

Algoritmul de retro-propagare (*backpropagation*)

1. Obiectiv

Obiectivul acestui laborator este prezentarea algoritmului de retro-propagare (engl. “backpropagation”), utilizat pentru antrenarea rețelelor neuronale de tip perceptron multi-strat.

2. Descrierea algoritmului

Rețelele neuronale au capacitatea de a învăța, însă modalitatea concretă în care se realizează acest proces este dată de algoritmul folosit pentru antrenare. O rețea se consideră antrenată dacă aplicarea unui vector de intrare conduce la obținerea unei ieșiri dorite, sau foarte apropiate de aceasta.

Antrenarea constă în aplicarea secvențială a diferiți vectori de intrare și ajustarea ponderilor din rețea în raport cu o procedură predeterminată. În acest timp, ponderile conexiunilor converg gradual spre anumite valori pentru care fiecare vector de intrare produce vectorul de ieșire dorit. După aplicarea unei intrări, se compară ieșirea calculată cu ieșirea dorită, după care diferența este folosită pentru modificarea ponderilor cu scopul minimizării erorii la un nivel acceptabil.

Algoritmul backpropagation este cel mai cunoscut și utilizat algoritm de antrenare pentru rețele neuronale de tip perceptron multi-strat. Numit și *algoritmul delta generalizat*, el se bazează pe minimizarea diferenței dintre ieșirea dorită și ieșirea reală, prin metoda gradientului descendent. Se definește *gradientul* unei funcții $F(x,y,z,...)$ drept:

$$\vec{\nabla} F = \frac{\partial F}{\partial x} \vec{i} + \frac{\partial F}{\partial y} \vec{j} + \frac{\partial F}{\partial z} \vec{k} + \dots$$

Gradientul ne spune cum variază funcția în diferite direcții. Ideea algoritmului este găsirea minimului funcției de eroare în raport cu ponderile conexiunilor. Eroarea este dată de diferența dintre ieșirea dorită și ieșirea efectivă a rețelei. Cea mai utilizată funcție de eroare este eroarea medie pătratică (engl. “mean squared error”, MSE). Dacă avem în mulțimea de antrenare K vectori iar rețeaua are N ieșiri, atunci:

$$MSE = \frac{\sum_k \sum_n e_{kn}^2}{K \cdot N} = \frac{\sum_k \sum_n (y_{kn}^d - y_{kn})^2}{K \cdot N},$$

unde y_{kn}^d este valoarea dorită la ieșirea n a rețelei pentru vectorul k , iar y_{kn} este ieșirea efectivă a rețelei.

Algoritmul pentru un perceptron multi-strat cu un singur strat ascuns este următorul:

Pasul 1: Inițializarea

Toate ponderile și pragurile rețelei sunt inițializate cu valori aleatorii *nenule*, distribuite uniform într-un mic interval, de exemplu $(-0.1, 0.1)$.

Dacă valorile acestea sunt 0, gradientii care vor fi calculați pe parcursul antrenării vor fi tot 0 (dacă nu există o legătură directă între intrare și ieșire) și rețeaua nu va învăța. Este chiar indicată încercarea mai multor antrenări, cu ponderi inițiale diferite, pentru găsirea celei mai bune valori pentru funcția cost (minimul erorii). Dimpotrivă, dacă valorile inițiale sunt mari, ele tind să satureze unitățile respective. În acest caz, derivata funcției sigmoide este foarte mică. Ea acționează ca un factor de multiplicare în timpul învățării și deci unitățile saturate vor fi aproape blocate, ceea ce face învățarea foarte lentă.

Pasul 2: O nouă epocă de antrenare

O epocă reprezintă prezentarea tuturor exemplelor din mulțimea de antrenare. În majoritatea cazurilor, antrenarea rețelei presupune mai multe epoci de antrenare.

Pentru a păstra rigoarea matematică, ponderile vor fi ajustate numai după ce *toți* vectorii de test vor fi aplicați rețelei. În acest scop, gradientii ponderilor trebuie memorați și ajustați după fiecare vector din mulțimea de antrenare, iar la sfârșitul unei epoci de antrenare, se vor modifica ponderile *o singură dată*. Există și varianta “on-line”, mai simplă, în care ponderile sunt actualizate direct; în acest caz, poate conta ordinea în care sunt prezentați rețelei vectorii de test.

Se inițializează corecțiile ponderilor și eroarea curentă cu 0: $\Delta w_{ij} = 0$, $E = 0$.

Pasul 3: Propagarea semnalului înainte

3.1. La intrările rețelei se aplică un vector din mulțimea de antrenare.

3.2. Se calculează ieșirile neuronilor din stratul ascuns:

$$y_j = f\left(\sum_{i=1}^n x_i \cdot w_{ij} - \theta_j\right),$$

unde n este numărul de intrări ale neuronului j din stratul ascuns, iar f este funcția de activare sigmoidă.

Pentru simplitate, pragul este considerat ca fiind ponderea unei conexiuni suplimentare a cărei intrare este întotdeauna egală cu 1. Expresia de mai sus devine acum:

$$y_j = f\left(\sum_{i=1}^{n+1} x_i \cdot w_{ij}\right).$$

3.3. Se calculează ieșirile reale ale rețelei:

$$y_k = f\left(\sum_{j=1}^{m+1} y_j \cdot w_{jk}\right),$$

unde m este numărul de intrări ale neuronului k din stratul de ieșire.

3.4. Se actualizează eroarea medie pătratică pe epocă:

$$MSE = \frac{\sum_k \sum_n e_{kn}^2}{K \cdot N} = \frac{\sum_k \sum_n (y_{kn}^d - y_{kn})^2}{K \cdot N}.$$

Pasul 4: Propagarea erorilor înapoi și ajustarea ponderilor

4.1. Se calculează gradientii erorilor pentru neuronii din stratul de ieșire:

$$\delta_k = f' \cdot e_k,$$

unde f' este derivata funcției de activare iar eroarea $e_k = y_k^d - y_k$.

Dacă folosim sigmoida unipolară, derivata acesteia este:

$$f'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = f(x) \cdot (1 - f(x)).$$

Vom avea:

$$\delta_k = y_k \cdot (1 - y_k) \cdot e_k.$$

4.2. Se actualizează corecțiile ponderilor dintre stratul ascuns și stratul de ieșire:

$$\Delta w_{jk} = \Delta w_{jk} + \alpha \cdot y_j \cdot \delta_k,$$

unde α este rata de învățare.

4.3. Se calculează gradientii erorilor pentru neuronii din stratul ascuns:

$$\delta_j = y_j \cdot (1 - y_j) \cdot \sum_{k=1}^l \delta_k \cdot w_{jk},$$

unde l este numărul de ieșiri ale rețelei.

4.4. Se actualizează corecțiile ponderilor dintre stratul de intrare și stratul ascuns:

$$\Delta w_{ij} = \Delta w_{ij} + \alpha \cdot x_i \cdot \delta_j.$$

Pasul 5. O nouă iterație

Dacă mai sunt vectori de test în epoca de antrenare curentă, se trece la pasul 3.

Dacă nu, se actualizează ponderile tuturor conexiunilor pe baza corecțiilor ponderilor:

$$w_{ij} = w_{ij} + \Delta w_{ij}.$$

Dacă s-a încheiat o epocă, se testează dacă s-a îndeplinit criteriul de terminare ($E < E_{max}$ sau atingerea unui număr maxim de epoci de antrenare). Dacă nu, se trece la pasul 2. Dacă da, algoritmul se termină.

3. Considerente practice

În implementarea algoritmului, apar o serie de probleme practice, legate în special de alegerea parametrilor antrenării și a configurației rețelei.

În primul rând, o rată de învățare mică determină o convergență lentă a algoritmului, pe când una prea mare poate cauza eșecul (algoritmul va „sări” peste soluție). Pentru probleme relativ simple, o rată de învățare $\alpha = 0.5$ este acceptabilă, însă în general se recomandă rate de învățare în jur de 0.1 – 0.2.

Este recomandată preprocesarea și postprocesarea datelor, astfel încât rețeaua să opereze cu valori scalate, de exemplu în intervalul [0.1, 0.9] pentru sigmoida unipolară.

O altă problemă caracteristică acestui mod de antrenare este dată de minimele locale. Într-un minim local, gradientii erorii devin 0 și învățarea nu mai continuă. O soluție este încercarea mai multor antrenări independente, cu ponderi inițializate diferit la început, ceea ce crește probabilitatea găsirii minimului global. Pentru probleme mari, acest lucru poate fi greu de realizat și atunci pot fi acceptate și minime locale, cu condiția ca erorile să fie suficient de mici. De asemenea, se pot încerca diferite configurații ale rețelei, cu un număr mai mare de neuroni în stratul ascuns sau cu mai multe straturi ascunse, care în general conduc la minime locale mai mici. Totuși, deși minimele locale sunt într-adevăr o problemă, în practică nu reprezintă dificultăți de nesoluționat.

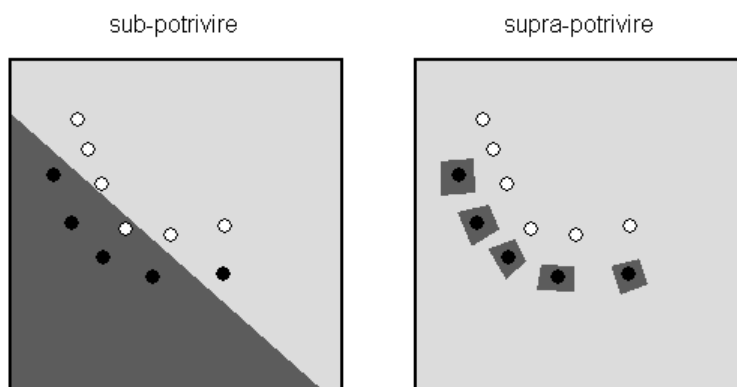
O chestiune importantă este alegerea unei configurații cât mai bune pentru rețea din punct de vedere al numărului de neuroni în straturile ascunse. În multe situații, un singur strat ascuns este suficient. Nu există niște reguli precise de alegere a numărului de neuroni. În general, rețeaua poate fi văzută ca un sistem în care numărul vectorilor de test înmulțit cu numărul de ieșiri reprezintă numărul de ecuații iar numărul ponderilor reprezintă numărul necunoscutelor. Ecuațiile sistemului sunt în general neliniare și foarte complexe și deci este foarte dificilă rezolvarea lor exactă prin mijloace convenționale. Algoritmul de antrenare urmărește tocmai găsirea unor soluții aproximative care să minimizeze erorile.

O rețea neuronală trebuie să fie capabilă de generalizare. Dacă rețeaua aproximează bine mulțimea de antrenare, aceasta nu este o garanție că va găsi soluții la fel de bune și pentru datele din altă mulțime, cea de test. Generalizarea presupune existența în date a unor regularități, a unui model

care *poate fi* învățat. În analogie cu sistemele liniare clasice, aceasta ar însemna niște ecuații redundante. Astfel, dacă numărul de ponderi este mai mic decât numărul de vectori de test, pentru o aproximare corectă rețeaua trebuie să se bazeze pe modelele intrinseci din date, modele care se vor regăsi și în datele de test. O regulă euristică afirmă că numărul de ponderi trebuie să fie în jur sau sub o zecime din produsul dintre numărul de vectori de antrenare și numărul de ieșiri. În unele situații însă (de exemplu, dacă datele de antrenare sunt relativ puține), numărul de ponderi poate fi chiar jumătate din produs.

Pentru un perceptron multi-strat se consideră că numărul de neuroni dintr-un strat trebuie să fie suficient de mare pentru ca acest strat să furnizeze trei sau mai multe laturi pentru fiecare regiune convexă identificată de stratul următor. Deci numărul de neuroni dintr-un strat ar trebui să fie aproximativ de trei ori mai mare decât cel din stratul următor.

După cum am menționat, un număr insuficient de ponderi conduce la „sub-potrivire” (engl. “underfitting”), în timp ce un număr prea mare de ponderi conduce la „supra-potrivire” (engl. “overfitting”), fenomene prezentate în figură. Același lucru apare dacă numărul de epoci de antrenare este prea mic sau prea mare.



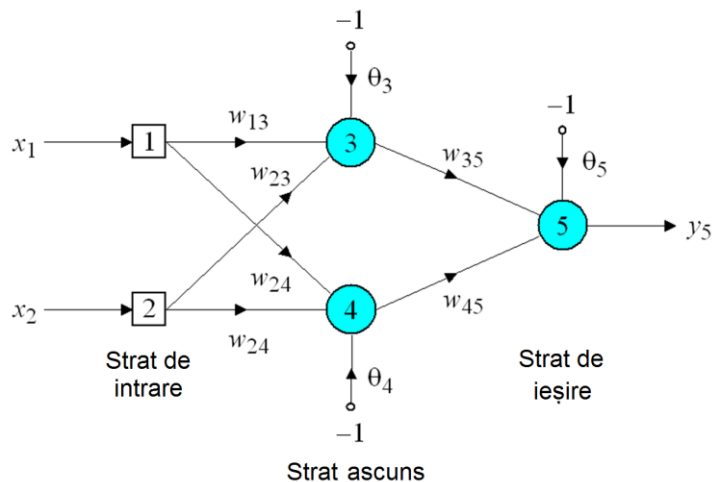
O metodă de rezolvare a acestei probleme este oprirea antrenării în momentul în care se atinge cea mai bună generalizare. Pentru o rețea suficient de mare, eroarea de antrenare scade în mod continuu pe măsură ce crește numărul epocilor de antrenare. Totuși, pentru date diferite de cele din mulțimea de antrenare (mulțimea de test), se constată că eroarea scade la început până la un punct în care începe din nou să crească. De aceea, oprirea antrenării trebuie făcută în momentul când eroarea pentru mulțimea de *test* este minimă.

Cea mai simplă modalitate de a estima capacitatea de generalizare este împărțirea mulțimii de date disponibile în 2/3 pentru antrenare și 1/3 pentru testare. Antrenarea se oprește atunci când începe să crească eroarea pentru mulțimea de validare, moment numit „punct de maximă generalizare”. În funcție de performanțele rețelei în acest moment, se pot încerca apoi diferite configurații, crescând sau micșorând numărul de neuroni din stratul (sau straturile) ascuns(e).

O metodă mai laborioasă, dar care este cea mai utilizată pentru evaluarea și compararea performanțelor modelelor de învățare, este *validarea încrucișată* (engl. “cross-validation”), în care datele sunt împărțite în n grupuri (de obicei 10), $n - 1$ grupuri se folosesc pentru antrenare, iar al n -lea pentru testare. Procesul se repetă de n ori, schimbând grupul de test. În final, se face media performanțelor obținute pentru cele n grupuri de test.

4. Aplicație

Implementați algoritmul *backpropagation* pentru o rețea neuronală de tip perceptron multi-strat, cu un singur strat ascuns și cu doi neuroni în stratul ascuns, cu scopul de a învăța funcții logice binare.



Programul trebuie să permită introducerea într-o formă convenabilă (de exemplu, dintr-un fișier sau direct de la consolă) a parametrilor algoritmului: număr maxim de epoci de antrenare, rata de învățare etc., precum și a funcției care trebuie aproximată. De exemplu, pentru funcția ȘI:

x_1	x_2	y_5
0	0	0
0	1	0
1	0	0
1	1	1

se va introduce “0 0 0 1”, pentru funcția XOR se va introduce “0 1 1 0” ș.a.m.d.

Modalitatea de notare

- (0.1 p) introducerea datelor
- (0.1 p) funcția de activare sigmoidă
- (0.1 p) propagarea înainte
- (0.1 p) gradientii și corecția ponderilor pentru stratul de ieșire (pentru propagarea înapoi)
- (0.1 p) gradientii și corecția ponderilor pentru stratul ascuns (pentru propagarea înapoi)
- (0.1 p) actualizarea ponderilor, procedura iterativă de antrenare, afișarea rezultatelor