

PROCESAREA CUNOȘTINȚELOR ȘI CALCUL INTELIGENT - PCCI

Rețele neuronale artificiale pentru învățare
supravegheată. Metode de îmbunătățire a învățării

Cuprins:

- Problema supraantrenării
- Metode de îmbunătățire a învățării
 - ❖ Optimizarea datelor de intrare
 - ❖ Optimizarea funcției eroare
 - ❖ Optimizarea ponderilor
 - ❖ Algoritmul Resilient Backpropagation – RProp
 - ❖ Algoritme de învățare alternative

Recapitulare

- Neuronul elementar:

$$y = b + w_1 \cdot x_1 + \dots + w_i \cdot x_i + \dots + w_n \cdot x_n =$$
$$= w_0 \cdot 1 + w_1 \cdot x_1 + \dots + w_i \cdot x_i + \dots + w_n \cdot x_n = \sum_{i=0}^n w_i \cdot x_i, x_0 = 1$$

- Algoritmul de retropropagare după metoda gradientului
- Eroarea:

$$E(w, v) = \frac{1}{2} \cdot \sum_{m=1}^M \left(d^{(m)} - o(w, v, x)^{(m)} \right)^2$$

- Ajustarea ponderilor:

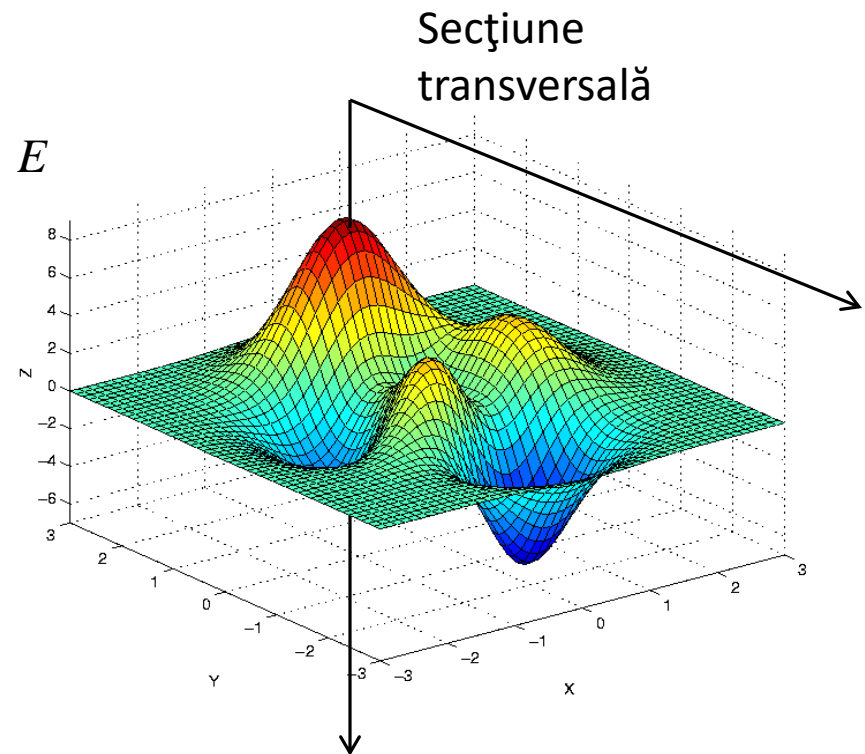
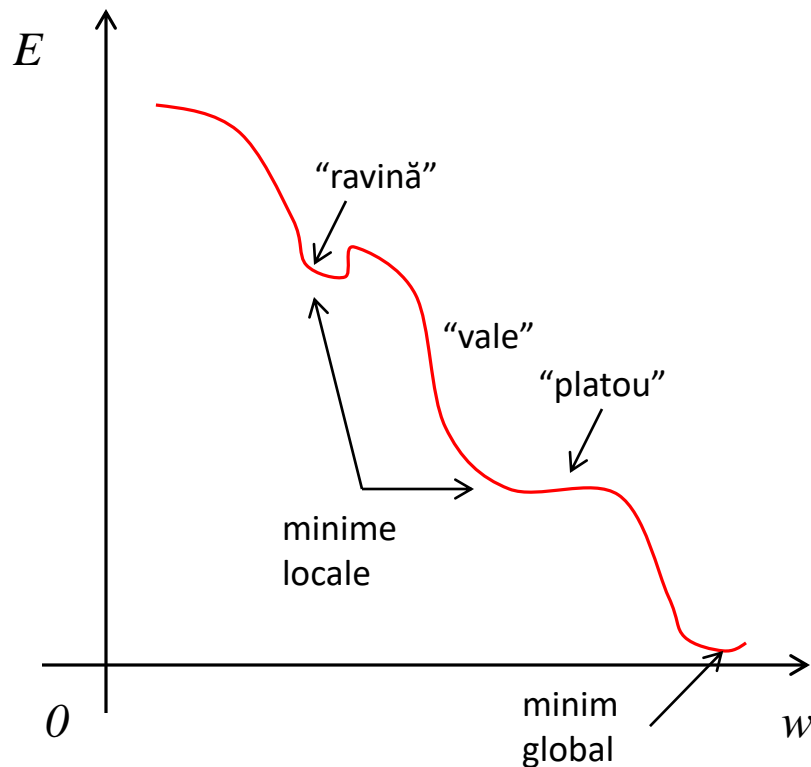
$$w_{ij}^{(t+1)} = w_{ij}^{(t)} - \eta \cdot \frac{\partial E}{\partial w_{ij}^{(t)}} \Big|_{w=w_{ij}^{(t)}}, \quad i = 1 \dots I, j = 1 \dots J$$

Recapitulare

- Algoritmul de retropropagare:
 - ❖ Definire problemă și pregătire date de intrare.
 - ❖ Definire arhitectură RNA.
 - ❖ Inițializează ponderi și praguri RNA.
 - ❖ Repetă până la îndeplinirea unui criteriu de oprire
 - Propagare înainte (calcul erori)
 - Propagare înapoi (corectare ponderi și praguri în funcție de eroare).

Recapitulare

- Folosirea unor modele de intrare cu mai mulți parametri și suprapunerea efectelor generate de mai mulți neuroni interconectați fac suprafața erorii non-convexă, cu “dealuri” și “văi”, cu “platouri” locale și “ravine”, care pot încetini sau bloca atingerea obiectivului antrenării, a minimului global.



Antrenare și generalizare

- În cazul unei antrenări ideale, eroarea scade repede, valorile de ieșire approximate de RNA se apropie suficient de aproape de valorile de ieșire și RNA este capabilă să **generalizeze**, adică să calculeze ieșiri rezonabil corecte pentru date de intrare noi.
- Scopul antrenării RNA este generalizarea.
- Soluțiile date de RNA nu trebuie să fie perfecte, cu rezonabil de aproape de cele exacte, cu eroare cât mai mică.

Posibile probleme de antrenare și generalizare:

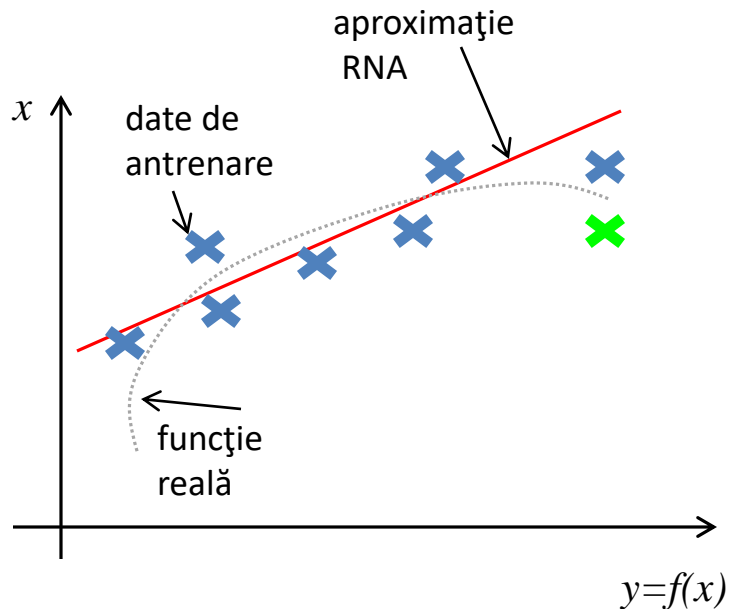
- Antrenarea are eroare mică, dar RNA nu generalizează bine, nu calculează corect ieșirea pentru modele de intrare noi.
 - ❖ RNA se află în unul din două cazuri: subantrenare (underfitting) sau supraantrenare (overfitting).
 - Datele de intrare: sunt suficiente și relevante ?
 - Arhitectura rețelei (număr de straturi și de neuroni) este adecvată ?
 - Funcția de aproximare modelată de RNA este adecvată ?
- Învățare insuficientă - eroarea de antrenare rămâne mare (blocaj în minim local)
 - ❖ Adaptarea ponderilor nu se desfășoară corespunzător sau eroarea se blochează în minime locale sau “ravine”.
 - Ponderile și pragurile au valori rezonabile ?
 - Rata de învățare are o valoare adecvată ?
- Învățare lentă
 - ❖ Eroarea nu trece suficient de repede peste minimele locale, “platouri” sau “ravine”.
 - ❖ Corecțiile ponderilor sunt prea mici.
 - Prelucrarea suplimentară a datelor de antrenare (uniformizare)
 - Factori de corecție suplimentari
 - Algoritme mai bune de adaptare a ponderilor și pragurilor.

Subantrenarea

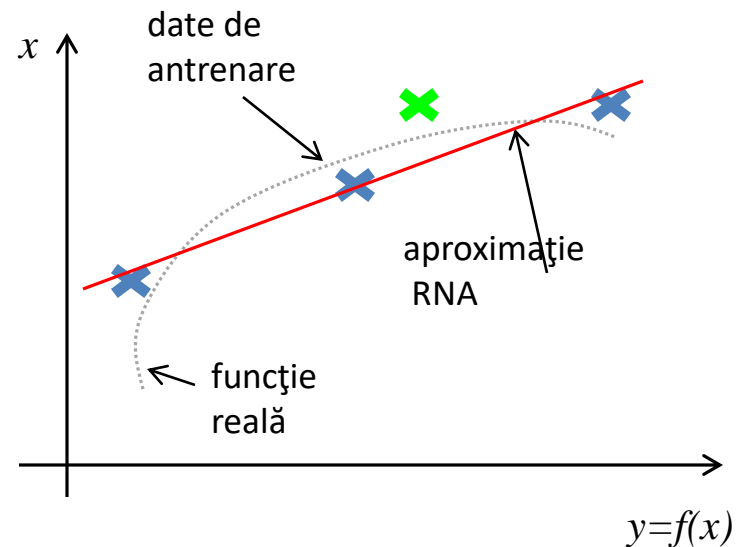
Subantrenare: funcția cu care RNA aproximează ieșirea dorită este prea simplă

❖ RNA nu va generaliza corect pentru punctele de tipul ✕

– Curbă prea netedă față de necesar



– RNA are la dispoziție prea puține date, care nu reflectă variația reală a funcției

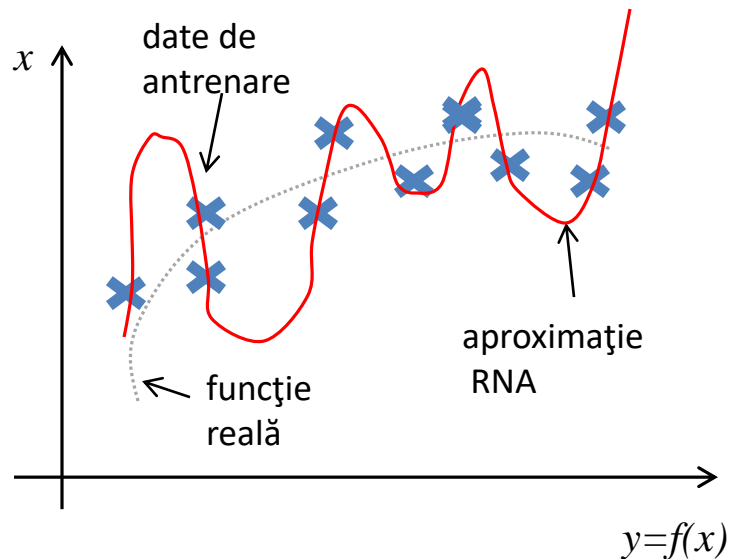


Supraantrenarea

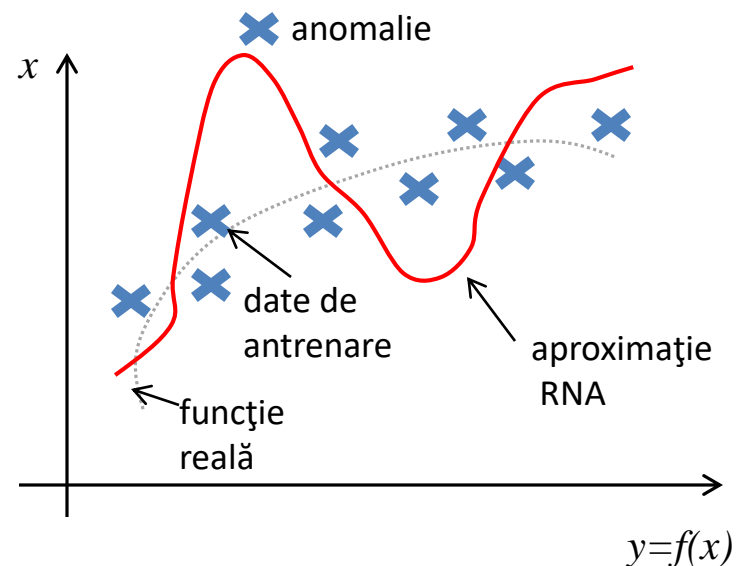
Subantrenare: funcția cu care RNA aproximează ieșirea dorită este prea complexă.

- ❖ RNA nu va generaliza corect punctele neincluse în setul de antrenare

– Curbă prea întortocheată față de necesar



– Prezența unor anomalii în datele de intrare

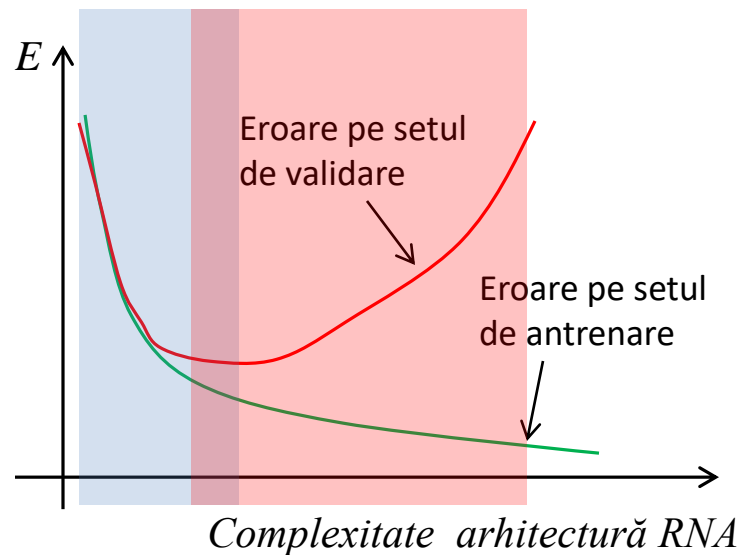


Soluții

- Pentru evitarea sau prevenirea subantrenării:
 - ❖ Antrenarea unor RNA cu arhitectură complexă (mulți neuroni ascunși, multe straturi ascunse) - permite aproximarea unor funcții complexe
 - ❖ Folosirea unui număr mare de modele de antrenare – pentru a reprezenta mai fidel funcția căutată
 - ❖ Creșterea numărului de caracteristici al unui model de intrare – crește numărul de intrări al RNA și numărul de conexiuni între neuroni (ponderi), permițând definirea unor funcții de aproximat complexe
- Pentru evitarea supraantrenării:
 - ❖ Antrenarea unor RNA cu arhitectură simplă (puțini neuroni ascunși, un strat ascuns) - reduce complexitatea funcției ce poate fi aproximată
 - ❖ Eliminarea anomaliilor din datele de antrenare, prin analiză vizuală sau statistică.
 - ❖ Folosirea unor modele de intrare cu un număr mic de caracteristici - reduce numărul de ponderi.

Detectarea supraantrenării

- Datele disponibile pentru setul de antrenare se împart în trei:
 - ❖ Setul de antrenare (60%) – pentru aplicarea algoritmului de învățare care realizează adaptarea ponderilor.
 - ❖ Setul de validare (20%) - pentru verificarea progresului antrenării
 - ❖ Setul test (20%) – pentru verificarea generalizării



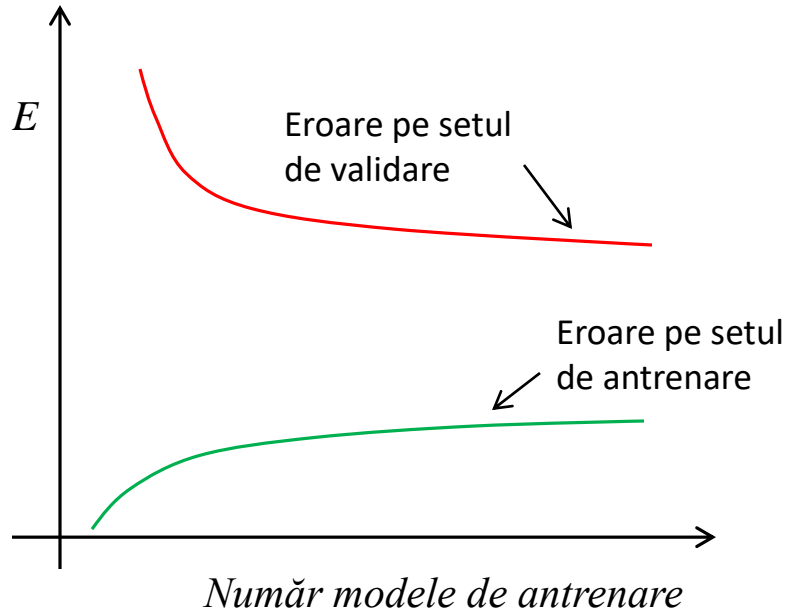
- Eroarea pe setul de antrenare scade, deoarece pe aceste date se corectează ponderile. RNA aproximează cel mai bine aceste puncte.
- Eroarea pe setul de validare crește, deoarece ponderile (aproximarea RNA) nu se potrivește și pe datele noi din setul de validare, deci generalizarea nu este asigurată
- **Antrenarea se oprește când eroarea pe setul de validare începe să crească.**

subantrenare

supraantrenare

Detectarea supraantrenării

- Câte modele de antrenare să folosesc ? Când sunt prea multe ?
- Supraantrenare:

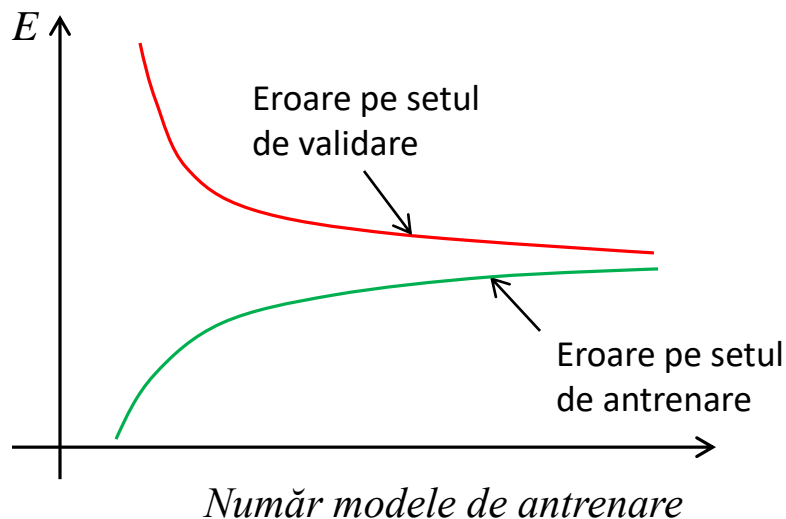


- Eroarea pe setul de antrenare crește, deoarece crește numărul de modele și erorile se cumulează.
- Eroarea pe setul de validare scade, dar nu se apropie de cea obținută pe setul de antrenare, deoarece funcția de aproximare a RNA nu se potrivește și pentru modelele din setul de validare .

[Andrew NG, scikits]

Detectarea subantrenării

- Câte modele de antrenare să folosesc ? Când sunt prea puține ?
- Subantrenare:



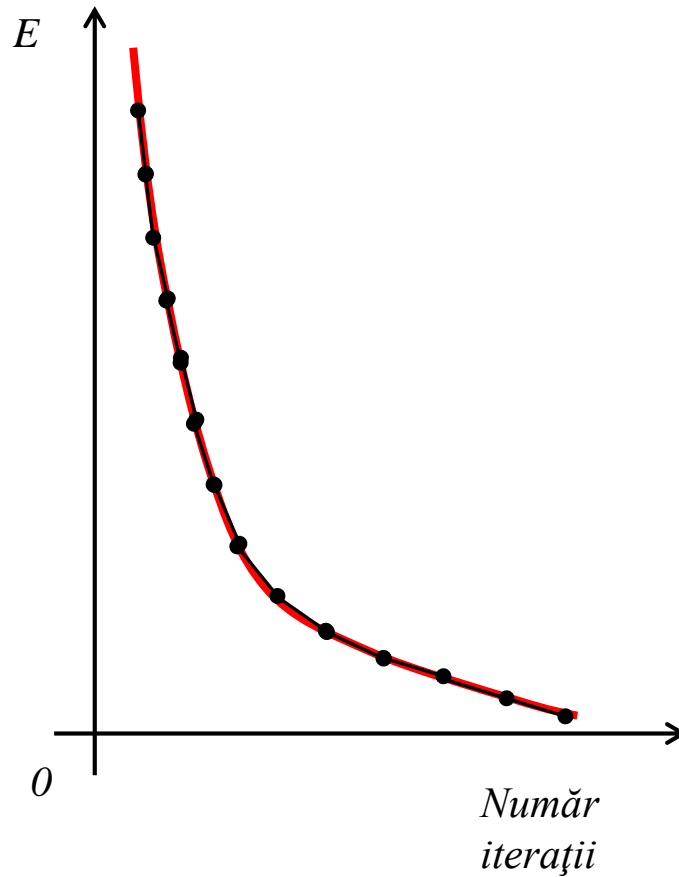
- Eroarea pe setul de antrenare crește, deoarece crește numărul de modele și erorile se cumulează
- Eroarea pe setul de validare scade, deoarece aproximarea RNA e mai bună cu mai multe modele, și se apropie de cea a setului de antrenare
- Curbele se întâlnesc la o valoare mare a erorii

Influența ratei de învățare asupra antrenării

$$w_{ij}^{(t+1)} = w_{ij}^{(t)} - \eta \cdot \frac{\partial E}{\partial w_{ij}^{(t)}} \Big|_{w=w_{ij}^{(t)}}, \quad i = 1 \dots I, j = 1 \dots J$$

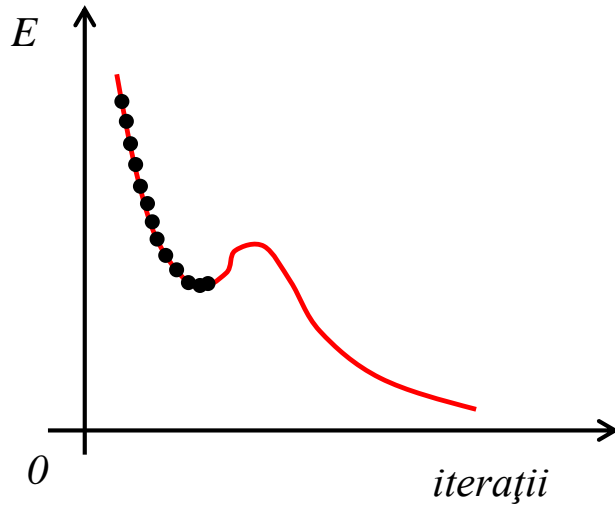
- O rată de învățare prea mică înseamnă corecții mici pentru ponderi
 - ❖ Antrenarea poate dura mult sau se poate bloca într-o “ravenă” sau pe un “platou” local.
- O rată de învățare prea mare înseamnă corecții mai mari pentru ponderi, care pot cauza ratarea minimului global și, în caz extrem, divergență.
 - ❖ Antrenarea nu găsește minimul global, ceea ce probabil va influența negativ și generalizarea
- Valori recomandate: $\eta=0.1 \dots 0.2$
- Nu există o rețetă garantată de succes, trebuie făcute teste

Rată de învățare prea mică = învățare lentă

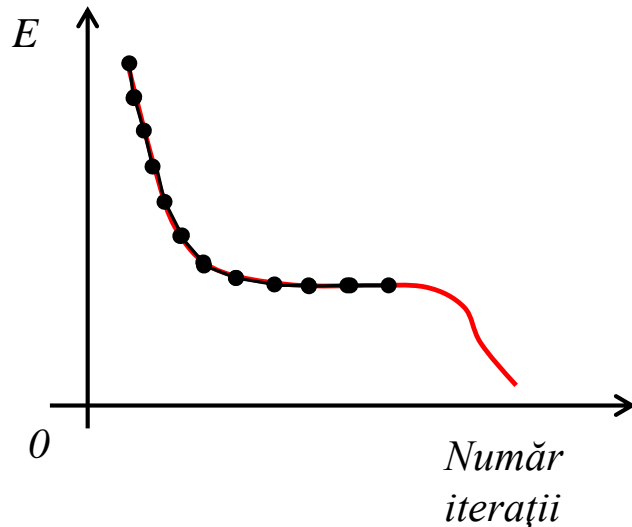


- Minimul global al erorii se atinge după multe iterații
- Caz fericit, când suprafața erorii e convexă
- Se poate ajunge la zeci de mii, chiar sute de mii de iterații

Rată de învățare prea mică = blocaj

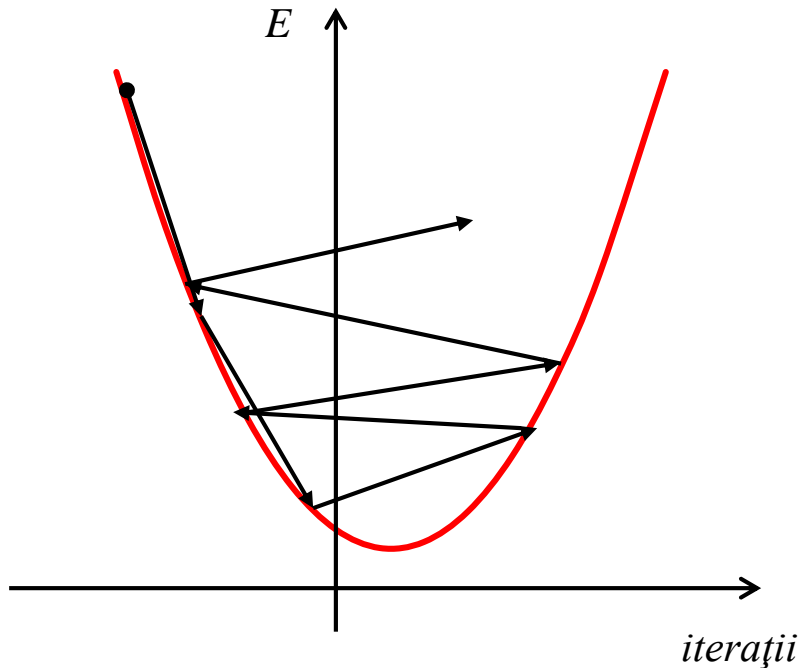


- Gradientul nu e suficient de mare pentru a ieși din “ravenă”.



- Dacă criteriul de oprire al antrenării e un prag de variație a erorii pe parcursul unui număr mic de iterații succesive, antrenarea se poate opri înainte de ieșirea din “platou”.

Rată de învățare prea mare = divergență



- Când gradientul este prea mare, corecția erorii este prea mare
- Antrenarea poate intra în divergență

Antrenarea cu rate de învățare adaptive

- Reducerea progresivă a ratei de învățare de la o iterație la alta:

$$\eta^{(t+1)} = p \cdot \eta^{(t)}, p = [0.1...1]$$

- ❖ În prima parte a antrenării, rata de învățare mare permite apropierea mai rapidă de minimul global, sărind peste “ravenele” mici și minimele globale
 - ❖ La apropierea de minimul global, converge către acesta, reducând riscul depășirii lui.
-
- Cu cât corecția este mai mare (p este mai mic), cu atât mă apropii de cazul η prea mic.
 - Cu cât corecția este mai mică (p este mai mare), cu atât mă apropii de cazul η prea mare.

Antrenarea cu rate de învățare adaptive

- Viteza și acuratețea învățării pot fi îmbunătățite dacă fiecărei ponderi i se atașează o rată de învățare distinctă, adaptată diferit, după următorul principiu [Hinton]:
 - ❖ Valoarea inițială a ponderilor este 1
 - ❖ De la o iterație la alta, se calculează variația (gradientul) ponderii
 - ❖ Dacă, în două iterații succesive, gradientul își păstrează semnul (eroarea scade și apoi iar scade ori crește și iar crește), rata de învățare este amplificată – sunt pe drumul cel bun.
 - ❖ Dacă, în două iterații succesive, gradientul își schimbă semnul (eroarea scade, apoi crește ori crește, apoi scade, adică oscilează), rata de învățare este diminuată – sunt pe drumul greșit.

$$w_{ij}^{(t+1)} = w_{ij}^{(t)} - \eta \cdot \text{corr}_{ij} \frac{\partial E}{\partial w_{ij}^{(t)}}$$

$$\text{test} = \left(\frac{\partial E}{\partial w_{ij}} \right)^{(t+1)} \cdot \left(\frac{\partial E}{\partial w_{ij}} \right)^{(t)}$$

$$\text{test} > 0, \quad \text{corr}_{ij}^{(t+1)} = \text{corr}_{ij}^{(t)} + 0.05$$

$$\text{test} \leq 0, \quad \text{corr}_{ij}^{(t+1)} = \text{corr}_{ij}^{(t)} \cdot 0.95$$

$$\text{corr} = [0..100]$$

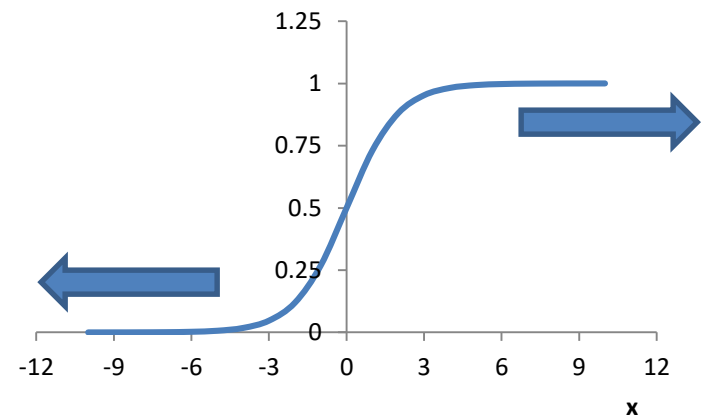
Inițializarea ponderilor

- Ponderile nu pot fi inițializate cu 0, deoarece toate ieșirile calculate ar fi 0, deci eroarea și corecțiile ar fi aceleași la fiecare iterație, pentru toate ponderile.
- Există numeroase metode dezvoltate pe bază empirică, care au dat rezultate bune:
 - ❖ Inițializare aleatorie cu valori mici, în intervalul (-1, 1)
 - ❖ Regula Russo:
$$-\frac{2.4}{I} \leq w_{ij} \leq \frac{2.4}{I}$$

I - numărul de intrări ale RNA
 - ❖ ...și altele

Optimizarea antrenării

- Pe parcursul antrenării, prin adăugarea succesivă a corecțiilor Δw , ponderile tind să aibă valori foarte mari (pozitive) sau foarte mici (negative).
- Dacă se folosește funcția de activare sigmoid logistic, valori depărtate de 0 ale intrării nete a neuronului saturează ieșirea sigmoidului către 0 sau 1, ceea ce scade performanțele antrenării
- Este de dorit menținerea intrării nete în zona de variație amplă a sigmoidului, pentru diferențierea fină în cadrul claselor identificate.
 - Se pot modifica ponderile w
 - Se pot modifica intrările x



$$net = w_1 \cdot x_1 + \dots + w_i \cdot x_i + \dots + w_n \cdot x_n \quad y = \frac{1}{1 + e^{-net}}$$

Optimizarea valorii ponderilor în timpul antrenării

- Deprecierea ponderilor (weight decay, regularization):

$$E = E + pen = \frac{1}{2} \cdot \sum_{m=1}^M \left(d^{(m)} - o^{(m)} \right)^2 + \frac{1}{2} \cdot \lambda \cdot \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K w^2$$

$$\Delta w = -\eta \cdot \left(\frac{\partial E}{\partial w} \right) - \lambda \cdot w$$

- Efect: ponderile care trebuie să scadă devin și mai mici dacă eroarea E nu este mare, pentru că se scade din ele și valoarea $\lambda \cdot w$. Ponderile cu valoare absolută mică (inutile, nu afectează semnificativ antrenarea) tind spre 0.
- Coeficientul λ se alege de către utilizator.
- Altă variantă: reducerea ponderilor cu o valoare scalară.

Optimizarea valorii ponderilor în timpul antrenării

- Folosirea unui termen de moment

$$w_i^{(t+1)} = w_i^{(t)} - \eta \cdot \frac{\partial E}{\partial w_i^{(t)}} \Big|_{w=w_i^{(t)}} + \mu \cdot (w_i^{(t)} - w_i^{(t-1)}), \quad i = 0 \dots n$$

- Efect: Corecția ponderii ține cont de evoluția erorii în iterațiile anterioare (inerție), nu se modifică concomitent cu variația suprafeței erorii, ci cu o anumită întârziere, permițând trecerea mai rapidă peste minime locale și “ravine”.
 - ❖ La începutul antrenării, eroarea scade mai rapid către minimul global
 - ❖ Către sfârșitul antrenării, se reduc oscilațiile în apropierea minimului.
- Valoarea coeficientului (subunitară, 0.1...0.9) depinde de problemă și se determină prin încercări.

Preprocesarea datelor de intrare

- Când se folosesc funcții de activare care dau valori într-un interval limitat (de ex. sigmoid logistic între 0 și 1, tangent hiperbolic între -1 și 1), este benefică scalarea datelor de antrenare în același interval (de ex, pentru sigmoid logistic, scalarea între 0.15 și 0.85, pentru a menține antrenarea în zona de variație amplă a sigmoidului).
- În absența scalării, componentele $w_i \cdot x_i$ în care x are valori mai mari (de exemplu $x_i=3500$ vs $x_j=1.2$) influențează mai mult ieșirea neuronului.
- Scalarea se face pe coloane (elemente x_i ale unui model).

Preprocesarea datelor de intrare

- Pentru un vector de date (o coloană reprezentând valoarea unui parametru x_m dintr-un set de antrenare cu M modele):

$$X = (x_1, \dots, x_m, \dots, x_M)^T$$

- ❖ Scalarea min-max (între o valoare minimă și maximă):

$$x_m^{scalat} = \frac{x_m^{nescalat} - x_{\min}^{nescalat}}{x_{\max}^{nescalat} - x_{\min}^{nescalat}} \cdot (\lim_{scalare}^{\max} - \lim_{scalare}^{\min}) + \lim_{scalare}^{\min}$$

- ❖ Scalarea în intervalul [0,1]:

$$x_m^{scalat} = \frac{x_m^{nescalat} - x_{\min}^{nescalat}}{x_{\max}^{nescalat} - x_{\min}^{nescalat}}$$

Preprocesarea datelor de intrare

- ❖ Scalarea prin aducerea valorilor la medie 0 (centrarea pe medie)

$$x_m^{scalat} = x_m^{nescalat} - \bar{x} = x_m^{nescalat} - \frac{1}{M} \sum_{m=1}^M x_m^{nescalat}$$

- ❖ Standardizarea (scalarea după varianță) transformarea astfel încât media să fie 0 și deviația standard să fie 1 (transformarea z):

- Deviația standard (standard deviation): rădăcina pătrată din varianță (media pătratului diferențelor față de medie): variația datelor față de medie.

$$x_m^{scalat} = \frac{x_m^{nescalat} - \bar{x}}{std} \quad \text{unde} \quad std = \sqrt{\frac{1}{M-1} \sum_{m=1}^M (x_m^{nescalat} - \bar{x})^2}$$

- Caz particular - distribuția (clopotul) Gauss, caracteristică multor serii întâlnite în practică – cu număr egal de valori de o parte și de alta a mediei.
 - Mediile școlare.

Exemplu de scalare

- Scalarea după varianță:

Suprafață locuință [mp]	Venit lunar [lei]	Putere instalată [kW]	Consum mediu [kWh/lună]
x_1	x_2	x_3	y
110	3350	10	482
33	3100	8	208
42	5480	15	250
138	4450	20	550
28.5	2500	6	108
153	6600	12	573



Suprafață locuință [mp]	Venit lunar [lei]	Putere instalată [kW]	Consum mediu [kWh/lună]
x_1	x_2	x_3	y
0.461	-0.572	-0.361	0.608
-0.909	-0.732	-0.755	-0.779
-0.749	0.787	0.624	-0.566
0.959	0.130	1.609	0.953
-0.989	-1.115	-1.149	-1.285
1.226	1.502	0.033	1.069

Preprocesarea datelor de intrare

- Analiza statistică a setului de antrenare și:
 - ❖ Eliminarea anomaliilor (vezi supraantrenare)
 - ❖ Eliminarea parametrilor (coloanelor din matrice) corelați
 - Corelare – tendință aproximativ sincronă a valorilor

Îmbunătățirea convergenței – algoritmul RProp

- RProp corectează ponderile în funcție de semnul erorii în două iterații succesive, astfel:

$$w_{ij}^{t+1} = w_{ij}^t + \Delta w_{ij}^{t+1}$$

unde

$$\Delta w_{ij}^{t+1} = \begin{cases} -\operatorname{sign}\left[\frac{\partial E^{t+1}}{\partial w_{ij}}\right] \cdot \delta_{ij}^{t+1} & \text{if } \frac{\partial E^{t+1}}{\partial w_{ij}} \cdot \frac{\partial E^t}{\partial w_{ij}} \geq 0 \\ -\operatorname{sign}[\Delta w_{ij}^t] \cdot \delta_{ij}^{t+1} & \text{if } \frac{\partial E^{t+1}}{\partial w_{ij}} \cdot \frac{\partial E^t}{\partial w_{ij}} < 0 \end{cases}$$

$$\delta_{ij}^{t+1} = \begin{cases} \eta^+ \cdot \delta_{ij}^t & \text{if } \frac{\partial E^{t+1}}{\partial w_{ij}} \cdot \frac{\partial E^t}{\partial w_{ij}} > 0 \\ \eta^- \cdot \delta_{ij}^t & \text{if } \frac{\partial E^{t+1}}{\partial w_{ij}} \cdot \frac{\partial E^t}{\partial w_{ij}} < 0 \\ \delta_{ij}^t & \text{if } \frac{\partial E^{t+1}}{\partial w_{ij}} \cdot \frac{\partial E^t}{\partial w_{ij}} = 0 \end{cases}$$

sign = semn(+ sau -)

Îmbunătățirea convergenței – algoritmul RProp

- Dacă eroarea nu schimbă semnul (continuă să scadă), δ este amplificat cu un coeficient supraunitar η^+ și corecția ponderii este accelerată.
- Dacă eroarea schimbă semnul (oscilație, minim depășit) corecția erorii este diminuată cu coeficientul η^- pentru a se reveni în minim.
- Dacă una din derivate se anulează (am atins minimul, panta derivatei este 0), δ rămâne neschimbat.
- Valori recomandate: $\eta^+=1.2$, $\eta^-=0.5$, $\delta_{init}=0.5$.

va urma...

- Rețele neuronale artificiale pentru clasificare
 - ❖ RNA Kohonen