (\forall) k = 1, 2, \dots, p. \quad (3.1)$$

În cele ce urmează se va prezenta o metodă de instruire globală, care determină toate ponderile simultan. Deoarece rezolvarea unui sistem de inecuații este mai complicată decât a unui sistem de ecuații, se înlocuiește problema rezolvării sistemului (3.1) cu rezolvarea unui sistem de ecuații “echivalent”. În acest scop se adoptă un vector \mathbf{b} cu p componente, numere pozitive convenabil alese,

$$\mathbf{b} = [b_1 \ b_2 \ \dots \ b_p]^T, \text{ cu } b_k > 0 \quad (\forall) k = \overline{1, p}, \quad (3.2)$$

și se înlocuiește inegalitatea

$$\mathbf{w}^T \mathbf{z}^k > 0, \quad (3.3)$$

cu egalitatea

$$\mathbf{w}^T \mathbf{z}^k = b_k. \quad (3.4)$$

Astfel, determinarea ponderilor rețelei, respectiv a vectorului \mathbf{w} se reduce la problema rezolvării sistemului liniar de ecuații

$$\mathbf{w}^T \mathbf{z}^k = b_k \quad (\forall) k = 1, 2, \dots, p. \quad (3.5)$$

Rezolvarea exactă a acestui sistem nu este posibilă în toate cazurile, pentru că, de multe ori, numărul p al ecuațiilor este mai mare decât numărul necunoscutelor egal cu $N + 1$ (se consideră că valoarea pragului este absorbită în mulțimea ponderilor pe mărimea de intrare x_0 cu valoarea 1).

Fie \mathbf{w} o soluție aproximativă a sistemului (3.5). Eroarea pătratică asociată vectorului de instruire \mathbf{z}^k este

$$e_k^2 = (\mathbf{w}^T \mathbf{z}^k - b_k)^2. \quad (3.6)$$

Se introduce acum funcția criteriu $J : \Re^{N+1} \rightarrow \Re$ definită ca suma erorilor pătratice a tuturor vectorilor de instruire

$$J = \sum_{k=1}^p (\mathbf{w}^T \mathbf{z}^k - b_k)^2 \quad (3.7)$$

Pentru descrierea matriceală a funcției criteriu introducem matricea \mathbf{Z} de dimensiune $p \times (N+1)$ care conține pe linia k vectorul transpus \mathbf{z}^k

$$\mathbf{Z} = [\mathbf{z}^{1T} \mathbf{z}^{2T} \dots \mathbf{z}^{kT} \dots \mathbf{z}^{pT}]^T. \quad (3.8)$$

Sistemul de ecuații (3.5) a cărui soluție aproximativă o căutăm se poate scrie acum sub forma

$$\mathbf{Z}\mathbf{w} = \mathbf{b}. \quad (3.9)$$

Pe de altă parte, având în vedere că $\mathbf{z}^T \mathbf{w} = \mathbf{w}^T \mathbf{z}$, vectorul $\mathbf{Z}\mathbf{w}$ se poate exprima prin

$$\mathbf{Z}\mathbf{w} = [\mathbf{w}^T \mathbf{z}^1 \mathbf{w}^T \mathbf{z}^2 \dots \mathbf{w}^T \mathbf{z}^p]^T, \quad (3.10)$$

deci

$$\sum_{k=1}^p (\mathbf{w}^T \mathbf{z}^k - b_k)^2 = \|\mathbf{Z}\mathbf{w} - \mathbf{b}\|^2 = (\mathbf{Z}\mathbf{w} - \mathbf{b})^T (\mathbf{Z}\mathbf{w} - \mathbf{b}). \quad (3.11)$$

Rezultă că funcția criteriu J se poate scrie și conform expresiei

$$J(\mathbf{w}) = \|\mathbf{Z}\mathbf{w} - \mathbf{b}\|^2 = (\mathbf{Z}\mathbf{w} - \mathbf{b})^T (\mathbf{Z}\mathbf{w} - \mathbf{b}). \quad (3.12)$$

Soluția aproximativă poate fi determinată prin minimizarea funcției criteriu, problemă care se rezolvă prin anularea derivatei în raport cu \mathbf{w} :

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{Z}^T (\mathbf{Z}\mathbf{w} - \mathbf{b}) = 0, \quad (3.13)$$

de unde rezultă:

$$\mathbf{Z}^T \mathbf{Z}\mathbf{w} = \mathbf{Z}^T \mathbf{b}. \quad (3.14)$$

Dacă matricea pătratică $\mathbf{Z}^T \mathbf{Z}$ este nesingulară, atunci sistemul precedent are soluția:

$$\mathbf{w}^* = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T \mathbf{b}. \quad (3.15)$$

Matricea $\mathbf{Z}^+ = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T$ se numește matrice inversă generalizată sau pseudo-inversă a matricei $\mathbf{Z}_{p \times (N+1)}$. Este evidentă proprietatea

$$\mathbf{Z}^+ \mathbf{Z} = \mathbf{I}, \quad (3.16)$$

unde \mathbf{I} este matricea unitate de dimensiune $(N+1) \times (N+1)$. Această proprietate arată că \mathbf{Z}^+ este inversa la stânga a matricei \mathbf{Z} . În particular dacă \mathbf{Z} este o matrice pătratică nesingulară, atunci avem

$$\mathbf{Z}^+ = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T = \mathbf{Z}^{-1} (\mathbf{Z}^T)^{-1} \mathbf{Z}^T = \mathbf{Z}^{-1}. \quad (3.17)$$

Așadar, dacă inversa unei matrice pătratică există, atunci inversa și inversa generalizată coincid. În general, nu orice matrice dreptunghiulară admite o inversă generalizată. Condiția de existență a inversei generalizate este dată de următoarea teoremă:

Fie \mathbf{A} o matrice de dimensiune $p \times q$ cu $p > q$. Matricea \mathbf{A} admite o inversă generalizată dacă și numai dacă

$$\text{rang } \mathbf{A} = q \quad (3.18)$$

Din cele prezentate rezultă că vectorul pondere soluție a problemei de instruire se calculează cu relația matriceală

$$\mathbf{w}^* = \mathbf{Z}^+ \mathbf{b}, \quad (3.19)$$

unde \mathbf{Z}^+ este inversa generalizată a matricei \mathbf{Z} .

Deoarece funcția criteriu J este convexă, minimul \mathbf{w}^* al acestei funcții este un minim global. Vectorul \mathbf{w}^* reprezintă o soluție de minimizare a erorii pătratice (MEP) pentru problema de instruire. Această soluție depinde în mod esențial de vectorul \mathbf{b} .

Observație: Chiar dacă clasele de instruire sunt liniar separabile, soluția MEP nu constituie neapărat un vector de separare în situația unei alegeri arbitrare a vectorului \mathbf{b} .

Conform teoremei de mai înainte, matricea \mathbf{Z} de dimensiune $p \times (N+1)$ admite o inversă generalizată la stânga dacă și numai dacă

$$\text{rang } \mathbf{Z} = (N+1). \quad (3.20)$$

Acest lucru are loc dacă cel puțin $(N+1)$ dintre vectorii de instruire \mathbf{z}^k , cu $k = \overline{1, p}$, $p > (N+1)$ sunt liniar independenți. Geometric aceasta înseamnă că există $(N+1)$ vectori de instruire care nu aparțin aceluiași hiperplan.

3.2. REȚEAUA ADALINE. PRINCIPIUL ALGORITMULUI WIDROW-HOFF DETERMINIST

Privind critic, metoda MEP are inconvenientul că necesită folosirea și inversarea unor matrice de dimensiuni mari, cu pericolul ca matricea $\mathbf{Z}^T \mathbf{Z}$ să fie singulară (deci neinvertibilă). Un alt dezavantaj de natură intrinsecă, mai puțin evident, provine din faptul că această metodă calculează ponderile într-un singur pas. Așadar ea nu permite ameliorarea valorilor ponderilor în pași succesivi, prin utilizarea de mai multe ori a fiecărui obiect din mulțimea de instruire.

Algoritmul Widrow-Hoff este un algoritm iterativ de instruire, care minimizează suma erorilor pătratice. Există mai multe variante de aplicare a acestui algoritm. Vom prezenta varianta de bază a algoritmului Widrow-Hoff aplicat unei rețele cu un singur neuron cu funcție de activare liniară numită de cei doi (Widrow și Hoff) rețea de tip ADALINE (ADaptive LInear NEuron). Pentru fixarea notațiilor în fig. 3.1 se prezintă sub formă grafică structura acestei rețele. De aici se observă că funcția de activare este liniară, fiind descrisă de relația $f(s) = s$, unde s este activarea totală care se calculează cu expresia vectorială

$s = \mathbf{w}^T \mathbf{x}$. Din acest motiv mărimea de ieșire are valorile cuprinse în mulțimea numerelor reale ($y \in \mathbb{R}$), spre deosebire de perceptron la care ieșirea poate lua numai valorile ± 1 sau 0, 1. Diferă, de asemenea, metoda de instruire care se impune ca o consecință a faptului menționat mai înainte.

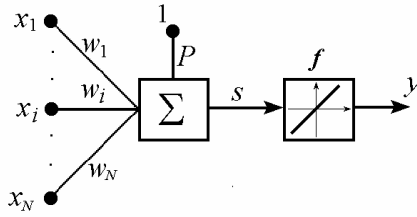


Fig. 3. 1 Structura unei rețele ADALINE.

Algoritmul Widrow-Hoff are la bază o metodă de optimizare de tip gradient care folosește ca funcție criteriu suma erorilor pătratice, considerând vectorii de intrare normalizați de semn. La metoda minimizării globale a erorii pătratice am obținut funcția criteriu:

$$J(\mathbf{w}) = (\mathbf{Z}\mathbf{w} - \mathbf{b})^T (\mathbf{Z}\mathbf{w} - \mathbf{b}), \quad (3.21)$$

în care:

$$\mathbf{Z} = [\mathbf{z}^{1T} \quad \mathbf{z}^{2T} \quad \dots \quad \mathbf{z}^{kT} \quad \dots \quad \mathbf{z}^{pT}]^T, \quad (3.22)$$

unde \mathbf{z}^k este vectorul normalizat din mulțimea de instruire care cuprinde, neglijând semnul, componentele $x_1^k \quad x_2^k \quad \dots \quad x_j^k \quad \dots \quad x_{N+1}^k$ ale mărimilor de intrare.

Vectorul pondere se calculează iterativ cu relația cunoscută

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \alpha_k \mathbf{p}^k, \quad (3.23)$$

unde

$$\mathbf{p}^k = \text{grad } J(\mathbf{w}) = \nabla J(\mathbf{w}) = \left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^k}. \quad (3.24)$$

La metoda MEP globală am stabilit că:

$$\mathbf{p}^k = \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{Z}^T (\mathbf{Z}\mathbf{w}^k - \mathbf{b}), \quad (3.25)$$

deci corecția la pasul k este de forma

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k \mathbf{Z}^T (\mathbf{b} - \mathbf{Z}\mathbf{w}^k), \quad (3.26)$$

unde $c_k = 2\alpha_k$.

În continuare vom modifica această regulă de instruire, considerând o ipoteză, general valabilă la instruirea rețelelor neuronale, potrivit căreia la pasul k se modifică numai un vector și anume \mathbf{z}^k . În acest sens matricea \mathbf{Z} se modifică adoptându-se conform expresiei

$$\mathbf{Z} = [0 \quad 0 \dots \mathbf{z}^{kT} \dots 0]^T, \quad (3.27)$$

în care toate liniile au componentele nule mai puțin linia k ce conține vectorul linie \mathbf{z}^{kT} . Altfel spus matricea \mathbf{Z} se înlocuiește în regula de corecție cu \mathbf{z}^{kT} (vector linie). De asemenea, în relația de calcul recurent se înlocuiește vectorul \mathbf{b} cu componenta b_k , obținându-se astfel expresia:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k \mathbf{z}^k [b_k - \mathbf{z}^{kT} \mathbf{w}^k], \quad (3.28)$$

care se poate rescrie sub forma

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k [b_k - (\mathbf{w}^k)^T \mathbf{z}^k] \mathbf{z}^k. \quad (3.29)$$

Metoda de căutare recurentă astfel obținută este cunoscută drept *regula de corecție Widrow-Hoff* sau *regula delta*. Numărul b_k poate fi interpretat ca fiind valoarea dorită a ieșirii. În acest caz valoarea corecției rezultă proporțională cu diferența (delta) dintre ieșirea dorită și ieșirea reală a rețelei:

$$\mathbf{w}^{k+1} - \mathbf{w}^k = c_k \delta_k \mathbf{z}^k, \quad (3.30)$$

unde

$$\delta_k = d^k - (\mathbf{w}^k)^T \mathbf{z}^k. \quad (3.31)$$

Pentru convergența metodei este necesar ca c_k să descrească cu creșterea lui k . O alegere convenabilă pentru c_k este

$$c_k = c_1 / k, \quad (3.32)$$

unde $c_1 > 0$.

Observație: În prezența unor condiții normale impuse coeficienților c_k procedeul iterativ Widrow-Hoff converge spre soluția de minim global a metodei MEP.

3. 3 ABORDAREA PROBABILISTĂ A PROBLEMEI DE MINIMIZARE A ERORII PĂTRATICE (MEP)

Formularea problemei. Funcția criteriu

Pentru a dezvolta o abordare statistică a instruirii vom considera vectorul mărimilor de intrare, la pasul k , o variabilă aleatoare vectorială $(N + 1)$ - dimensională, notată \mathbf{x}^k . Familia de variabile aleatoare

$$\{\mathbf{x}^k \mid k = 1, 2, \dots, p\} \quad (3.33)$$

constituie un proces aleator. În mod analog ieșirile, reală y^k și dorită d^k ale rețelei vor fi considerate drept variabile aleatoare. Așadar secvențele $\{d^k\}$ și $\{y^k\}$ definesc procese aleatoare sau secvențe de variabile aleatoare. Eroarea pătratică instantanee e_k^2 reprezintă de asemenea o variabilă aleatoare a cărei valoare medie statistică este

$$E\{e_k^2\} = E\{(d^k - y^k)^2\}, \quad (3.34)$$

unde $E\{\cdot\}$ reprezintă operatorul de mediere statistică. În cele ce urmează se va presupune că vectorul ponderilor \mathbf{w}^k este determinist. Din acest motiv vom avea proprietatea specifică operatorului de mediere statistică

$$E\{\mathbf{w}^{kT} \mathbf{x}^k\} = \mathbf{w}^{kT} E\{\mathbf{x}^k\}. \quad (3.35)$$

Înlocuind $y^k = \mathbf{w}^{kT} \mathbf{x}^k$ în expresia valorii medii a erorii pătratice și având în vedere proprietățile cunoscute ale operatorului de mediere statistică putem scrie

$$E\{e_k^2\} = E\{(d^k - \mathbf{w}^{kT} \mathbf{x}^k)^2\} = E\{(d^k)^2\} + E\{\mathbf{w}^{kT} \mathbf{x}^k\}^2 - 2E\{d^k \mathbf{w}^{kT} \mathbf{x}^k\} \quad (3.36)$$

În continuare, deoarece \mathbf{w}^k este un vector determinist vom avea

$$E\{\mathbf{w}^{kT} \mathbf{x}^k\}^2 = E\{\mathbf{w}^{kT} \mathbf{x}^k (\mathbf{w}^{kT} \mathbf{x}^k)\} = \mathbf{w}^{kT} E\{\mathbf{x}^k \mathbf{x}^{kT}\} \mathbf{w}^k \quad (3.37)$$

și

$$E\{d^k \mathbf{w}^{kT} \mathbf{x}^k\} = E\{d^k \mathbf{x}^{kT} \mathbf{w}^k\} = E\{d^k \mathbf{x}^{kT}\} \mathbf{w}^k. \quad (3.38)$$

În ultimele două relații s-a utilizat proprietatea produsului scalar $\mathbf{w}^{kT} \mathbf{x}^k = \mathbf{x}^{kT} \mathbf{w}^k$.

În continuare să analizăm mărimile $E\{\mathbf{x}^k \mathbf{x}^{kT}\}$ și $E\{d^k \mathbf{x}^{kT}\}$. Pentru simplificarea prezentării, renunțăm temporar la scrierea indicelui superior k . În aceste condiții $E\{\mathbf{x}^k \mathbf{x}^{kT}\}$ va fi o matrice notată:

$$\mathbf{R} = E\{\mathbf{x} \cdot \mathbf{x}^T\} = \begin{bmatrix} E\{x_1^2\} & E\{x_1 x_2\} & \dots & E\{x_1 x_{N+1}\} \\ \dots & \dots & \dots & \dots \\ E\{x_{N+1} x_1\} & \dots & \dots & E\{x_{N+1}^2\} \end{bmatrix} \quad (3.39)$$

de dimensiune $(N+1) \times (N+1)$, similară cu matricea de autocorelație a vectorului mărimilor de intrare, iar $E\{d^k \mathbf{x}^{kT}\}$ un vector

$$\mathbf{q} = E\{d \mathbf{x}\} = [E\{d_1 x_1\} \ E\{d_2 x_2\} \dots E\{d_{N+1} x_{N+1}\}]^T, \quad (3.40)$$

identic cu vectorul de intercorelație ieșire-intrare.

Cu notațiile astfel introduse, valoarea medie a erorii pătratice la momentul k se va putea scrie sub forma:

$$E\{e_k^2\} = E\{(d^k)^2\} + \mathbf{w}^{kT} \mathbf{R}^k \mathbf{w}^k - 2\mathbf{q}^{kT} \mathbf{w}^k, \quad (3.41)$$

sau dacă suprimăm în continuare, pentru simplitate, indicele k

$$J(\mathbf{w}) = E\{e^2\} = E\{d^2\} + \mathbf{w}^T \mathbf{R} \mathbf{w} - 2\mathbf{w}^T \mathbf{q}. \quad (3.42)$$

Soluția MEP probabilistă

Vectorul ponderilor sinaptice este o soluție a problemei de minim

$$\begin{cases} J(\mathbf{w}) \rightarrow \min \\ \mathbf{w} \in \Re^{N+1} \end{cases} \quad (3.43)$$

Punctul de minim al funcției criteriu rezultă rezolvând ecuația

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 0. \quad (3.44)$$

Având în vedere relația de derivare a unei forme pătratice și faptul că matricea de autocorelație \mathbf{R} este simetrică rezultă ecuația echivalentă:

$$2\mathbf{R}\mathbf{w} - 2\mathbf{q} = 0. \quad (3.45)$$

Dacă matricea \mathbf{R} este inversabilă se obține soluția analitică a problemei de instruire.

$$\mathbf{w}^* = \mathbf{R}^{-1}\mathbf{q}. \quad (3.46)$$

În plus, pentru ca extremul obținut să fie un punct de minim trebuie să fie îndeplinită condiția

$$\frac{\partial^2 J(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} > 0, \quad (3.47)$$

respectiv $\mathbf{R} > 0$.

Uneori în literatura de specialitate soluția MEP probabilistă se exprimă în funcție de pseudo-inversa matricei \mathbf{R} . Acest rezultat se obține simplu din condiția de extrem

$$\mathbf{R}\mathbf{w} = \mathbf{q}, \quad (3.48)$$

înmulțind la stânga cu \mathbf{R}^T , de unde rezultă

$$\mathbf{R}^T \mathbf{R} \mathbf{w} = \mathbf{R}^T \mathbf{q}. \quad (3.49)$$

Dacă matricea $\mathbf{R}^T \mathbf{R}$ este nesingulară atunci există matricea

$$\mathbf{R}^+ = (\mathbf{R}^T \mathbf{R})^{-1} \mathbf{R}^T \quad (3.50)$$

care este inversa generalizată a matricei \mathbf{R} . Soluția MEP probabilistă se scrie acum sub forma

$$\mathbf{w}^* = \mathbf{R}^+ \mathbf{q}. \quad (3.51)$$

Se constată că această soluție a problemei MEP probabilistă este similară formal cu soluția MEP deterministă. Remarcăm, totuși, că matricea \mathbf{R} care apare în această expresie este pătratică, pe când matricea \mathbf{Z} care apare în soluția MEP deterministă este dreptunghiulară de dimensiune $p \times (N + 1)$.

Procedeu iterativ de determinare pe bază statistică a vectorului pondere optim

Metoda MEP probabilistă prezentată mai înainte are dezavantajul major că necesită cunoașterea matricei de autocorelație statistică \mathbf{R} a vectorului de intrare și a vectorului de intercorelație statistică ieșire-intrare \mathbf{q} . Calculul explicit al lui \mathbf{R} și \mathbf{q} impune cunoașterea statistică a semnalelor de intrare și ieșire, adică a funcțiilor de densitate de probabilitate ale acestora la fiecare moment k . O altă dificultate este legată de calculul matricei \mathbf{R}^{-1} și de resursele de memorie pe care le implică acest calcul. O soluție alternativă o constituie folosirea unei metode iterative de optimizare de tip gradient. În acest caz vectorul pondere se deplasează pe suprafața funcției criteriu J , pe direcția celei mai mari pante.

Procedura iterativă se va scrie și în acest caz conform relației:

$$\mathbf{w}^{k+1} = \mathbf{w}^k - \alpha_k \mathbf{p}^k, \quad (3.52)$$

unde direcția de coborâre este gradientul funcției criteriu în punctul \mathbf{w}^k :

$$\mathbf{p}^k = \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} = 2(\mathbf{R}^k \mathbf{w}^k - \mathbf{q}^k). \quad (3.53)$$

Rezultă că:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k (\mathbf{R}^k \mathbf{w}^k - \mathbf{q}^k), \quad (3.54)$$

unde $c_k = 2\alpha_k$.

Determinarea direcției de coborâre ridică aproape aceleași probleme tehnice ca și determinarea, într-un singur pas, a soluției MEP probabiliste. Ar trebui să cunoaștem dinainte \mathbf{R} și \mathbf{q} , ceea ce înseamnă să știm a priori suprafața corespunzătoare funcției criteriu. Pentru a depăși aceste dificultăți vom folosi o versiune simplificată a funcției criteriu. În acest scop se înlocuiește valoarea medie statistică a erorii pătratice cu valoarea instantanee a acesteia, deci:

$$J(\mathbf{w}) = e_k^2. \quad (3.55)$$

Direcția de coborâre este în acest caz

$$\mathbf{p}^k = \left. \frac{\partial (e_k^2)}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^k}, \quad (3.56)$$

unde

$$e_k^2 = [d^k - \mathbf{w}^{kT} \mathbf{x}^k]^2. \quad (3.57)$$

Efectuând derivata funcției pătratice, rezultă

$$\mathbf{p}^k = -2e_k \mathbf{x}^k, \quad (3.58)$$

relație care se obține avându-se în vedere proprietatea de derivare a produsului scalar

$$\frac{\partial}{\partial \mathbf{w}} (\mathbf{w}^T \mathbf{x}) = \frac{\partial}{\partial \mathbf{w}} (\mathbf{x}^T \mathbf{w}) = \mathbf{x}. \quad (3.59)$$

Am obținut astfel următoarea aproximare:

$$\left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^k} = \frac{\partial}{\partial \mathbf{w}} E(e_k^2) \approx -2e_k \mathbf{x}^k. \quad (3.60)$$

Regula de corecție se poate scrie acum:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k e_k \mathbf{x}^k, \quad (3.61)$$

sau exprimând în mod explicit eroarea:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k [d^k - \mathbf{w}^{kT} \mathbf{x}^k] \mathbf{x}^k. \quad (3.62)$$

S-a regăsit, aşadar, procedeul iterativ de instruire Widrow-Hoff determinist sau regula delta. În legătură cu această formă a algoritmului sunt necesare câteva comentarii:

- i) vectorii \mathbf{x}^k nu au suferit vreo normalizare de semn;
- ii) valorile ieşirii d^k pot fi pozitive sau negative spre deosebire de b_k de la algoritmul determinist care era în mod obligatoriu pozitiv;
- iii) dacă ne situăm în cazul paradigmei clasificatorului, din regula de mai sus se poate deduce algoritmul Widrow-Hoff determinist:

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k [b_k - \mathbf{w}^{kT} \mathbf{z}^k] \mathbf{z}^k, \text{ unde } \mathbf{z}^k \text{ sunt vectori normalizaţi şi } b_k > 0;$$
- iv) sarcinile pe care le poate îndeplini o reţea ADALINE (neuron instruit cu algoritmul Widrow-Hoff) sunt mai generale decât stabilirea apartenenţei unui obiect la o clasă (din două posibile).

3. 4 ALGORITMUL WIDROW-HOFF

În metoda Widrow-Hoff, la fiecare pas al instruirii se prezintă reţelei doar un singur vector de intrare. Mai întâi se consideră cazul unei rate de învăţare constante, deci

$$c_k = c, \quad k = 1, 2, \dots$$

Descrierea algoritmului

P₀ Se iniţializează ponderile reţelei alegând un vector \mathbf{w}^1 arbitrar. Se fixează rata de învăţare $c > 0$. Se stabileşte numărul vectorilor de instruire, p . Se pune $k := 1$ şi $e^2 := 0$ (eroarea totală).

P₁ Se prezintă reţelei un vector de intrare \mathbf{x}^k şi ieşirea corectă (dorită) d^k .

P₂ Se calculează valoarea instantanee a erorii pătratice e_k^2

$$e_k^2 = [d^k - \mathbf{w}^{kT} \mathbf{x}^k]^2 \text{ şi eroarea totală, } e^2 := e^2 + e_k^2.$$

P₃ Se adaptează vectorul pondere folosind regula delta probabilistă

$$\mathbf{w}^{k+1} = \mathbf{w}^k + c_k [d^k - \mathbf{w}^{kT} \mathbf{x}^k] \mathbf{x}^k.$$

P₄ Se pune $k := k + 1$. Pentru $k \leq p$ se merge la pasul **P₁**, altfel la pasul **P₅**.

P₅ Se verifică dacă $e^2 \leq \varepsilon$ (de exemplu $\varepsilon = 10^{-5}$). Dacă da STOP, altfel se merge la pasul **P₁**.

Alegerea ratei de instruire c are o importanţă deosebită în cadrul algoritmului iterativ Widrow-Hoff. Această mărime controlează convergenţa spre vectorul pondere obţinut prin metoda MEP. Dacă c are o valoare prea mică atunci algoritmul converge foarte lent. Dacă c este prea mare, atunci algoritmul poate să nu sesizeze punctele de minim. În acest caz este posibil ca procesul de instruire să nu conveargă niciodată. Rata de învăţare se poate alege o constantă, sau variabilă, caz în care se calculează cu relaţia $c_k = 1/k$ sau $c_k = c/k$, $c < 0$, $k = 1, 2, \dots$ unde k este indexul (pasul curent) al algoritmului de instruire. Alegerea constantei c se face pe baza experienţei sau riguros folosindu-se unele noţiuni din teoria clasificării care nu vor fi expuse aici.

Algoritmul Widrow-Hoff versus algoritmul perceptronului

În această secțiune se va face un mic studiu comparativ al celor mai importanți algoritmi de instruire prezentați până acum: algoritmul perceptronului și procedura Widrow-Hoff.

- i) Algoritmul Widrow-Hoff a fost dezvoltat pentru a găsi o metodă de instruire pentru date neseparabile, ca o alternativă la algoritmul perceptronului;
- ii) O altă motivație a fost de a găsi, pentru clase separabile, o soluție mai bună decât cele furnizate de algoritmul perceptronului. Instruirea cu algoritmul perceptronului se oprește de îndată ce se obține o soluție care clasifică în mod corect fiecare formă de instruire. În cele mai multe cazuri, soluția găsită este abia acceptabilă (vezi figura 3. 2a). În figura 3. 2b se prezintă o soluție mai bună ce se obține cu metoda Widrow –Hoff;

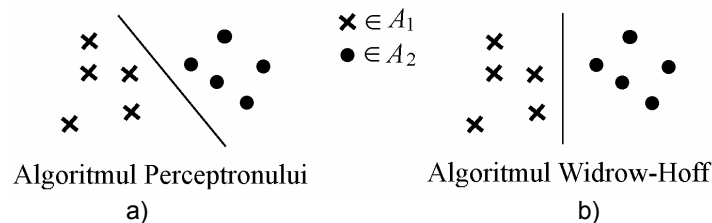


Fig. 3. 2 Comparația soluțiilor de clasificare realizate cu metoda perceptronului respectiv cu procedura W-H.

- iii) Algoritmul Widrow-Hoff asigură față de algoritmul perceptronului o viteză de instruire mai mare. Acest lucru se realizează prin ajustarea ratei de instruire la fiecare iterație a procedurii.
- iv) La metoda de instruire cu algoritmul Widrow-Hoff, ieșirile dorite d_k pot fi numere reale, nefiind limitate la valorile ± 1 ca la perceptron.
- v) Metoda MEP globală și algoritmul Widrow-Hoff furnizează o soluție atât pentru clase separabile cât și pentru clase neseparabile. Prezența însă a unui singur punct izolat poate face ca hiperplanul de separare să fie complet perturbat, astfel încât majoritatea datelor de instruire să fie eronat clasificate. Acest lucru reprezintă un dezavantaj serios deoarece punctul izolat poate fi un exemplu de instruire fals datorat zgomotului.

Nu există la această oră o comparație realizată în mod satisfăcător între algoritmul Widrow-Hoff și algoritmul buzunarului (Gallant).