

Artificial Neural Networks

(abbrev. ANNs, or simply NNs)

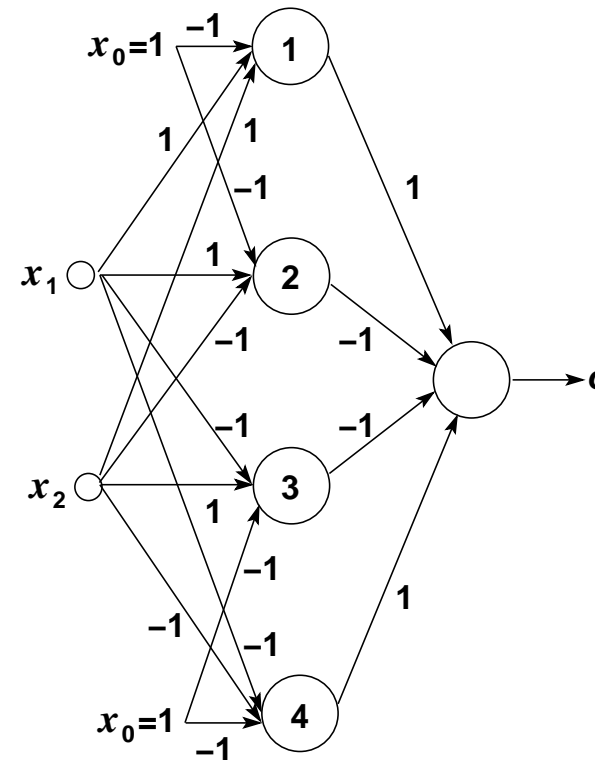
Calculation of the output of a neural network made of threshold perceptrons

CMU, 2010 fall, Aarti Singh, HW5, pr. 4.1.1

Considerăm rețeaua neuronală din figura alăturată. Toate unitățile acestei rețele neuronale sunt de tip prag, adică folosesc pentru activare funcția $sign$ definită prin $sign(z) = 1$ dacă $z \geq 0$ și -1 în rest. Pentru unitatea de pe nivelul de ieșire (lăsată nenumerotată), ponderea corespunzătoare termenului liber ($x_0 = 1$) este 0.

a. Scrieți funcția matematică calculată de fiecare dintre unitățile rețelei, în raport cu intrările x_1 și x_2 . Veți nota cu o_i (unde $i = 1, \dots, 4$) ieșirile unităților de pe nivelul ascuns și cu o ieșirea rețelei, adică valoarea produsă de către unitatea de pe nivelul de ieșire.

b. Calculați output-ul rețelei atunci când intrările x_1 și x_2 iau valori în mulțimea $\{-1, 1\}$.



Răspuns

a.

$$o_1(x_1, x_2) = \text{sign}(x_1 + x_2 - 1),$$

$$o_2(x_1, x_2) = \text{sign}(x_1 - x_2 - 1),$$

$$o_3(x_1, x_2) = \text{sign}(-x_1 + x_2 - 1),$$

$$o_4(x_1, x_2) = \text{sign}(-x_1 - x_2 - 1),$$

$$o(x_1, x_2) = \text{sign}(o_1(x_1, x_2) - o_2(x_1, x_2) - o_3(x_1, x_2) + o_4(x_1, x_2)).$$

b.

x_1	x_2	o_1	o_2	o_3	o_4	o
1	1	1	-1	-1	-1	1
1	-1	-1	1	-1	-1	-1
-1	1	-1	-1	1	-1	-1
-1	-1	-1	-1	-1	1	1

c. Indicați funcțiile boolene reprezentate de către unitățile de pe nivelul ascuns (1, 2, 3 și 4) atunci când $x_1, x_2 \in \{-1, 1\}$. Procedați similar pentru ieșirea o .

Răspuns:

Din tabelul obținut la punctul precedent este imediat că, atunci când $x_1, x_2 \in \{-1, 1\}$, ieșirile calculate de către unitățile 1, 2, 3 și 4 corespund conceptelor/funcțiilor $x_1 \wedge x_2$, $x_1 \wedge \neg x_2$, $\neg x_1 \wedge x_2$ și respectiv $\neg x_1 \wedge \neg x_2$. Ieșirea rețelei, o , corespunde conceptului $\neg(x_1 \text{ XOR } x_2)$.

Observație: Se poate remarca faptul că în această rețea un neuron (și doar unul!) de pe nivelul ascuns este activat la fiecare combinație de valori (din cele 4 posibile) pentru $x_1, x_2 \in \{-1, 1\}$.

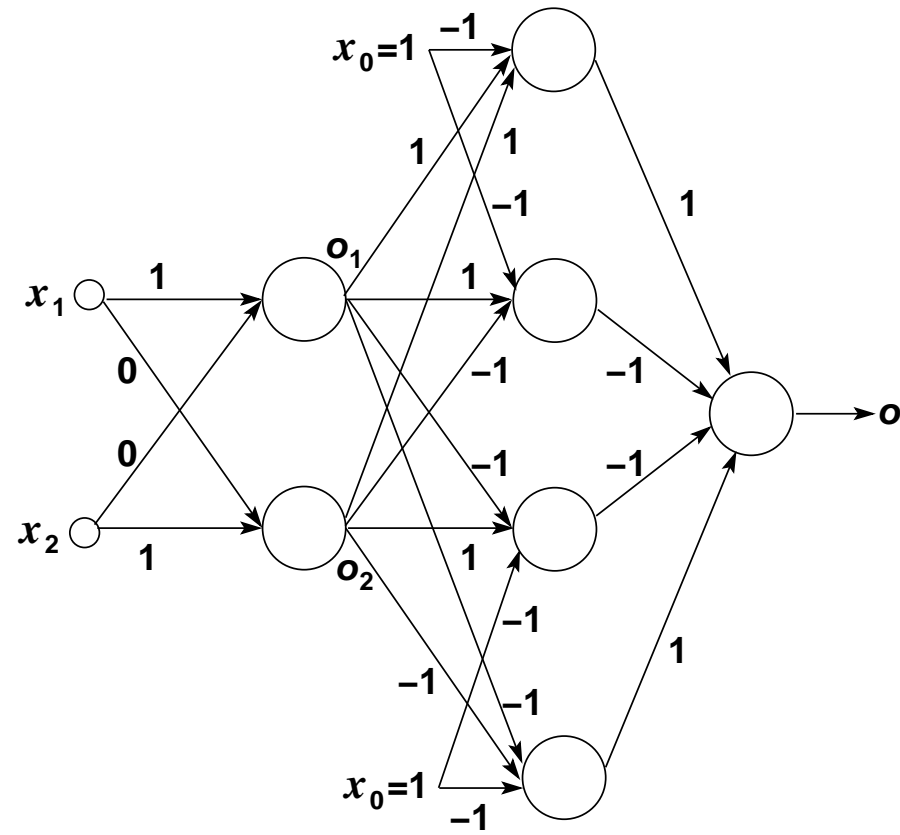
d. Specificați cum anume ar putea fi modificată rețeaua dată astfel încât noua variantă să calculeze funcția de variabile reale (nu boolene ca mai înainte!) x_1 și x_2 , a cărei relație de definiție este: $f(x_1, x_2) = 1$ dacă $x_1, x_2 \geq 0$ sau $x_1, x_2 < 0$, și -1 în caz contrar.

Răspuns:

Este imediat că funcția indicată în enunț se poate scrie sub forma

$$f(x_1, x_2) = \begin{cases} 1 & \text{dacă } \text{sign}(x_1) \cdot \text{sign}(x_2) \geq 0 \\ -1 & \text{dacă } \text{sign}(x_1) \cdot \text{sign}(x_2) < 0. \end{cases}$$

Se observă imediat că aceasta coincide cu funcția $\neg(\text{sign}(x_1) \text{ XOR } \text{sign}(x_2))$. Prin urmare, este suficient să adăugăm la rețeaua dată în enunț un nivel ascuns suplimentar, format din două unități cu funcție de activare de tip prag, care să transforme intrările x_1 și x_2 în $\text{sign}(x_1)$ și respectiv $\text{sign}(x_2)$. Vom obține ca rezultat rețeaua din slide-ul următor.



Expressivity of neural networks: boolean functions

CMU, 2010 fall, Aarti Singh, HW5, pr. 4.3

Consider the set of binary functions over d Boolean (binary) variables: $F = \{f : \{0, 1\}^d \rightarrow \{0, 1\}\}$. Show that any function in F can be represented by a neural network with a single hidden layer.

Note: You can work with values other than $\{0, 1\}$, such as $\{1, -1\}$, as long as the function is still binary.

Solution

Let S denote the set of all possible d -bit binary vectors. For any function $f \in F$, define $T_f = \{b \in \{0,1\}^d \mid f(b) = 1\}$. We construct a neural network with one hidden layer to represent f as follows:

The neural network has $|T_f|$ hidden units, each of which corresponds to a distinct binary vector in T_f . The weights from the input layer to the i -th hidden unit are determined by the values of the corresponding binary vector b_i : the j -th weight is 1 if b_{ij} is 1 and -1 otherwise. The activation functions in the hidden units are threshold functions:

$$h_i(v) = \begin{cases} 1 & \text{if } v \geq \sum_{j=1}^d b_{ij}, \\ 0 & \text{otherwise.} \end{cases}$$

It is easy to see that h_i outputs 1 if and only if the input is b_i . The weights from all hidden units to the output layer are 1, and the activation function in the output unit is the identity function.

To show that the above neural network represents f , we first note that for every b in T_f exactly one hidden unit, the one corresponding to b , outputs 1, and therefore the neural net outputs 1. For any $b \notin T_f$, all hidden units output 0, and the neural network outputs 0. This completes the proof.

Exemplification:

CMU, 2011 spring, Roni Rosenfeld, HW4, pr. 1.c

A neural network with only one hidden layer that implements $(A \vee \neg B) \text{ XOR } (\neg C \vee D)$.

Solution (in Romanian)

Explicitând definiția funcției logice XOR și apoi aplicând regulile lui DeMorgan, expresia booleană din enunț devine:

$$\begin{aligned}(A \vee \neg B) \text{ XOR } (\neg C \vee D) &= [(A \vee \neg B) \wedge \neg(\neg C \vee D)] \vee [\neg(A \vee \neg B) \wedge (\neg C \vee D)] \\ &= [(A \vee \neg B) \wedge (C \wedge \neg D)] \vee [(\neg A \wedge B) \wedge (\neg C \vee D)]\end{aligned}$$

Întrucât operatorii logici \vee și \wedge sunt distributivi unul față de celălalt, rezultă că:

$$(A \vee \neg B) \text{ XOR } (\neg C \vee D) = (A \wedge C \wedge \neg D) \vee (\neg B \wedge C \wedge \neg D) \vee (\neg A \wedge B \wedge \neg C) \vee (\neg A \wedge B \wedge D)$$

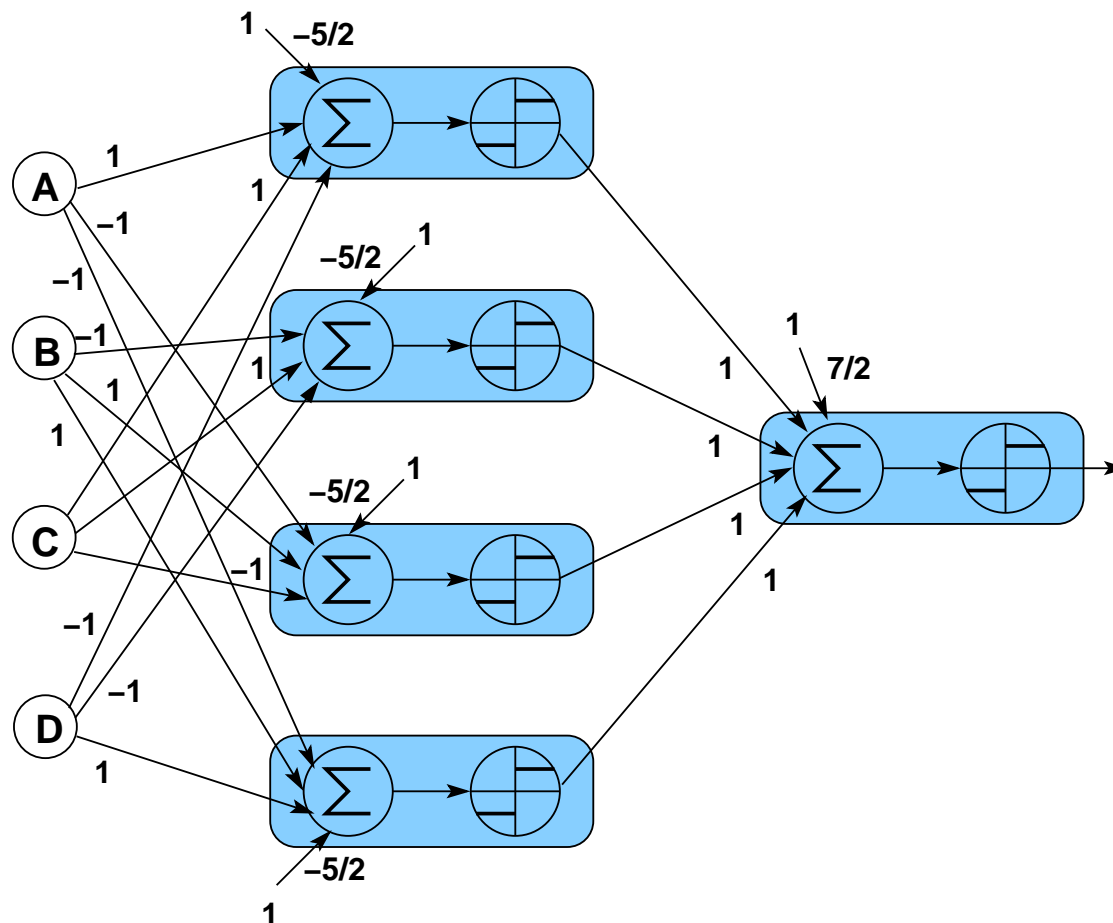
Fiecare din cele patru paranteze din partea dreaptă a egalității de mai sus poate fi reprezentată cu ajutorul unui perceptron-prag cu patru intrări, și anume câte o intrare pentru fiecare din cele trei variabile boolene din paranteză, plus o intrare pentru termenul liber.

De *exemplu*, conjuncției $A \wedge C \wedge \neg D$ îi putem asocia inegalitatea $x + z - t > 5/2$, deci ponderile intrărilor perceptronului corespunzător vor fi $-5/2, 1, 1$ și -1 .)

Cei patru perceptroni vor fi plasați pe primul nivel ascuns al rețelei pe care o construim.

Apoi, ieșirile acestor patru perceptroni vor constitui intrări pentru un alt perceptron-prag, care să reprezinte disjuncția a patru variabile boolene. (Acestui perceptron, situat pe nivelul de ieșire, îi putem asocia, de exemplu, inegalitatea $x + y + z + t > -7/2$.)

Rețeaua neuronală rezultată este cea din figura alăturată.



A negative expressivity result
concerning networks of threshold perceptrons

CMU, 2010 fall, Aarti Singh, HW5, pr. 4.1.2

Consider the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined as $f(x_1, x_2) = 1$ for $x_1, x_2 \geq 0$ or $x_1, x_2 < 0$, and -1 otherwise. It was shown (see CMU, 2010 fall, Aarti Singh, HW5, pr. 4.1.1) that it can be represented using a network of threshold perceptrons with two hidden layers.

Show that no neural network made of threshold units, with only one hidden layer, can represent f . For simplicity, you shall assume that the number of hidden units in a hidden layer is finite, but it can be arbitrarily large.

Hint

When there is only one hidden layer, the separating border corresponding to each hidden unit is a line in the 2-d plane such that the hidden unit outputs $+1$ on one side of the line and -1 on the other.

Consider a circular neighborhood (V) of the origin such that if a separating line determined by a hidden unit that crosses the neighborhood V , then it passes through the origin. (Since the number of hidden units is finite, we can always find such a neighborhood.)

Can any neural network made only of threshold units, with one hidden layer, represent f in the neighborhood V ?

Solution

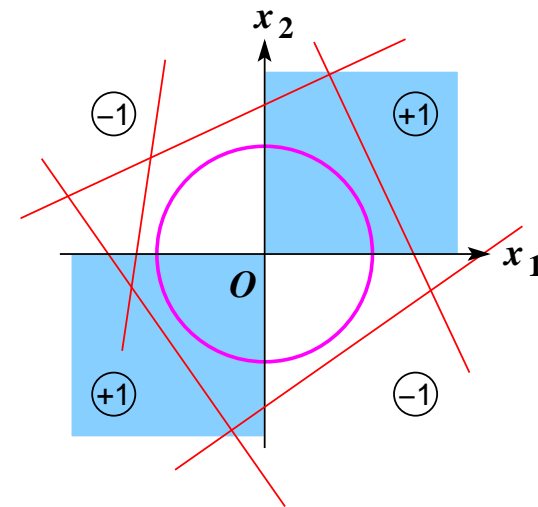
Given any neural network with one hidden layer of finite hidden units, we will consider a circular neighborhood of the origin

$$N_\varepsilon = \{(x_1, x_2) | \sqrt{x_1^2 + x_2^2} \leq \varepsilon\},$$

with ε being [less then] the minimum of the distances from the origin to the separating lines (hidden units) which do not pass through the origin.

We will consider two cases:

Case 1: No line crosses the neighborhood. In this case, the neural network must output the same value in this neighborhood, but f outputs 1 or -1 . Therefore, neural networks in this case cannot represent f .



Case 2

16.

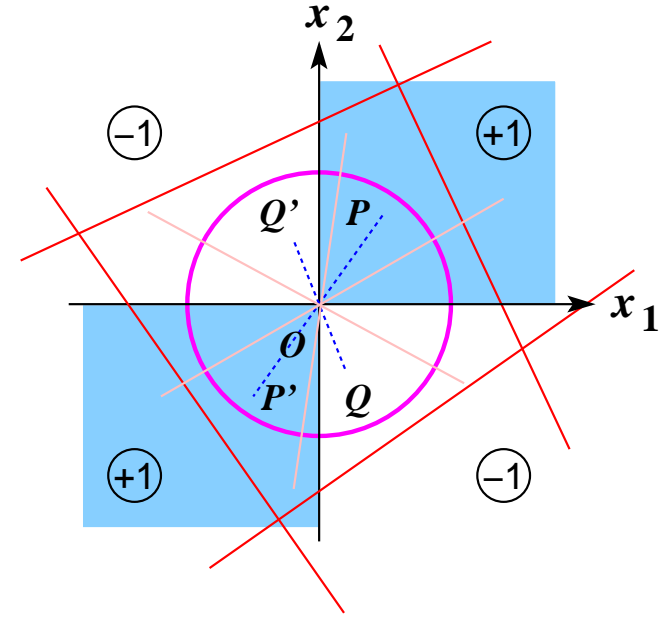
Some lines cross the neighborhood. By construction, these lines must all pass through the origin. Denote the set of the hidden units corresponding to these lines as O . Let us consider a point $P = (p_1, p_2)$ situated in this neighborhood and in the first quadrant, that does not fall on any line (see the figure). We can always find such a point because the number of the hidden units, hence the number of lines, is finite. By symmetry, the point $P' = (-p_1, -p_2)$, also in the neighborhood, does not fall on any line either, and $f(P) = f(P') = 1$. Let $y_h(P) \in \{1, -1\}$ denote the output of the hidden unit h on P .

It is easy to see the following:

$$y_h(P) = -y_h(P') \quad \forall h \in O, \quad \text{and} \quad y_h(P) = y_h(P') \quad \forall h \notin O.$$

Let w_h denote the weight from the hidden unit h to the output unit. We thus have

$$\begin{aligned} \sum_h w_h y_h(P) &= \sum_{h \in O} w_h y_h(P) + \sum_{h \notin O} w_h y_h(P) = - \sum_{h \in O} w_h y_h(P') + \sum_{h \notin O} w_h y_h(P') \\ &= - \sum_h w_h y_h(P') + 2 \sum_{h \notin O} w_h y_h(P'). \end{aligned} \tag{1}$$



Similarly, we can pick Q and $Q' = -Q$ in the fourth and the second quadrantal parts of the neighborhood, and have that

$$\sum_h w_h y_h(Q) = - \sum_h w_h y_h(Q') + 2 \sum_{h \notin O} w_h y_h(Q'). \quad (2)$$

Since $f(P) = f(P') = 1$ and $f(Q) = f(Q') = -1$, in order to represent f , the free weight w_0 (corresponding to $x_0 = 1$) for the output unit of the neural net must satisfy the inequalities

$$\begin{aligned} \sum_h w_h y_h(P) &\geq -w_0, \quad \sum_h w_h y_h(P') \geq -w_0 \stackrel{(1)}{\Rightarrow} \sum_{h \notin O} w_h y_h(P') \geq -w_0, \\ \sum_h w_h y_h(Q) &< -w_0, \quad \sum_h w_h y_h(Q') < -w_0 \stackrel{(2)}{\Rightarrow} \sum_{h \notin O} w_h y_h(Q') < -w_0, \end{aligned}$$

which cannot happen because $\sum_{h \notin O} w_h y_h(P') = \sum_{h \notin O} w_h y_h(Q')$, since $y_h(P') = y_h(Q')$ for every $h \notin O$. This completes the proof.

Expressivity of neural networks:

Any Lipschitz continuous function
defined on a bounded interval from \mathbb{R}
can be approximated by a neural network
with a single hidden layer

CMU, 2011 fall, Tom Mitchell, Aarti Singh, HW5, pr. 2.3

Let $f(x)$ be any function whose domain is $[C, D)$, for real values $C < D$. Suppose that the function is Lipschitz continuous, that is,

$$\forall x, x' \in [C, D), |f(x') - f(x)| \leq L|x' - x|,$$

for some constant $L \geq 0$.

Construct a neural network with one hidden layer that approximates this function within $\varepsilon > 0$, that is, $\forall x \in [C, D), |f(x) - \text{out}(x)| \leq \varepsilon$, where $\text{out}(x)$ is the output of your neural network given the input x .

Hint: You may use as building blocks the network(s) designed at CMU, 2007 fall, Carlos Guestrin, HW2, pr. 4.

Note: Your network should use only the activation functions g_I and g_S (the identity and respectively the step function) mentioned at above referred problem.

You need to specify the number K of hidden units, the activation function for each unit, and a formula for calculating each weight $w_0, w_k, w_0^{(k)}$, and $w_1^{(k)}$, for each $k \in \{1, 2, \dots, K\}$. These weights may be specified in terms of C, D, L and ε , as well as the values of $f(x)$ evaluated at a finite number of x values of your choosing (you need to explicitly specify which x values you use). You do not need to explicitly write the $out(x)$ function.

Why does your network attain the given accuracy?

Răspuns

Este imediat că din ipoteza $|f(x) - f(x')| \leq L|x - x'|$, în cazul în care $|x - x'| \leq \frac{\varepsilon}{L}$, rezultă

$$|f(x) - f(x')| \leq \varepsilon \quad (3)$$

Așadar, este de dorit să „acoperim” intervalul $[C, D)$ cu intervale de lungime $\frac{2\varepsilon}{L}$ (sau mai mică):

$$[C, C + \frac{2\varepsilon}{L}), [C + \frac{2\varepsilon}{L}, C + 2\frac{2\varepsilon}{L}), \dots, [C + (K - 1)\frac{2\varepsilon}{L}, D))$$

și să luăm punctele x' de forma $\xi_i = (\alpha_i + \beta_i)/2$, pentru cu $i = 1, \dots, K$, unde

$$K \stackrel{\text{not.}}{=} \left\lceil (D - C)\frac{L}{2\varepsilon} \right\rceil,$$

$$\alpha_1 = C, \beta_1 = \alpha_2 = C + \frac{2\varepsilon}{L}, \dots, \beta_{K-1} = \alpha_K = C + (K - 1)\frac{2\varepsilon}{L}, \beta_K = D.$$

Așadar, prin ξ_i am notat mijlocul intervalului i de mai sus, $[\alpha_i, \beta_i)$.

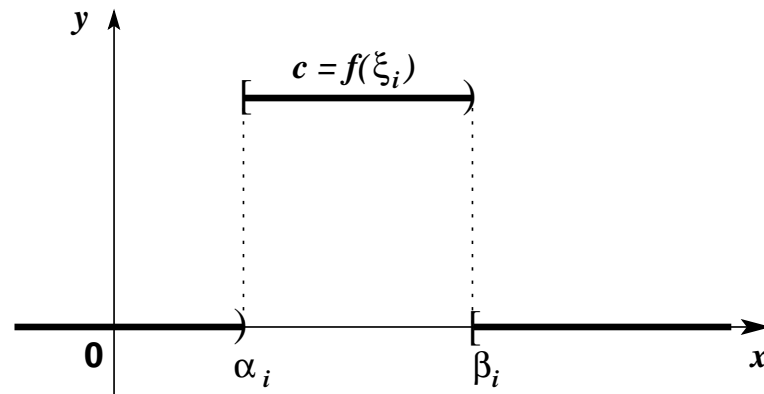
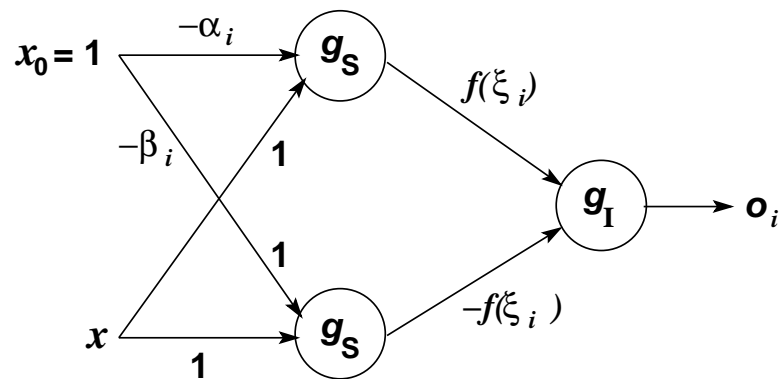
În continuare, vom arăta că putem construi o rețea neuronală cu un singur nivel ascuns și care folosește doar unități de tip liniar sau de tip prag, astfel încât

$$\text{pentru } \forall x \in [C, D), \text{ dacă } x \in [\alpha_i, \beta_i), \text{ atunci } out(x) \stackrel{def.}{=} f(\xi_i). \quad (4)$$

În consecință, dacă în relația (3) vom înlocui x' cu ξ_i , va rezulta imediat că $|f(x) - out(x)| \leq \varepsilon$.

Revenind acum la proprietatea (4), vom arăta că există o rețea având K unități pe nivelul ascuns (și o singură unitate de ieșire), astfel încât, dacă $x \in [\alpha_i, \beta_i)$, atunci unitatea i de pe nivelul ascuns se activează (producând valoarea $f(\xi_i)$), iar toate celelalte unități de pe nivelul ascuns rămân neactivate (adică produc output 0). Prin urmare, rezultatul va fi exact cel dorit.

Conform problemei CMU, 2007 fall, Carlos Guestrin, HW2, pr. 4.b, există o rețea neuronală care produce valoarea $c = f(\xi_i)$ pentru orice input $x \in [\alpha_i, \beta_i)$ și 0 în rest, și care folosește funcții de activare g_S pe nivelul ascuns și g_I pe nivelul de ieșire:

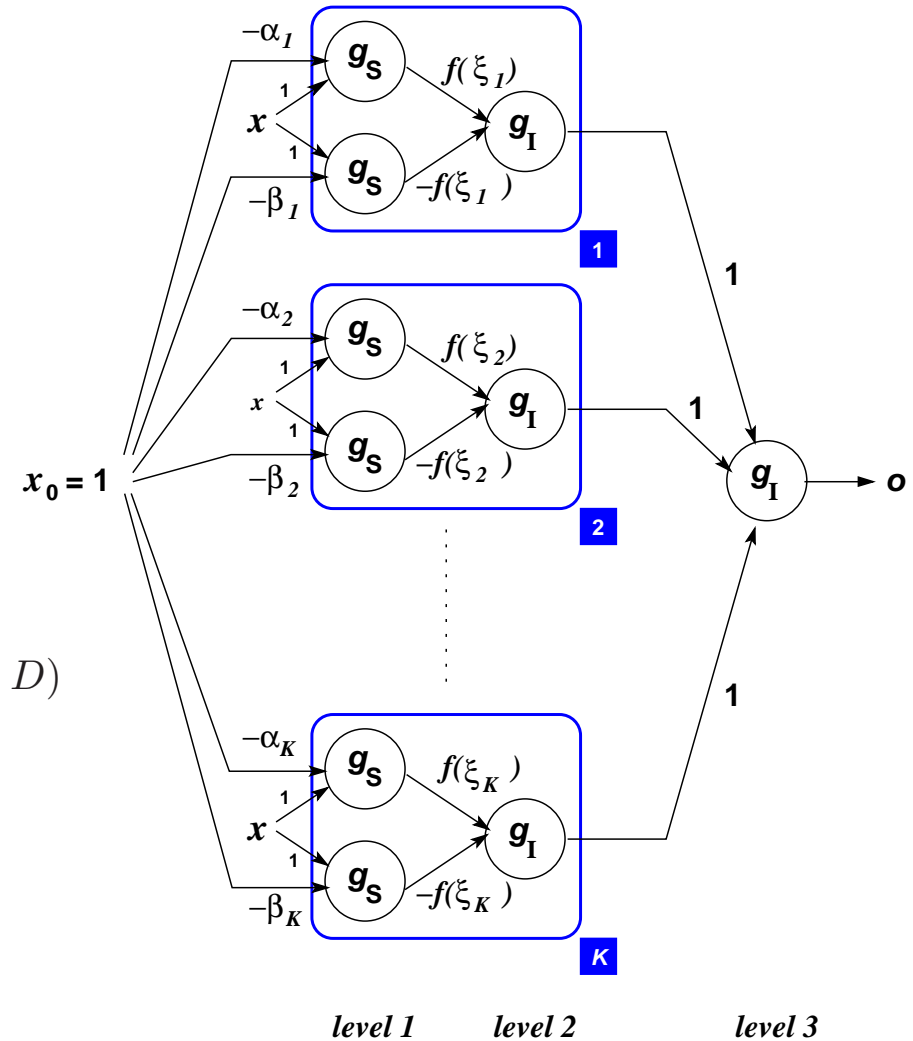


„Ansamblând“ K astfel de rețele, vom obține rețeaua din figura alăturată.

Output-ul acestei rețele este exact cel dorit (deci satisface condiția din enunț):

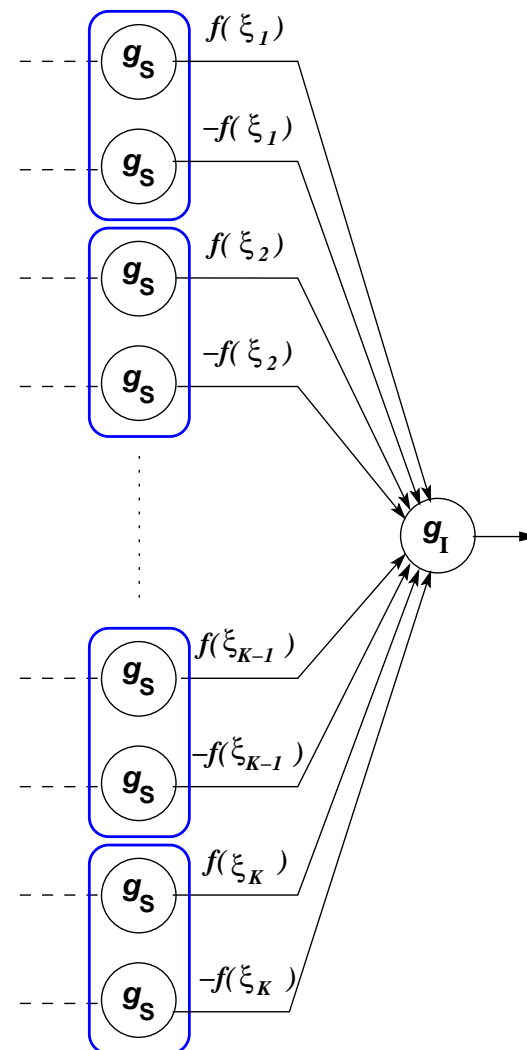
$$\mathbf{out}(x) = \begin{cases} f(\xi_1) & \text{pt. } x \in [\alpha_1 = C, \beta_1) \\ f(\xi_2) & \text{pt. } x \in [\alpha_2 = \beta_1, \beta_2) \\ \dots & \dots \dots \\ f(\xi_K) & \text{pt. } x \in [\alpha_K = \beta_{K-1}, \beta_K = D) \end{cases}$$

$$\Rightarrow |f(x) - \mathbf{out}(x)| \leq \varepsilon, \forall x \in [C, D).$$

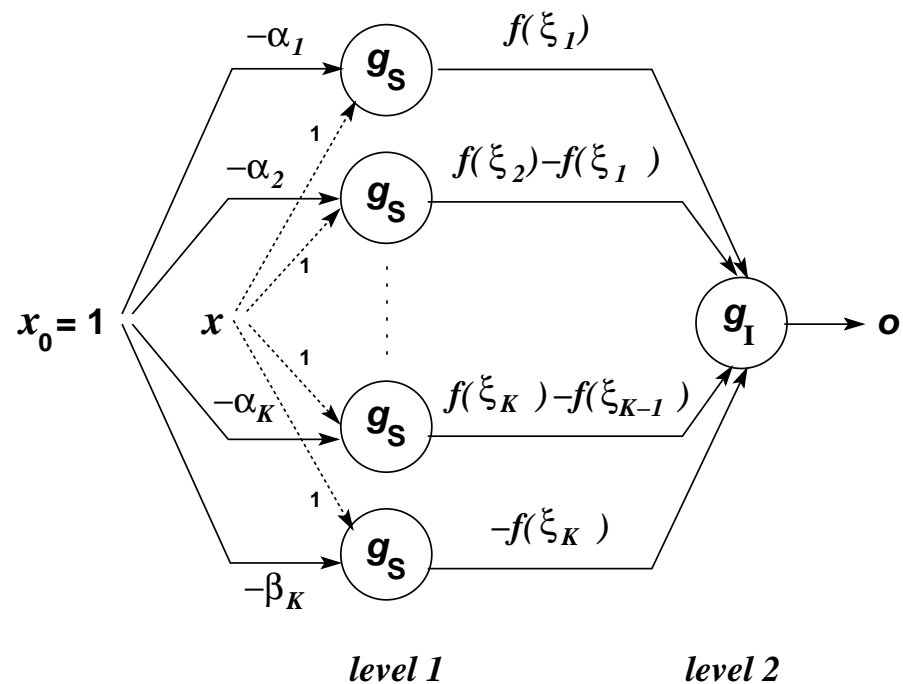


Neajunsul este că rețeaua aceasta are două niveluri ascunse, în loc de unul singur, așa cum se cere în enunț.

Totuși, el poate fi „corectat” imediat, comasând nivelurile 2 și 3 — formate doar din unități liniare — într-o singură unitate liniară (care va constitui nivelul de ieșire al noii rețele), ca în figura alăturată.



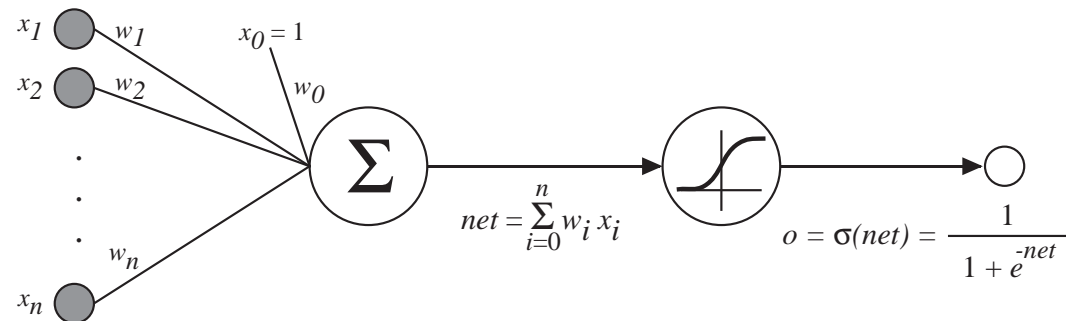
De asemenea, datorită faptului că $\beta_1 = \alpha_2, \beta_2 = \alpha_3, \dots, \beta_{K-1} = \alpha_K$, nivelul ascuns (al acestei noi rețele), format din cele $2K$ unități de tip prag, poate fi echivalat cu un altul, compus din doar $K + 1$ unități de tip prag. La final, obținem rețeaua din figura alăturată.



The sigmoidal / logistic perceptron:

Exemplification: setting up its weights so as to represent
a given boolean expression

CMU, 2003 fall, T. Mitchell, A. Moore, midterm, pr. 6.a



Consider a single sigmoid threshold unit with three inputs, x_1 , x_2 , and x_3 .

$$y = \sigma(w_0 + w_1x_1 + w_2x_2 + w_3x_3) \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}.$$

We input values of either 0 or 1 for each of these inputs. Assign values to weights w_0 , w_1 , w_2 and w_3 so that the output of the sigmoid unit is greater than 0.5 if and only if $(x_1 \text{ AND } x_2) \text{ OR } x_3$ is true.

Solution (in Romanian)

Vom scrie mai întâi tabela de valori a funcției / expresiei boolene date. Vom adăuga în această tabelă o coloană care va conține forma [particulară a] sumei ponderate $w_0 + w_1x_1 + w_2x_2 + w_3x_3$ — adică, *argumentul* pe care-l va lua funcția sigmoidală (σ) — pentru fiecare combinație de valori ale variabilelor x_1, x_2, x_3 . În plus, la fiecare linie vom preciza și *condiția* care trebuie să fie satisfăcută de către suma ponderată de pe linia respectivă, în așa fel încât perceptronul să realizeze codificarea cerută în enunț.

x_1	x_2	x_3	$(x_1 \wedge x_2) \vee x_3$	$w_0 + w_1x_1 + w_2x_2 + w_3x_3$	$\vdots 0$
0	0	0	0	w_0	< 0
0	0	1	1	$w_0 + w_3$	> 0
0	1	0	0	$w_0 + w_2$	< 0
0	1	1	1	$w_0 + w_2 + w_3$	> 0
1	0	0	0	$w_0 + w_1$	< 0
1	0	1	1	$w_0 + w_1 + w_3$	> 0
1	1	0	1	$w_0 + w_1 + w_2$	> 0
1	1	1	1	$w_0 + w_1 + w_2 + w_3$	> 0

(5) Așadar, rezumând, pentru ca unitatea sigmoidală să codifice expresia booleană dată, trebuie ca ponderile w_i să satisfacă sistemul de inecuații scrise în ultima coloană a tabelului alăturat.

În vederea găsirii unei soluții pentru acest sistem, putem observa că în expresia $(x_1 \wedge x_2) \vee x_3$ variabilele logice x_1 și x_2 joacă rol simetric în raport cu operatorul \wedge . Deci este natural să considerăm că ponderile alocate lui x_1 și x_2 [ca input-uri] în perceptronul sigmoidal pot fi egale. Introducând deci restricția $w_1 = w_2$, sistemul (5) va deveni cel scris în dreapta.

$$\left\{ \begin{array}{l} w_0 < 0 \\ w_0 + w_3 > 0 \\ w_0 + w_1 < 0 \\ w_0 + w_1 + w_3 > 0 \\ w_0 + 2w_1 > 0 \\ w_0 + 2w_1 + w_3 > 0 \end{array} \right. \quad (6)$$

Mai departe, analizând din nou expresia $(x_1 \wedge x_2) \vee x_3$, este natural să gândim că, în raport cu operatorul \vee , ponderea variabilei x_3 ar trebui să fie aceeași cu ponderile „cumulate” ale variabilelor x_1 și x_2 . Prin urmare, „injectând” și restricția $w_3 = w_1 + w_2 = 2w_1$, sistemul de inecuații (6) va deveni cel scris în dreapta.

$$\left\{ \begin{array}{l} w_0 < 0 \\ w_0 + 2w_1 > 0 \\ w_0 + w_1 < 0 \\ w_0 + 3w_1 > 0 \\ w_0 + 4w_1 > 0 \end{array} \right. \quad (7)$$

Din forma sistemului (7), apare natural să explorăm ce anume se întâmplă dacă restricționăm spațiul de soluții impunând condiția $w_1 > 0$. În această ipoteză, sistemul (7) va avea aceeași soluție ca și inecuația dublă

$$w_0 + w_1 < 0 < w_0 + 2w_1,$$

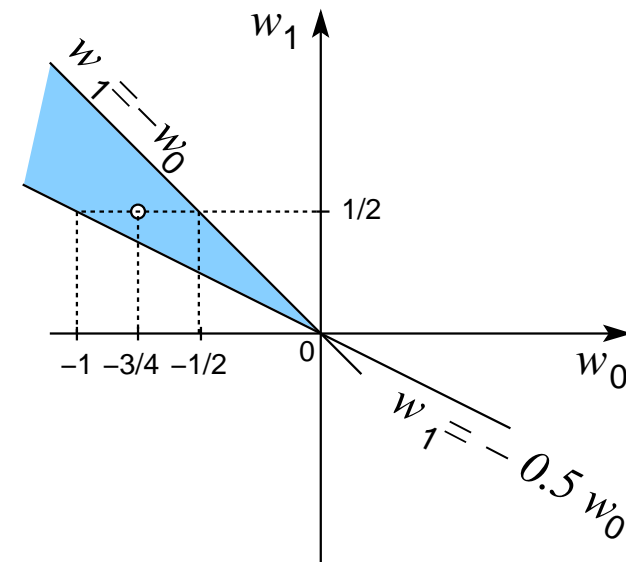
care este echivalentă cu

$$w_1 < -w_0 < 2w_1. \quad (8)$$

Pentru această ultimă inecuație dublă este foarte ușor să indicăm soluții. Iată una dintre ele: $w_1 = \frac{1}{2}$ și $w_0 = -\frac{3}{4}$. Prin urmare, $w_2 = \frac{1}{2}$ și $w_3 = 1$.

Observație (1)

De fapt, dacă lucrăm într-un reper de coordonate $w_0 O w_1$, orice punct (w_0, w_1) din cadranul al doilea [și] care este situat între dreptele de ecuații $w_1 = -w_0$ și respectiv $w_1 = -\frac{1}{2}w_0$ este o soluție a inecuației duble (8), deci și a sistemului de inecuații (5), la care, așa cum am văzut, au fost adăugate restricțiile $w_1 = w_2 = \frac{1}{2}w_3$.



(Invers, adică a preciza care este întreg spațiul de soluții al sistemului (5) este mult în afara obiectivului acestui exercițiu.)

Observație (2)

Rostul acestui [tip de] exercițiu este să arate că pentru a găsi valori corespunzătoare ponderilor unui perceptron în așa fel încât el să reprezinte/codifice o anumită funcție, putem apela la [rezolvarea de] sisteme de restricții / inecuații.

Mai general, așa cum precizează și Tom Mitchell în Machine Learning la pag. 95, putem folosi metode de programare liniară sau ne-liniară.

Pe de altă parte, acest exercițiu pune în evidență în mod indirect faptul că ar fi de dorit ca (alternativ) să dispunem de o procedură de calcul generală, cât mai simplă, prin care să atingem obiectivul menționat anterior.

O astfel de procedură — bazată pe metoda gradientului descendent — va face obiectul problemelor următoare. În raport cu programarea liniară sau cea ne-liniară, această procedură are avantajul că se scalează în mod elegant/convenabil la nivel de rețea neuronală.

**Rosenblatt's *Perceptron* Algorithm:
some simple properties**

CMU, 2015 spring, Tom Mitchell, Nina Balcan, HW6, pr 3.ab

CMU, 2017 spring, Barnabas Póczos, HW3, pr. 1.bc

Consider running the *Perceptron* algorithm (see the nearby pseudo-code) on some sequence of examples S . (Remember that an example is a data point and its label.)

```
initialize  $\bar{w} \leftarrow \bar{0}$ 
for  $i = 1, \dots, n$  do
  if  $y_i (\bar{w} \cdot \bar{x}_i) \leq 0$  then
     $\bar{w} \leftarrow \bar{w} + y_i \bar{x}_i$ 
  end if
end for
```

Let S' be the same set of examples as S , but presented in a different order.

a. Does the *Perceptron* algorithm necessarily make the same number of mistakes on S as it does on S' ?

If so, why? If not, show such an S and S' where the *Perceptron* algorithm makes a different number of mistakes on S' than it does on S .

Solution (in Romanian)

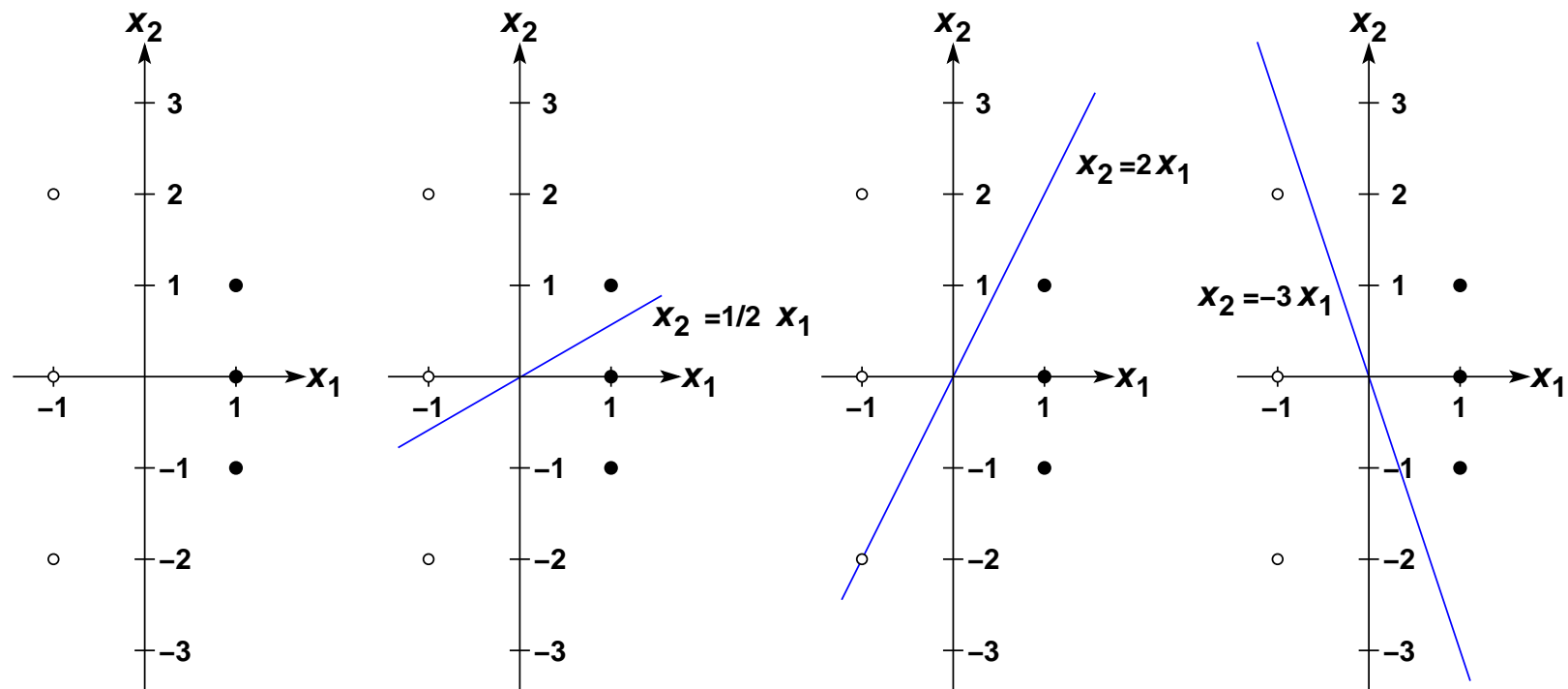
a. Considerăm următorul set de exemple (S), în ordinea indicată:

Exemplul	1	2	3	4	5	6
Instanța (x_1, x_2)	$(-1, 2)$	$(1, 0)$	$(1, 1)$	$(-1, 0)$	$(-1, -2)$	$(1, -1)$
Eticheta y	-1	$+1$	$+1$	-1	-1	$+1$

Sintetizăm execuția algoritmului *Perceptron* pe aceste date, fără a folosi termenul “bias” (w_0), alcătuind tabelul următor:

Iterația(i)	\bar{x}_i	y_i	$y_i \bar{w} \cdot \bar{x}_i$	\bar{w}	greșeală	ec. separatorului
inițializare	—	—	$(0, 0)$	—	—	
1	$(-1, 2)$	-1	$0 \leq 0$	$(1, -2)$	da	$x_2 = 0.5x_1$
2	$(1, 0)$	$+1$	$1 > 0$	$(1, -2)$	nu	$x_2 = 0.5x_1$
3	$(1, 1)$	$+1$	$-1 \leq 0$	$(2, -1)$	da	$x_2 = 2x_1$
4	$(-1, 0)$	-1	$2 > 0$	$(2, -1)$	nu	$x_2 = 2x_1$
5	$(-1, -2)$	-1	$0 \leq 0$	$(3, 1)$	da	$x_2 = -3x_1$
6	$(1, -1)$	$+1$	$2 > 0$	$(3, 1)$	nu	$x_2 = -3x_1$

Așadar, algoritmul a comis trei greșeli. Cei trei separatori obținuți în cursul aplicării algoritmului *Perceptron* pe șirul de exemple S sunt prezentați în figura de mai jos.



Este imediat că modul în care se construiește separatorul — a cărui expresie analitică este o combinație liniară de instanțe greșit catalogate — este dependent de exemplele prezentate. Întrebarea care a fost pusă în enunț este dacă pentru două succesiuni/ordini diferite, rezultatele finale (i.e., expresiile separatorilor obținuți) coincid (întotdeauna) sau nu. În continuare vom arăta că în mod obișnuit rezultatele finale (i.e., pozițiile finale ale separatorilor obținuți) *nu* coincid.

Considerăm S' secvența de exemple obținută din S prin inversarea [ordinii] exemplelor \bar{x}_1 și \bar{x}_2 . Este imediat că la prima iterație perceptronul greșește și apoi calculează $\bar{w} = (1, 0)$, ceea ce conduce la ecuația separatorului $x_1 = 0$, adică axa Ox_1 . Evident, aceasta dreaptă este un separator perfect pentru întregul set de exemple, deci algoritmul nu va mai comite până la final nicio [altă] greșală. Concluzia este că, pe un același set de exemple de antrenament, algoritmul *Perceptron* poate comite un număr diferit de greșeli (și, de asemenea, un alt separator) dacă exemplele îi sunt prezentate în două succesiuni diferite.

Observație importantă

Deși în ambele cazuri de mai sus (S și S') algoritmul *Perceptron* găsește un separator liniar pentru datele de intrare, acest fapt nu este garantat în cazul general, chiar dacă datele sunt liniar separabile. Vezi de exemplu problema CMU, 2015 spring, Alex Smola, midterm, pr. 4. (Motivele sunt: rata de învățare (implicită) prea mare și/sau ne-parcurerea setului de date decât o singură dată.)

b. The *Perceptron* algorithm classifies data points by their position relative to a hyper-plane. The weight vector, w , learned by Perceptron will be normal to this hyperplane? (True or False?)

Răspuns:

Cunoaștem din geometria analitică faptul că orice doi vectori \bar{w} și \bar{x} pentru care are loc relația $\bar{w} \cdot \bar{x} = 0$ sunt ortogonali (adică, perpendiculari unul pe celălalt).

În cazul perceptronilor (de tip liniar, prag sau sigmoidal), ecuația separatorului este tocmai de forma $\bar{w} \cdot \bar{x} = 0$, unde \bar{x} este un punct arbitrar de pe dreapta [sau, în general, hiper-planul] separator. Considerând \bar{x}_1 și \bar{x}_2 două astfel de puncte, diferența lor, $\bar{x}_1 - \bar{x}_2$, va fi un vector care are direcția dreptei-separator. Din relațiile $\bar{w} \cdot \bar{x}_1 = 0$ și $\bar{w} \cdot \bar{x}_2 = 0$ rezultă $\bar{w} \cdot (\bar{x}_1 - \bar{x}_2) = 0$. Prin urmare, vectorul de ponderi \bar{w} este perpendicular pe direcția dreptei-separator.

Proprietatea aceasta se poate verifica și în mod direct/particular pe datele de la punctele precedente. De exemplu, la punctul a , pe șirul de exemple S , la iterația 1 avem $\bar{w} = (1, -2)$ și se poate observa direct pe grafic că acest vector este perpendicular pe direcția dreptei $y = \frac{1}{2}x$.

Așadar, răspunsul este *Adevărat*.

Rosenblatt' *Perceptron* Algorithm:
Convergence / Mistake bounds

MIT, 2009 fall, Tommy Jaakkola, lecture notes 2

Presupunem că folosim perceptronul de tip prag în varianta Rosenblatt și vrem să învățăm un concept, folosind instanțele de antrenament $x_1, \dots, x_n, \dots \in \mathbb{R}^d$ împreună cu etichetele corespunzătoare $y_1, \dots, y_n, \dots \in \{-1, 1\}$.

Demonstrați că în cazul în care sunt îndeplinite condițiile *i-iv* de mai jos, algoritmul de actualizare a ponderilor perceptronului termină într-un număr finit de pași. Formal, exprimăm acest fapt astfel: $\exists m \in \mathbb{N}$ astfel încât $y_t w^{(m)} \cdot x_t \geq 0$ pentru orice $t \in \{1, \dots, n, \dots\}$, unde $w^{(m)}$ este vectorul de ponderi obținut de perceptron la iterația m .

Iată acum *condițiile* menționate mai sus:

i. Instanțele x_1, \dots, x_n, \dots sunt separabile liniar prin originea sistemului de coordonate, cu o margine finită $\gamma > 0$. Din punct de vedere formal, aceasta înseamnă că există $w^* \in \mathbb{R}^d$ astfel încât $y_t w^* \cdot x_t \geq \gamma$ pentru $t = 1, \dots, n, \dots$

ii. Toate instanțele x_1, \dots, x_n, \dots sunt conținute într-o sferă din \mathbb{R}^d cu centrul în origine, adică $\exists R > 0$ astfel încât $\|x_t\| \stackrel{\text{def.}}{=} \sqrt{x_t \cdot x_t} \leq R$ pentru orice t .

iii. Învățarea se face în manieră *incrementală*, folosind următoarea regulă de actualizare a ponderilor:

$$w^{(k+1)} = w^{(k)} + y_{t_k} x_{t_k} \text{ pentru un } t_k \in \{1, \dots, n, \dots\} \text{ a.î. } y_{t_k} w^{(k)} \cdot x_{t_k} \leq 0, \quad (9)$$

ceea ce înseamnă că instanța x_{t_k} este clasificată eronat de către perceptron la iterația k .

iv. Startarea procesului de învățare se face cu $w^{(0)} = 0 \in \mathbb{R}^d$.

Indicație: Arătați că la fiecare iterație (k) a algoritmului, sunt satisfăcute următoarele proprietăți:

- a. $w^* \cdot w^{(k)} \geq k\gamma$;
- b. $\|w^{(k)}\|^2 \leq kR^2$;
- c. $k \leq \left(\frac{\|w^*\|}{\gamma} R \right)^2$.

Ultima inegalitate indică [o margine superioară pentru] numărul maxim de iterații executate de către perceptron.

Observația 1: Notând cu θ_t unghiul format de vectorii x_t și w^* în \mathbb{R}^d și ținând cont de faptul că

$$\cos \theta_t = \cos(x_t, w^*) = \frac{x_t \cdot w^*}{\|x_t\| \|w^*\|},$$

deci

$$y_t w^* \cdot x_t = y_t \|w^*\| \|x_t\| \cos \theta_t,$$

ceea ce, coroborat cu condiția i , implică $y_t \|x_t\| \cos \theta_t \|w^*\| \geq \gamma$, adică $y_t \|x_t\| \cos \theta_t \geq \frac{\gamma}{\|w^*\|}$. Din punct de vedere geometric, aceasta înseamnă că în spațiul \mathbb{R}^d distanța de la orice instanță x_t la vectorul w^* (care trece prin originea sistemului de coordonate) este mai mare sau egală cu $\frac{\gamma}{\|w^*\|}$.

Observația 2: Separatorul $w^{(k)}$ este o combinație liniară de instanțele x_i , întrucât regula de actualizare este $w^{(k+1)} = w^{(k)} + y_i x_i$. Observația aceasta este valabilă și la antrenarea perceptronului-prag clasic, bazat pe regula delta.

Solution (in Romanian)

Ideea de bază

Algoritmul de actualizare a ponderilor perceptronului determină schimbarea poziției separatorului la fiecare iterație (k) la care avem de a face cu un exemplu clasificat greșit.

Intuitiv, ne așteptăm ca poziția separatorului la iterația $k + 1$ (adică hiperplanul de ecuație $w^{(k+1)} \cdot x = 0$) să se apropie de poziția separatorului w^* care definește conceptul de învățat.

În consecință, cosinusul unghiului dintre $w^{(k)}$ și w^* ar trebui să crească de la o iterație la alta.

Pentru a dovedi/verifica în mod riguros aceasta, vom ține cont că, prin definiție,

$$\cos(w^{(k)}, w^*) = \frac{w^{(k)} \cdot w^*}{||w^{(k)}|| ||w^*||}.$$

Mai întâi vom compara valorile produsului scalar de la numărătorul fracției de mai sus, la două iterații succesive. Folosind regula (9), putem scrie:

$$w^{(k+1)} \cdot w^* = (w^{(k)} + y_{t_k} x_{t_k}) \cdot w^* = w^{(k)} \cdot w^* + y_{t_k} x_{t_k} \cdot w^*.$$

Întrucât $y_{t_k} x_{t_k} \cdot w^* \geq \gamma$ (vezi condiția *i* din enunț), rezultă că valoarea produsului scalar $w^{(k)} \cdot w^*$ crește la fiecare iterație cu o cantitate cel puțin egală cu γ . Cum $w^{(0)} = 0$ conform restricției *iv*, este imediat că $w^{(k)} \cdot w^* \geq k\gamma$ la iterația k . Așadar, numărătorul fracției care definește $\cos(w^{(k)}, w^*)$ crește cel puțin liniar în raport cu k . Cu aceasta, tocmai am demonstrat relația *a*.

A analiza evoluția numitorului fracției care definește $\cos(w^{(k)}, w^*)$ revine la a compara $\|w^{(k+1)}\|$ cu $\|w^{(k)}\|$, întrucât $\|w^*\|$ este constant în raport cu k .

$$\begin{aligned}
\|w^{(k+1)}\|^2 &\stackrel{def.}{=} w^{(k+1)} \cdot w^{(k+1)} \stackrel{(9)}{=} (w^{(k)} + y_{t_k} x_{t_k}) \cdot (w^{(k)} + y_{t_k} x_{t_k}) \\
&= w^{(k)} \cdot w^{(k)} + 2y_{t_k} w^{(k)} \cdot x_{t_k} + y_{t_k}^2 x_{t_k} \cdot x_{t_k} \\
&= \|w^{(k)}\|^2 + 2y_{t_k} w^{(k)} \cdot x_{t_k} + y_{t_k}^2 \|x_{t_k}\|^2 \\
&= \|w^{(k)}\|^2 + 2y_{t_k} w^{(k)} \cdot x_{t_k} + \|x_{t_k}\|^2 \\
&\leq \|w^{(k)}\|^2 + R^2.
\end{aligned}$$

Inegalitatea de mai sus derivă din faptul că $y_{t_k} w^{(k)} \cdot x_{t_k} \leq 0$ (i.e., exemplul t_k este clasificat greșit la iterația k , conform condiției *iii* din enunț) și din ipoteza că orice instanță de antrenament este conținută în sfera de rază R și având centrul în originea spațiului \mathbb{R}^d , conform condiției *ii*. Cum $w^{(0)} = 0$, este imediat că $\|w^{(k)}\|^2 \leq k R^2$, deci $\|w^{(k)}\| \leq \sqrt{k} R$ la iterația k . Cu aceasta, am demonstrat relația *b*.

Acum, combinând rezultatele a și b , obținem următoarea inegalitate:

$$\cos(w^{(k)}, w^*) = \frac{w^{(k)} \cdot w^*}{\|w^{(k)}\| \|w^*\|} \geq \frac{k\gamma}{\sqrt{k} R \|w^*\|} = \sqrt{k} \frac{\gamma}{R \|w^*\|}.$$

Știm că valoarea funcției \cos este întotdeauna cel mult egală cu 1. În consecință,

$$1 \geq \cos(w^{(k)}, w^*) \geq \sqrt{k} \frac{\gamma}{R \|w^*\|} \Rightarrow k \leq \left(R \frac{\|w^*\|}{\gamma} \right)^2,$$

deci am demonstrat relația c .

Aceasta înseamnă că în condițiile specificate în enunț, antrenarea perceptronului-prag se va face în cel mult $\left\lfloor \left(R \frac{\|w^*\|}{\gamma} \right)^2 \right\rfloor$ iterații, unde perechea de simboluri $\lfloor \rfloor$ desemnează funcția parte întreagă inferioară.

Observația 3: Marginea superioară calculată mai sus este remarcabilă, întrucât ea nu depinde nici de instanțele x_i și nici de dimensiunea (d) a spațiului din care sunt selectate aceste instanțe.

Observația 4

Restricția referitoare la separabilitatea prin origine a instanțelor de antrenament (vezi [condiția i](#), prima parte) — și anume, componenta „liberă” w_0 asociată separatorului w^* trebuie să fie 0 — nu este de fapt limitativă. Cazul general al separabilității liniare în \mathbb{R}^d , adică $y_t(w_0^* + w^* \cdot x_t) \geq 0$ pentru orice $t \in \{1, 2, \dots, n \dots\}$, poate fi redus la cazul particular al separabilității prin origine în \mathbb{R}^{d+1} dacă pe de o parte se consideră $w' = (w^0, w^1, \dots, w^d)$, cu $w^* = (w^1, \dots, w^d)$, iar pe de altă parte se mapează toate instanțele de antrenament $x_t = (x_t^1, \dots, x_t^d)$ din \mathbb{R}^d în \mathbb{R}^{d+1} astfel: $x'_k = (x_t^0, x_t^1, \dots, x_t^d)$, cu $x_t^0 = 1$ pentru orice t . Cu această mapare, rezultă $y_t w' \cdot x'_t \geq 0$ pentru orice t (respectiv $y_t w' \cdot x'_t \geq \gamma$ când se lucrează cu margine finită, ca în enunțul acestei probleme).

Nici [restricția iv](#) ($w^{(0)} = 0$) nu este cu adevărat limitativă; în general algoritmi de antrenare a unităților / rețelelor neuronale inițializează ponderile w la valori mici (în modul).

Similar, este imediat că restricția potrivit căreia sfera în care sunt conținute instanțele x_1, \dots, x_n, \dots trebuie să aibă centrul în originea lui \mathbb{R}^d (vezi [condiția ii](#), partea a doua) poate fi eliminată fără ca rezultatul de convergență să fie afectat.

Observația 5

Deducerea marginii superioare de la punctul c de mai sus în cazul în care se folosește o rată de învățare oarecare $\eta > 0$ se face extinzând în mod natural demonstrația de mai sus (vezi CMU, 2013 spring, A. Smola, B. Póczos, HW2, pr. 2.a).

În concluzie, mai rămân de examinat **două condiții**: separabilitatea liniară cu margine $\gamma > 0$ și conținerea instanțelor de antrenament într-o sferă de rază finită. Se poate demonstra că ambele condiții **sunt esențiale** pentru convergența perceptronului Rosenblatt în regim de învățare online (vezi problema CMU, 2008 fall, Eric Xing, midterm exam, pr. 3.2).

Exemplifying

The gradient descent method:

finding the minimum of a real function of second degree

University of Utah, 2008 fall, Hal Daumé III, HW4, pr. 1

Suppose we are trying to find the minimum of the function $f(x) = 3x^2 - 2x + 1$, for uni-variate (scalar) x .

First, verify that this function is convex (and therefore it has a *global* minimum).

Second, find the minimum of this function using calculus.

Finally, perform (by hand – show your work) three steps of gradient descent with $\eta = 0.1$ and the initial point $x_0 = 1$. How close does it get to the true solution?

Solution (in Romanian)

Pentru a studia convexitatea funcției $f(x)$ se calculează derivata a doua:

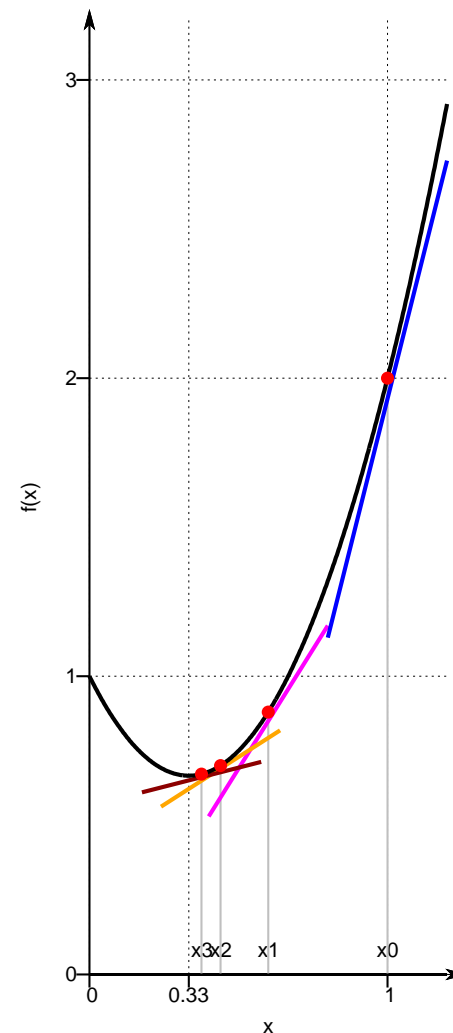
$$f''(x) = 6 > 0, \forall x \in \mathbb{R}.$$

Rezultă că funcția f este convexă pe întreg domeniul ei de definiție, deci are un (singur) punct de minim.

Pentru a calcula minimul funcției $f(x) = 3x^2 - 2x + 1$ utilizăm derivata de ordinul întâi:

$$f'(x) = 6x - 2.$$

Punctul de minim este dat de soluția ecuației $f'(x) = 0$, și anume: $x = \frac{1}{3} \approx 0.33$.



Observații

Se constată ușor că

1. apropierea de punctul de optim este [mai] rapidă atâta timp cât valoarea primei derivate (i.e., panta tangentei la graficul funcției) este mare în valoare absolută;

2. dacă rata de învățare η a fost fixată la o valoare prea mare, atunci este posibil ca la un moment dat să depășim punctul de optim și apoi să „pendulăm” în jurul lui. Acest *punct slab* al metodei gradientului poate fi contracarat reducând în mod dinamic mărimea lui η .

Altminteri metoda gradientului are *avantajul* de a fi o *tehnică de optimizare* foarte simplă din punct de vedere conceptual și ușor de implementat.

Alte două *puncte slabe* ale metodei gradientului descendent sunt:
imposibilitatea de a garanta găsirea optimul global și
numărul mare de iterații care trebuie executate pe unele seturi de date reale.

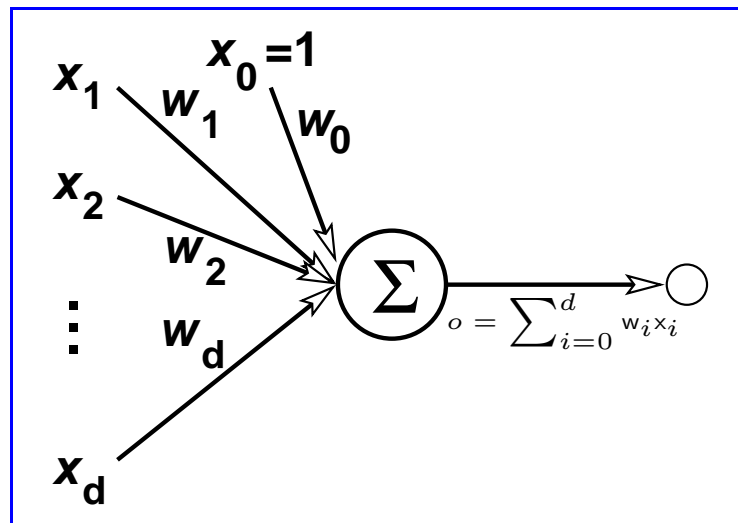
Theoretical foundations for the perceptron training algorithm

Liviu Ciortuz, 2017

following Tom Mitchell, *Machine Learning* book, 1997, p. 89-93

În acest exercițiu vom elabora/reda [mai întâi] partea de fundamentare teoretică pentru antrenarea perceptronului liniar, date fiind instanțele $(\bar{x}_1, t_1), \dots, (\bar{x}_n, t_n)$, cu $\bar{x}_i \in \mathbb{R}^d$ și $t_i \in \mathbb{R}$ pentru $i = 1, \dots, n$.

Așadar, vom considera perceptronul liniar având intrările $x_0 = 1, x_1, \dots, x_d$ și ponderile w_0, w_1, \dots, w_d .



a. În linie cu *metoda gradientului descendent*, deduceți care este forma *regulii de actualizare* a vectorului de ponderi $\bar{w} \in \mathbb{R}^{d+1}$ (sau, echivalent, forma regulilor de actualizare a ponderilor w_i , cu $i = 0, \dots, d$) pentru găsirea acelei valori care minimizează semi-suma pătratelor erorilor comise de perceptron,^a adică

$$E(\bar{w}) \stackrel{\text{def.}}{=} \frac{1}{2} \sum_{i=1}^n (t_i - o_i)^2.$$

În această ultimă expresie, o_i este output-ul perceptronului liniar pentru intrarea $\bar{x}_i \stackrel{\text{renot.}}{=} (x_0 = 1, x_{i,1}, \dots, x_{i,d})$, adică $o_i = w_0 + \sum_{j=1}^d w_j x_{i,j}$.

Vă reamintim că la aplicarea metodei gradientului descendent pentru găsirea unui punct de minim al unei funcții derivabile $f : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$, regula de actualizare a parametrilor \bar{w} ai funcției f are forma $\bar{w} \leftarrow \bar{w} + \Delta \bar{w}$, cu $\Delta \bar{w} \stackrel{\text{def.}}{=} -\eta \nabla_{\bar{w}} f(\bar{w})$, unde $\eta > 0$ este un număr real mic, numit *rata de învățare*, iar $\nabla_{\bar{w}} f(\bar{w})$ este vectorul *gradient*, adică $\left(\frac{\partial}{\partial w_0} f(\bar{w}), \frac{\partial}{\partial w_1} f(\bar{w}), \dots \right)$.

^aConform problemei CMU, 2001 fall, T. Mitchell, A. Moore, final exam, pr. 10.2, am putea scrie — în ipoteza că datele au fost generate probabilist, cu o componentă „zgomot” urmând o distribuție gaussiană uni-variată de medie 0, așa cum s-a precizat acolo —,

$$\bar{w}_{\text{ML}} = \arg \min_{\bar{w}} E(\bar{w}) = \arg \min_{\bar{w}} \frac{1}{2} \sum_{i=1}^n (t_i - \bar{w} \cdot \bar{x}_i)^2.$$

Answer

$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{j=1}^n (t_j - o_j)^2 = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial w_i} (t_j - o_j)^2 \\
 &= \frac{1}{2} \sum_{j=1}^n 2(t_j - o_j) \frac{\partial}{\partial w_i} (t_j - o_j) = \sum_{j=1}^n (t_j - o_j) \frac{\partial}{\partial w_i} (t_j - \bar{w} \cdot \bar{x}_j) \\
 &= \sum_d (t_j - o_j) (-x_{j,i})
 \end{aligned}$$

Therefore,

$$\Delta w_i = \eta \sum_{j=1}^n (t_j - o_j) x_{j,i}$$

b. Dacă în locul perceptronului liniar vom considera perceptronul-prag, respectiv perceptronul sigmoidal, prin ce va diferi forma regulii de actualizare a ponderilor față de rezultatul de la punctul a? (Pentru perceptronul-prag, se va presupune că $t_i \in \{-1, +1\}$, în vreme ce pentru perceptronul sigmoidal vom considera $t_i \in \{0, 1\}$, pentru $i = 1, \dots, n$.)

Answer:

$$\begin{aligned}
 \frac{\partial E}{\partial w_i} &= \frac{\partial}{\partial w_i} \frac{1}{2} \sum_{j=1}^n (t_j - o_j)^2 = \frac{1}{2} \sum_{j=1}^n \frac{\partial}{\partial w_i} (t_j - o_j)^2 = \frac{1}{2} \sum_{j=1}^n 2(t_j - o_j) \frac{\partial}{\partial w_i} (t_j - o_j) \\
 &= \sum_{j=1}^n (t_j - o_j) \frac{\partial}{\partial w_i} (t_j - \underbrace{\sigma(\bar{w} \cdot \bar{x}_j)}_{o_j}) = - \sum_{j=1}^n (t_j - o_j) \frac{\partial}{\partial w_i} \sigma(\bar{w} \cdot \bar{x}_j) \\
 &= - \sum_{j=1}^n (t_j - o_j) \sigma(\bar{w} \cdot \bar{x}_j) (1 - \sigma(\bar{w} \cdot \bar{x}_j)) \frac{\partial}{\partial w_i} \bar{w} \cdot \bar{x}_j = - \sum_{j=1}^n (t_j - o_j) o_j (1 - o_j) \frac{\partial}{\partial w_i} \bar{w} \cdot \bar{x}_j \\
 &= - \sum_{j=1}^n (t_j - o_j) o_j (1 - o_j) x_{j,i}
 \end{aligned}$$

Note: Here above we used the fact that

$$\sigma'(z) = \left(\frac{1}{1 + e^{-z}} \right)' = - \frac{(1 + e^{-z})'}{(1 + e^{-z})^2} = \frac{e^{-z}}{(1 + e^{-z})^2} = \frac{1 + e^{-z} - 1}{(1 + e^{-z})^2} = \frac{1}{1 + e^{-z}} - \frac{1}{(1 + e^{-z})^2} = \sigma(z)(1 - \sigma(z)), \forall z \in \mathbb{R}$$

c. Elaborați algoritmul de antrenare a perceptronului liniar (și apoi, dacă doriți, și a celui sigmoidal, dar veți preciza doar diferențele!). La *inițializare*, ponderile w_i vor primi ca valori numere reale mici. În ce privește *condiția de oprire* a algoritmului, se poate opta pentru una din următoarele variante (eventual combinate):

- toate instanțele de antrenament sunt corect clasificate (în ipoteza că datele \bar{x}_i , cu $i = 1, \dots, n$, sunt liniar separabile);
- efectuarea unui număr prestabilit de iterații;
- verificarea condiției $|E(\bar{w}^{(t+1)}) - E(\bar{w}^{(t)})| < \varepsilon$, unde pragul numeric $\varepsilon > 0$ este stabilit inițial, iar $E(\bar{w}^{(t)})$ este valoarea *funcției de eroare* E (i.e., ceea ce mai sus am numit semi-suma pătratelor erorilor) pe setul de date de antrenament, la finalul iterației t .

The gradient descent algorithm for the linear unit

GRADIENT-DESCENT(*training_examples*, η)

where each *training example* is a pair of the form $\langle \bar{x}, t \rangle$, with

\bar{x} – the vector of input values

t – the target output value.

η is the learning rate (e.g., .05).

- initialize each w_i to some small random value
- until the *termination condition* is met

initialize each Δw_i to zero;

for each $\langle \bar{x}, t \rangle$ in *training_examples*,

input the instance \bar{x} to the unit and compute the output o ;

for each linear unit weight w_i ,

$$\Delta w_i \leftarrow \Delta w_i + \eta(t - o)x_i;$$

for each linear unit weight w_i ,

$$w_i \leftarrow w_i + \Delta w_i$$

d. Ce cunoașteți (de la curs) despre convergența algoritmului de la punctul c? Dar despre convergența perceptronului-prag?

Answer:

[Hertz et al., 1991] has proven that the the gradient descent-based training algorithm [centered on the weight updating rule] for the **linear unit** is guaranteed to converge to a hypothesis that minimizes the sum of squared errors

- given a sufficiently small learning rate η ,
- even when the training data contains noise,
- even when the training data is not separable by H .

Prior to that, it was proven by [Minsky & Papert, 1969] that [the same algorithm, applied to] the **threshold perceptron** converges

- if the training data is linearly separable,
- and η is sufficiently small.

Using another *loss-function*
in conjunction with the linear perceptron
and the gradient descent method:

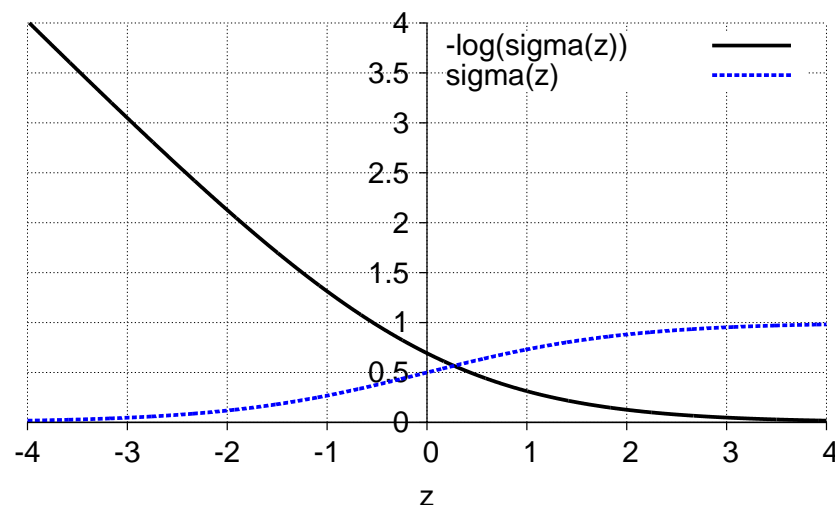
The log-sigmoidal function;
deduction of the *update rules*

University of Utah, 2008 fall, Hal Daumé III, HW4, pr. 2-4

The 0/1 loss function is hard to optimize using gradient descent because it is poorly behaved (discontinuous, non-differentiable, etc.). A common “smooth” version of 0/1 loss is the *sigmoid* loss, which makes use of our favorite function, the sigmoid: $\sigma(z) = 1/(1 + \exp(-z))$.

In this exercise we will work with another loss function, the log-sigmoid:

$$\ell(y, f(x)) = -\log(\sigma(yf(x))).$$



- a. First, verify that this is a reasonable loss function: when $f(x)$ is correct, the loss is low, and when $f(x)$ is incorrect, the loss is high. You may assume that $f(x)$ produces *real values*, where positive values mean class +1 and negative values mean class -1.
- b. Consider optimizing a linear function under log-sigmoid loss, so that we have $f(\bar{w}) = \sum_n -\log(\sigma(y_n(\bar{w} \cdot \bar{x}_n + b)))$. Compute the gradient of this function with respect to \bar{w} and with respect to b so that we might construct a gradient descent algorithm.
- c. Show that the loss function from the previous two questions is convex in both \bar{w} and b .
- d. Finally, state the *update rules* for both w_i and b which serve for training a linear perceptron, based on the gradient descent method in conjunction with the log-sigmoid loss function.

Solution (in Romanian)

a. Funcția l se poate scrie sub forma:

$$l(y, h(x)) = -\ln \sigma(y h(x)) = -\ln \frac{1}{1 + e^{-y h(x)}} = \ln(1 + e^{-y h(x)})$$

Din proprietățile logaritmului știm că $\ln 1 = 0$, $\ln a < 0$ dacă $0 < a < 1$, și $\ln a > 0$ dacă $a > 1$.

Dacă ipoteza $h(x)$ este corectă, adică are același semn cu y , atunci $e^{-y h(x)}$ este o valoare mică (și cu atât mai mică cu cât $|h(x)|$ este mai mare), deci și valoarea funcției $l(y, h(x))$ este mică.

Dacă ipoteza $h(x)$ este incorectă, adică are semn contrar semnului lui y , atunci $e^{-y h(x)}$ este o valoare mare (și cu atât mai mare cu cât $|h(x)|$ este mai mare), deci și valoarea funcției $l(y, h(x))$ este mare.

De *exemplu*, dacă $y = 1$ și $h(x) = -1$ (sau $y = -1$ și $h(x) = 1$) atunci pierderea/costul este $-\ln \frac{1}{1+e} \approx 1.31$, iar dacă $y = h(x) = 1$ (sau, ambele, -1) atunci costul este $-\ln \frac{1}{1+e^{-1}} = -\ln \frac{e}{1+e} \approx 0.31$.

b. Calculăm derivatele parțiale ale funcției $f(b, \bar{w}) = \sum_n -\ln \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))$ în raport cu w_i și respectiv b . Vom folosi *proprietatea* $\sigma'(x) = \sigma(x)(1 - \sigma(x))$.

$$\begin{aligned}
 \frac{\partial f}{\partial w_i} &= \frac{\partial}{\partial w_i} \sum_n -\ln \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) = - \sum_n \frac{\partial}{\partial w_i} \ln \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) \\
 &= - \sum_n \frac{1}{\sigma(y_n(\bar{w} \cdot \bar{x}_n + b))} \frac{\partial}{\partial w_i} \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) \\
 &= - \sum_n \frac{1}{\sigma(y_n(\bar{w} \cdot \bar{x}_n + b))} \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))(1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n x_{n,i} \\
 &= - \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n x_{n,i}.
 \end{aligned}$$

Similar, vom obține $\frac{\partial f}{\partial b} = - \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n$.

c. Funcția $f(b, \bar{w})$ este convexă în raport cu variabila w_i dacă $\frac{\partial^2 f}{\partial w_i^2} \geq 0$. Similar,

$f(b, \bar{w})$ este convexă în raport cu b dacă $\frac{\partial^2 f}{\partial b^2} \geq 0$. Așadar, trebuie să calculăm aceste derivate parțiale de ordinul 2. Pentru aceasta, vom folosi derivatele parțiale de ordinul întâi care au fost calculate la punctul precedent.

$$\begin{aligned}
 \frac{\partial^2 f}{\partial w_i^2} &= \frac{\partial}{\partial w_i} \left(\frac{\partial f}{\partial w_i} \right) = \frac{\partial}{\partial w_i} \left(- \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n x_{n,i} \right) \\
 &= \sum_n y_n x_{n,i} \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) \frac{\partial}{\partial w_i} (y_n(\bar{w} \cdot \bar{x}_n + b)) \\
 &= \sum_n y_n x_{n,i} \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n x_{n,i} \\
 &= \sum_n \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n^2 x_{n,i}^2
 \end{aligned}$$

Cum $\sigma(x) \in (0, 1)$ pentru orice $x \in \mathbb{R}$, rezultă că și $1 - \sigma(x) \in (0, 1) \forall x$. Așadar, $\sigma(x) > 0$, $1 - \sigma(x) > 0$ și, evident, $y_n^2 \geq 0$, deci $\frac{\partial^2 f}{\partial b^2} \geq 0$, adică f este convexă în w_i , pentru orice i .

Analog, calculăm derivata parțială de ordin secund a lui f în raport cu b :

$$\begin{aligned}
 \frac{\partial^2 f}{\partial b^2} &= \frac{\partial}{\partial b} \left(\frac{\partial f}{\partial b} \right) = \frac{\partial}{\partial b} \left(- \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n \right) \\
 &= \sum_n y_n \frac{\partial}{\partial b} \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) \\
 &= \sum_n y_n \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) \frac{\partial}{\partial b} (y_n(\bar{w} \cdot \bar{x}_n + b)) \\
 &= \sum_n y_n \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n \\
 &= \sum_n \sigma(y_n(\bar{w} \cdot \bar{x}_n + b)) (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n^2 \geq 0
 \end{aligned}$$

Așadar, f este convexă și în raport cu variabila b .

Observație: La curs, b a fost asimilat cu w_0 , iar $x_{n,0} = 1$ prin definiție. Dacă s-ar fi procedat la fel și aici, atunci nu am fi avut de calculat decât o derivată la punctul a și una la punctul b .

d. În *concluzie*, valoarea optimă a funcției f este un minim, iar algoritmul de învățare automată bazat pe metoda gradientului descendent va găsi (la limită, eventual) acest minim. Antrenarea acestui perceptron se face similar cu antrenarea unității liniare care a fost prezentată la curs, folosind reguli de actualizare de forma:

$$w_i \leftarrow w_i + \Delta w_i \text{ și } b \leftarrow w_i + \Delta b.$$

În cazul de față, Δw_i și Δb sunt definite de expresiile

$$\Delta w_i = -\eta \frac{\partial f}{\partial w_i} = \eta \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_n + b))) y_n \bar{x}_{n,i}$$

și respectiv

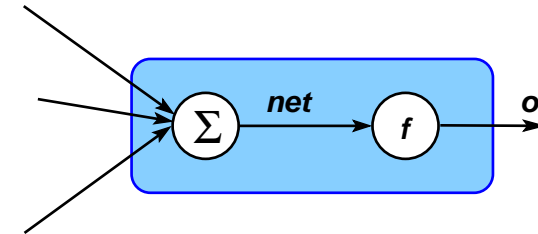
$$\Delta b = -\eta \frac{\partial f}{\partial b} = \eta \sum_n (1 - \sigma(y_n(\bar{w} \cdot \bar{x}_{n,i} + b))) y_n$$

Algoritmul de retro-propagare:
deducerea regulilor de actualizare a ponderilor
pentru o rețea feed-forward cu 2 niveluri,
folosind o funcție de activare oarecare (derivabilă) f

CMU, 2008 fall, Eric Xing, HW2, pr. 2.2

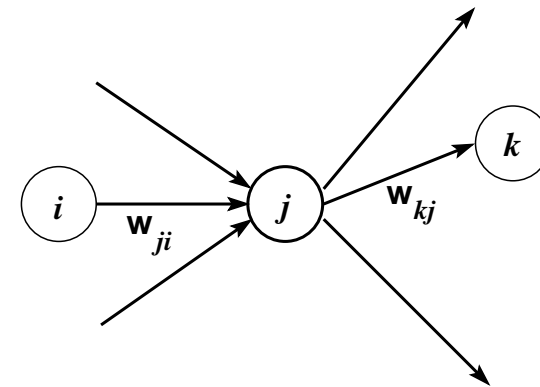
La acest exercițiu veți deriva regula de actualizare pentru algoritmul de retro-propagare (variantea stochastică) pentru o rețea neuronală artificială de tip “feed-forward”, cu două niveluri de unități sigmoidale. Se consideră D unități de intrare, H unități ascunse și K unități de ieșire.

Funcția de eroare cu care se lucrează este $E(\bar{w}) = \frac{1}{2} \sum_k (t_k - o_k)^2$, unde $o_k = f(net_k) = \frac{1}{1 + e^{-net_k}}$, iar c este o constantă.



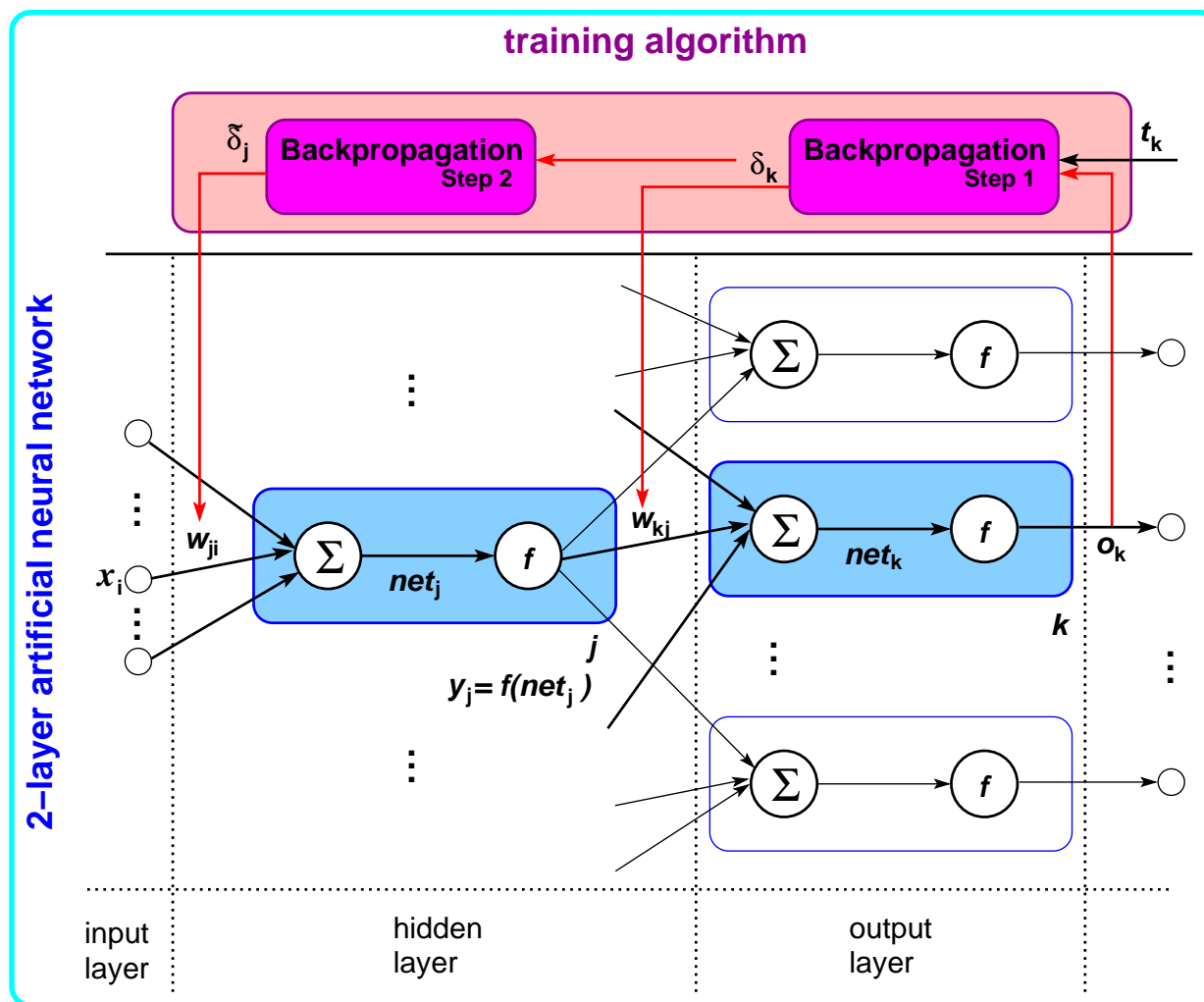
a. Fie w_{kj} ponderea conexiunii de la unitatea ascunsă j către unitatea de ieșire k . Arătați că regula de actualizare pentru w_{kj} este de forma $w_{kj} \leftarrow w_{kj} + \eta \delta_k y_j$, unde y_j este ieșirea unității ascunse j , iar $\delta_k = f'(net_k)(t_k - o_k)$, cu $net_k = \sum_{j'} w_{kj'} y_{j'}$.

b. Arătați că regulile de actualizare pentru ponderile care corespund conexiunilor input-to-hidden sunt de forma $w_{ji} \leftarrow w_{ji} + \eta \tilde{\delta}_j x_i$, unde x_i este intrarea i , iar $\tilde{\delta}_j = f'(net_j) [\sum_{k=1}^K w_{kj} \delta_k]$, cu $net_j = \sum_{i'} w_{ji'} x_{i'}$.



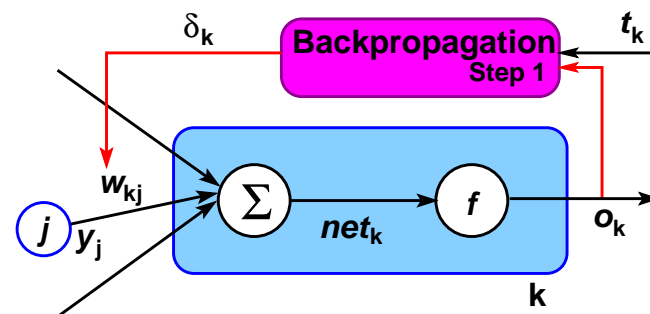
Răspuns

Observație: Vom face rezolvarea problemei în raport cu o funcție oarecare f derivabilă, lăsată nespecificată. Forma particulară dată în enunț pentru funcția f nu are nicio consecință particulară asupra raționamentului dezvoltat în continuare.



a. Conform metodei gradientului descendent,

$w_{kj} \leftarrow w_{kj} - \eta \frac{\partial E}{\partial w_{kj}}$, unde η este o constantă pozitivă (rata de învățare).



Intuitiv, ponderea w_{kj} influențează valoarea funcției de eroare E (doar) prin intermediul lui net_k . (S-a notat cu net_k suma $\sum_{j'} w_{kj'} y_{j'}$.) Așadar, este necesară aplicarea formulei pentru derivarea unei compuneri de funcții derivabile:

$$\begin{aligned} \frac{\partial E}{\partial w_{kj}} &= \frac{\partial E}{\partial net_k} \frac{\partial net_k}{\partial w_{kj}} = \left(\frac{\partial}{\partial net_k} \frac{1}{2} \sum_{k'=1}^K (t_{k'} - o_{k'})^2 \right) \left(\frac{\partial}{\partial w_{kj}} \sum_{j'=0}^{n_k} w_{kj'} y_{j'} \right) = \left(\frac{\partial}{\partial net_k} \frac{1}{2} (t_k - f(net_k))^2 \right) y_j \\ &= y_j \left(\frac{1}{2} 2(t_k - f(net_k)) \frac{\partial}{\partial net_k} (t_k - f(net_k)) \right) = -y_j (t_k - f(net_k)) f'(net_k) \end{aligned}$$

unde n_k este numărul de intrări ale unității k . În particular, $n_k = H$ dacă se consideră toate conexiunile posibile (hidden-to-output).

Dacă notăm $\delta_k = (t_k - f(net_k)) f'(net_k)$, atunci avem $\frac{\partial E}{\partial w_{kj}} = -\delta_k y_j$, și

$$w_{kj} \leftarrow w_{kj} - \eta \frac{\partial E}{\partial w_{kj}} = w_{kj} + \eta \delta_k y_j$$

b. Ca și la punctul precedent, $w_{ji} \leftarrow w_{ji} - \eta \frac{\partial E}{\partial w_{ji}}$. Ponderea w_{ji} influențează valoarea funcției de eroare E (doar) prin intermediul lui net_j . Folosim din nou formula pentru derivarea unei compuneri de funcții derivabile și ținem cont că $net_j = \sum_{i'} w_{ji'} x_{i'}$:

$$\begin{aligned} \frac{\partial E}{\partial w_{ji}} &= \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ji}} = \left(\frac{\partial}{\partial net_j} \frac{1}{2} \sum_{k=1}^K (t_k - o_k)^2 \right) \left(\frac{\partial}{\partial w_{ji}} \sum_{i'=0}^D w_{ji'} x_{i'} \right) \\ &= \left(\frac{1}{2} \sum_{k=1}^K 2(t_k - o_k) \frac{\partial}{\partial net_j} (t_k - o_k) \right) x_i = -x_i \sum_{k=1}^K (t_k - o_k) \frac{\partial o_k}{\partial net_j} \end{aligned}$$

Pe de altă parte, valoarea net_j influențează valoarea o_k (doar) prin intermediul lui net_k . Așadar,

$$\begin{aligned} \frac{\partial o_k}{\partial net_j} &= \frac{\partial o_k}{\partial net_k} \frac{\partial net_k}{\partial net_j} = f'(net_k) \frac{\partial}{\partial net_j} \sum_{j'=0}^{n_k} w_{kj'} y_{j'} \\ &= f'(net_k) \frac{\partial}{\partial net_j} \sum_{j'=0}^{n_k} w_{kj'} f(net_{j'}) = f'(net_k) f'(net_j) w_{kj} \end{aligned}$$

Înlocuind acest rezultat în egalitatea precedentă, vom avea:

$$\begin{aligned}\frac{\partial E}{\partial w_{ji}} &= -x_i \sum_{k=1}^K (t_k - o_k) f'(net_k) f'(net_j) w_{kj} \\ &= -x_i f'(net_j) \sum_{k=1}^K (t_k - o_k) f'(net_k) w_{kj} = -x_i f'(net_j) \sum_{k=1}^K \delta_k w_{kj}\end{aligned}$$

Folosind notația $\tilde{\delta}_j = f'(net_j) \sum_{k=1}^K \delta_k w_{kj}$, egalitatea de mai sus devine

$\frac{\partial E}{\partial w_{ji}} = -x_i \tilde{\delta}_j$, iar regula de actualizare a ponderii se transformă în:

$$w_{ji} \leftarrow w_{ji} + \eta \tilde{\delta}_j x_i$$

Observație

În mod similar cu demonstrația din cartea *Machine Learning* de Tom Mitchell, pag. 101-103, demonstrația de aici poate fi extinsă în mod facil la rețele neuronale feed-forward cu un număr oarecare de niveluri ascunse.

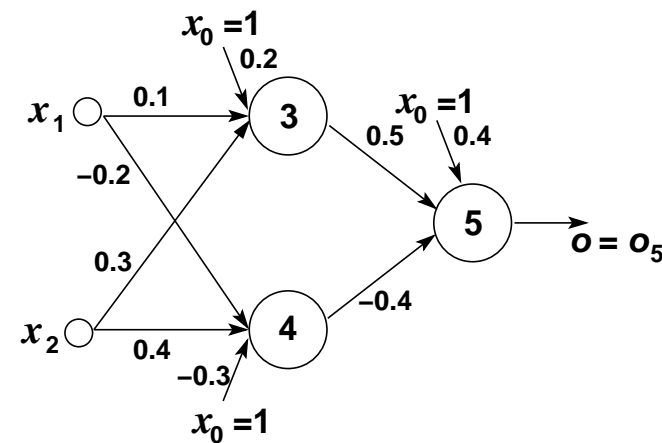
De asemenea, mai este posibilă *încă o generalizare*: funcția de activare f poate fi diferită de la o unitate la alta (sigmoidală, generalizat-sigmoidală (cu o anumită constantă c fixată), liniară etc). În demonstrație va fi suficient să atașăm simbolului f indicele unității neuronale respective.

**Aplicarea algoritmului de retro-propagare
pe o rețea feed-forward de unități sigmoidale
cu 2 niveluri**

prelucrare de Liviu Ciortuz, după un exemplu din
„Apprentissage artificiel“, A. Cornuéjols, L. Miclet, 2010, pag. 332, 341

Rețeaua neuronală din figura alăturată este formată din unități sigmoide.

a. Care este rezultatul produs de această rețea pentru intrarea $x \stackrel{not.}{=} (x_1, x_2) = (1, 1)$?



b. Executați manual prima iterație a algoritmului de retro-propagare pe rețeaua dată. Presupunem că în cazul intrării $x = (1, 1)$ ieșirea produsă de rețea ar trebui să fie $t = 0$. Luând rata de învățare $\eta = 1$, precizați care vor fi

- valorile ponderilor $w_{30}, w_{31}, w_{32}, w_{40}, w_{41}, w_{42}, w_{50}, w_{51}, w_{52}$, după aplicarea acestei prime iterații în antrenarea rețelei;^a
- noul output produs de rețea (după actualizarea ponderilor) pe aceeași intrare $x = (1, 1)$.

Sugestie: Pentru a vedea detaliile algoritmului de retro-propagare a erorii, consultați cartea *Machine Learning* de Tom Mitchell, pag. 98.

c. Comparați rezultatele de la punctele a și b. Ce constatați?

^aSemnificația pentru w_{ji} este cea din cartea *Machine Learning* a lui Tom Mitchell, la pag. 98: w_{ji} este ponderea de pe arcul/legătura de la unitatea (sau intrarea) i către unitatea j .

Solution (in Romanian)

a.

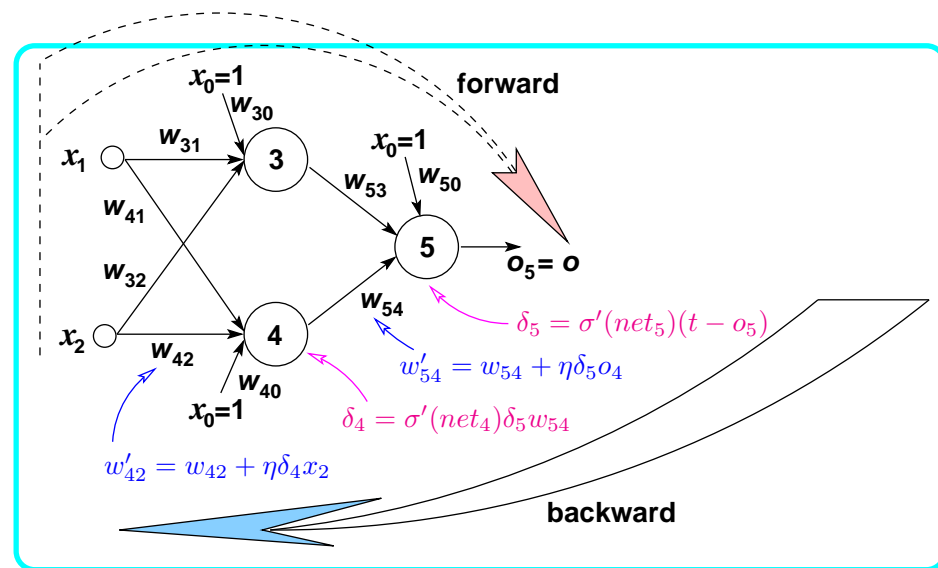
i	$net_i = \sum_j w_{ij}x_j$	$o_i = \sigma(net_i)$
3	$0.2 + 0.1 + 0.3 = 0.6$	$1/(1 + e^{-0.6}) \simeq 0.646$
4	$-0.3 - 0.2 + 0.4 = -0.1$	$1/(1 + e^{0.1}) \simeq 0.475$
5	$0.4 + 0.5 \cdot 0.646 - 0.4 \cdot 0.475 = 0.533$	$1/(1 + e^{-0.533}) \simeq 0.630$

Așadar, output-ul produs de rețea este 0.63.

b.

Precizare:

Ca să facilităm / sintetizăm înțelegerea modului în care se aplică acest algoritm, în schema alăturată am arătat cum anume vor fi calculate noile valori ale ponderilor w , precum și cantitățile δ implicate în execuția unei iterații a algoritmului. Ilustrarea se rezumă la parcurgerea (în ordinea output-to-input) unuia dintre drumurile din această rețea, și anume drumul 5, 4, 2. Extensia la celelalte căi este facilă.



În această imagine (dar numai aici) am notat cu w' noile valori pentru ponderi (calculate în ultimul pas al iterației), pentru a nu fi confundate cu vechile valori, care intervin în calculul cantităților δ .

Conform algoritmului de retro-propagare,

$$\begin{aligned}\delta_5 &\stackrel{\text{def.}}{=} -\frac{\partial E}{\partial \text{net}_5} = \sigma'(\text{net}_5)(t - o) = \sigma(\text{net}_5)(1 - \sigma(\text{net}_5))(t - o) \\ &= o(1 - o)(t - o) = 0.63(1 - 0.63)(0 - 0.63) = -0.147\end{aligned}$$

Output-ul neuronului 3 constituie unul din input-urile neuronului 5.

În notația folosită de Tom Mitchell, vom scrie $\text{Downstream}(3) = \{5\}$.

Așadar,

$$\delta_3 \stackrel{\text{def.}}{=} -\frac{\partial E}{\partial \text{net}_3} = o_3 \cdot (1 - o_3) \cdot \delta_5 \cdot w_{53} = 0.646 \cdot (1 - 0.646) \cdot (-0.147) \cdot 0.5 \simeq -0.017$$

În mod similar, $\text{Downstream}(4) = \{5\}$, și

$$\delta_4 \stackrel{\text{def.}}{=} -\frac{\partial E}{\partial \text{net}_4} = o_4 \cdot (1 - o_4) \cdot \delta_5 \cdot w_{54} = 0.475 \cdot (1 - 0.475) \cdot (-0.147) \cdot (-0.4) \simeq 0.015$$

Întrucât am terminat de calculat toate cantitățile de tip δ_j , vom trece acum la actualizarea ponderilor rețelei. Mai întâi vom calcula noile valori pentru ponderile de pe conexiunile hidden-to-output, deci vom avea:

$$w_{5j} \leftarrow w_{5j} + \Delta w_{5j} \text{ cu } \Delta w_{5j} = \eta \delta_5 o_j = \delta_5 o_j \text{ pentru } j \in \{0, 3, 4\}.$$

Așadar,

$$\begin{cases} \Delta w_{50} = -0.147 \cdot 1 = -0.147 \\ \Delta w_{53} = -0.147 \cdot 0.646 \simeq -0.095 \\ \Delta w_{54} = -0.147 \cdot 0.475 \simeq -0.070 \end{cases} \Rightarrow \begin{cases} w_{50} \leftarrow 0.4 - 0.147 \simeq 0.253 \\ w_{53} \leftarrow 0.5 - 0.1 = 0.405 \\ w_{54} \leftarrow -0.4 - 0.070 = -0.470 \end{cases}$$

Apoi vom actualiza ponderile care corespund conexiunilor de tip input-to-hidden.

Știm că $\Delta w_{3i} = \eta \delta_3 x_i$ pentru $i \in \{0, 1, 2\}$, așadar input-ul $x_0 = x_1 = x_2 = 1$ va implica $\Delta w_{30} = \Delta w_{31} = \Delta w_{32} = 1 \cdot (-0.017) \cdot 1 = -0.017$.

În consecință,

$$\begin{cases} w_{30} \leftarrow 0.2 - 0.017 = 0.183 \\ w_{31} \leftarrow 0.1 - 0.017 = 0.083 \\ w_{32} \leftarrow 0.3 - 0.017 = 0.283 \end{cases}$$

Similar, fiindcă $\Delta w_{4i} = \eta \delta_4 x_i$ pentru $i \in \{0, 1, 2\}$, rezultă că $\Delta w_{40} = \Delta w_{41} = \Delta w_{42} = 1 \cdot 0.015 \cdot 1 = 0.015$ și, prin urmare:

$$\begin{cases} w_{40} \leftarrow -0.3 + 0.015 = -0.285 \\ w_{41} \leftarrow -0.2 + 0.015 = -0.185 \\ w_{42} \leftarrow 0.4 + 0.015 = 0.415 \end{cases}$$

Acum putem calcula noua valoare produsă de rețeaua neuronală:

i	$net_i = \sum_j w_{ij}x_j$	$o_i = \sigma(net_i)$
3	$0.183 + 0.083 + 0.283 = 0.549$	$1/(1 + e^{-0.549}) \simeq 0.634$
4	$-0.285 - 0.185 + 0.415 = -0.055$	$1/(1 + e^{0.055}) \simeq 0.486$
5	$0.253 + 0.405 \cdot 0.634 - 0.47 \cdot 0.486 = 0.281$	$1/(1 + e^{-0.281}) \simeq 0.569$

c. Valoarea produsă de rețea după ce a fost aplicată o iterație din algoritmul de retro-propagare (0.569) este mai aproape de valoarea-target (0) decât fusese output-ul produs de rețea înainte de aplicarea acestei prime iterații (0.63). Este de așteptat ca rezultatul să se îmbunătățească la execuția următoarelor iterații ale algoritmului.

Regularizarea perceptronilor / rețelelor
pentru prevenirea overfitting-ului,
prin extinderea cu încă un termen
a funcției de cost

Tom Mitchell, *Machine Learning* book, 1997, pr. 4.10

Consider the following error function (which is an alternative to the sum of squares function):

$$E(\bar{w}) = \frac{1}{2} \sum_{d \in D} \sum_{k \in \text{Outputs}} (t_{kd} - o_{kd})^2 + \gamma \sum_{i,j} w_{ji}^2$$

Derive the gradient descent update rule for this definition of E . Show that it can be implemented by multiplying each weight by some constant before performing the standard gradient descent update ($w \leftarrow w + \Delta w$).

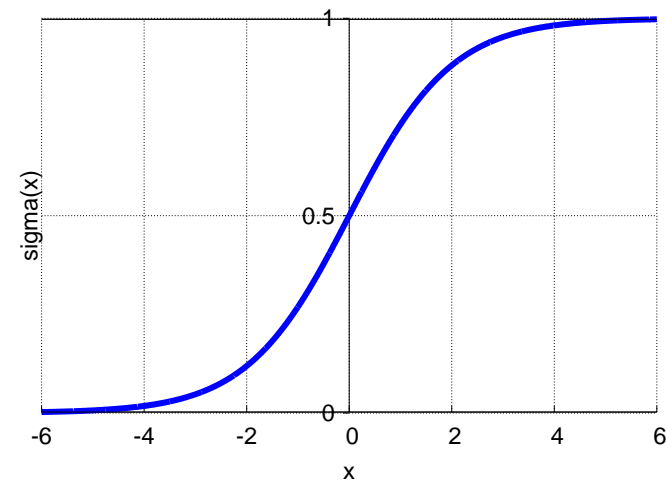
Convenție de notație (in Romanian):

Am folosit mai sus obișnuita [de acum] scriere a variabilelor w , și anume w_{ji} este ponderea care corespunde conexiunii dinspre unitatea i către unitatea j . Notațiile t_{kd} și o_{kd} corespund valorii-target și respectiv output-ului unității liniare k de pe nivelul de ieșire al rețelei, corespunzător input-ului x_d .

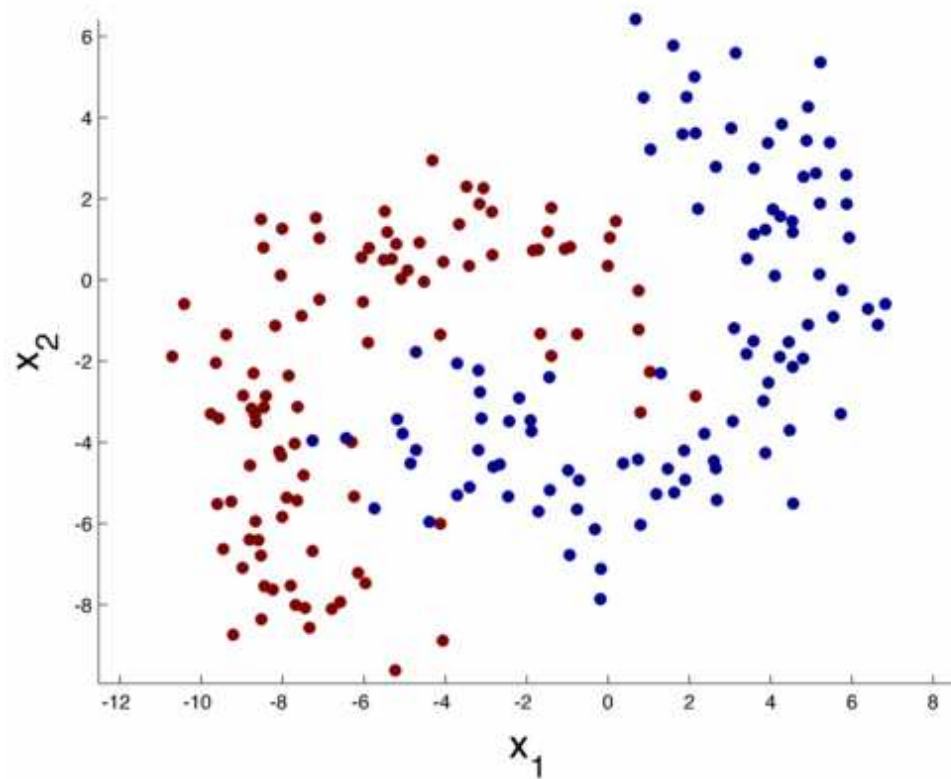
Remark (in Romanian)

Suma $\sum_i w_i^2$ se poate scrie ca $\|w\|^2$, unde $\|w\|$ notează norma vectorului w . Intuitiv, minimizarea funcției obiectiv $E(\bar{w})$ va implica menținerea lui $\|w\|^2$ la o valoare destul de redusă.

Efectul practic, în cazul folosirii de unități sigmoide este următorul: granița (suprafața) de decizie calculată de către rețea pentru ponderi w mici (în modul) este aproape liniară — a se vedea graficul funcției σ în jurul originii —, ceea ce o împiedică să se „muleze” în jurul neregularităților din datele de antrenament.

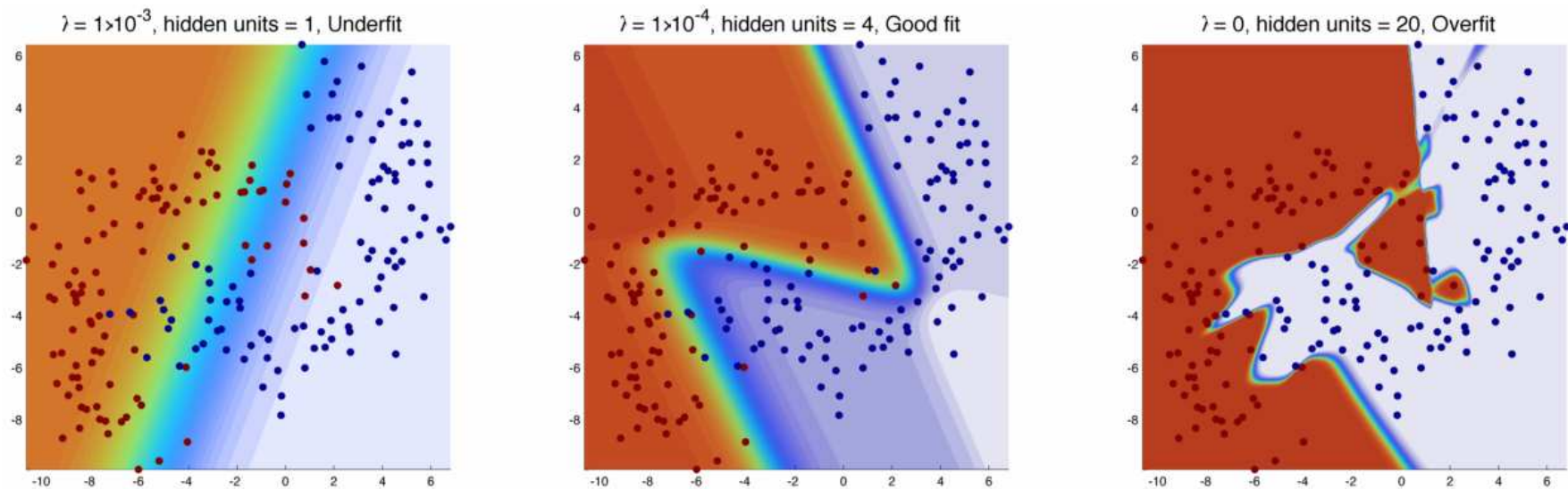


Exemplification: The Banana Dataset



Cf. CMU, 2014 fall, W. Cohen, Z. Bar-Joseph, HW3, pr. 2.

Exemplification (cont'd)



Cf. CMU, 2014 fall, W. Cohen, Z. Bar-Joseph, HW3, pr. 2.

Remark (in Romanian)

Pentru simplitate, vom considera că rețea noastră este de tip feed-forward și că are [doar] două niveluri, iar unitățile sunt liniare.

Din același motiv, vom lucra cu varianta stochastică/incrementală a algoritmului de retro-propagare, adică vom considera o singură instanță de antrenament pe ciclu/„epocă“ de antrenare.

Demonstrația următoare se poate extinde imediat la cazul unităților de diverse tipuri (de exemplu sigmoidal etc.), și la varianta “batch” a algoritmului de retro-propagare.

Solution (in Romanian)

În condițiile stabilite mai înainte, putem scrie:

$$E(\bar{w}) = \frac{1}{2} \sum_{k \in \text{Outputs}} (t_k - o_k)^2 + \gamma \|\bar{w}\|^2,$$

unde t_k și o_k corespund valorii-target și respectiv output-ului unității liniare k de pe nivelul de ieșire al rețelei, pentru instanța de intrare x considerată.

În cele ce urmează, w_{kj} va corespunde unei conexiuni hidden-to-output (unde j indică o unitate de pe nivelul ascuns, iar k o unitate de pe nivelul de ieșire), iar w_{ji} va corespunde unei conexiuni input-to-hidden (i – input, j – hidden).

Prin urmare, pentru ponderile de pe conexiunile hidden-to-output, vom avea:

$$\begin{aligned}
 \frac{\partial E}{\partial w_{kj}} &= \frac{\partial}{\partial w_{kj}} \left[\frac{1}{2}(t_k - o_k)^2 + \gamma w_{kj}^2 \right] + \sum_{k' \neq k} \frac{\partial}{\partial w_{kj}} \left[\frac{1}{2}(t_{k'} - o_{k'})^2 + \gamma w_{k'j}^2 \right] \\
 &= -\frac{1}{2}2(t_k - o_k) \frac{\partial}{\partial w_{kj}} o_k + 2\gamma w_{kj} \\
 &= -(t_k - o_k) \frac{\partial}{\partial w_{kj}} \sum_{j'} w_{kj'} y_{j'} + 2\gamma w_{kj} = -(t_k - o_k) y_j + 2\gamma w_{kj}
 \end{aligned}$$

Așadar,

$$w_{kj} \leftarrow w_{kj} - \eta \frac{\partial E}{\partial w_{kj}} = w_{kj} + \eta [(t_k - o_k) y_j - 2\gamma w_{kj}] = (1 - 2\eta\gamma) w_{kj} + \eta(t_k - o_k) y_j$$

Tratăm acum cazul ponderilor de pe conexiunile input-to-hidden.

Notăm $Downstream(j)$ mulțimea tuturor indicilor k cu proprietatea că există conexiune de la unitatea j către unitatea k .

$$\begin{aligned}
\frac{\partial E}{\partial w_{ji}} &= \frac{\partial}{\partial w_{ji}} \left\{ \left[\sum_{k \in Downstream(j)} \frac{1}{2} (t_k - o_k)^2 \right] + \gamma w_{ji}^2 \right\} = 2\gamma w_{ji} - \sum_{k \in Downstream(j)} (t_k - o_k) \frac{\partial}{\partial w_{ji}} o_k \\
&= 2\gamma w_{ji} - \sum_{k \in Downstream(j)} \left[(t_k - o_k) \frac{\partial}{\partial w_{ji}} \left(\sum_{j'} w_{kj'} y_{j'} \right) \right] \\
&= 2\gamma w_{ji} - \sum_{k \in Downstream(j)} \left\{ (t_k - o_k) \frac{\partial}{\partial w_{ji}} \left[\sum_{j'} w_{kj'} \left(\sum_{i'} w_{j'i'} x_{i'} \right) \right] \right\} \\
&= 2\gamma w_{ji} - \sum_{k \in Downstream(j)} \left[(t_k - o_k) \frac{\partial}{\partial w_{ji}} \left(\sum_{j'} \sum_{i'} w_{kj'} w_{j'i'} x_{i'} \right) \right] \\
&= 2\gamma w_{ji} - x_i \left(\sum_{k \in Downstream(j)} (t_k - o_k) w_{kj} \right)
\end{aligned}$$

Așadar,

$$\begin{aligned}
 w_{ji} \leftarrow w_{ji} - \eta \frac{\partial E}{\partial w_{ji}} &= w_{ji} + \eta \left[x_i \left(\sum_{k \in \text{Downstream}(j)} w_{kj} (t_k - o_k) \right) - 2\gamma w_{ji} \right] \\
 &= (1 - 2\eta\gamma)w_{ji} + \eta x_i \left(\sum_{k \in \text{Downstream}(j)} w_{kj} (t_k - o_k) \right)
 \end{aligned}$$

Concluzie

Analizând regulile de actualizare a ponderilor deduse mai sus, se observă că în ambele cazuri (input-to-hidden și hidden-to-output) ele sunt de forma $w \leftarrow (1 - 2\eta\gamma)w + \Delta w$.

Deosebirea față de cazul clasic ($w \leftarrow w + \Delta w$), este evidentă. Întrucât algoritmul de retro-propagare inițializează ponderile w cu valori apropiate de 0 — iar în practică se folosesc valori mici pentru constantele η și γ —, lucrând cu varianta funcției de eroare propusă în această problemă rezultă valori ale ponderilor w mai aproape de 0 (decât în varianta clasică).

**[Networks of] sigmoidal units:
The cross-entropy error function**

CMU, 2011 fall, Eric Xing, HW1, pr. 3.3

Fie datele de antrenament $D = (X, t) = \{(x_i, t_i)\}, i = 1, 2, \dots, N$. De exemplu, x_i (element dintr-un spațiu \mathbb{R}^d) poate fi imaginea unei fețe, în vreme ce t_i este o etichetă binară care are valoarea 0 dacă fața respectivă este a unui bărbat, respectiv 1 în cazul unei femei.

Vom considera o unitate neuronală de tip sigmoidal.^a Notând cu y valoarea reală produsă de această unitate, rezultă că $y \in (0, 1)$. Este natural să interpretăm acest output ca fiind probabilitatea ca eticheta booleană t să fie 1. Formal, putem scrie $y \stackrel{def.}{=} P(t = 1 \mid x; w)$.

În consecință, este natural să căutăm valori convenabile pentru ponderile w folosind principiul verosimilității maxime (MLE),^b sub forma următoare:

$$w_{MLE} = \operatorname{argmax}_w \prod_{i=1}^N P(t_i \mid x_i; w).$$

^aMai general, se poate considera o rețea formată din unități sigmoide. (De fapt, așa a fost formulată problema originală de la CMU.) Într-un astfel de caz este însă suficient să presupunem că tipul unității / unităților de pe nivelul de ieșire din rețea este cel sigmoidal.

^bDe fapt, în cazul de față vom lucra cu verosimilitate condițională.

a. Arătați că a maximiza expresia $\prod_{i=1}^N P(t_i \mid x_i; w)$ revine la a minimiza expresia următoare, despre care putem spune, după o observare atentă, că este o cross-entropie:

$$E(w) = - \sum_{i=1}^N [t_i \ln y_i + (1 - t_i) \ln(1 - y_i)],$$

unde y_i este output-ul produs de rețeaua neuronală pentru exemplul x_i .

b. Să presupunem că este posibil ca etichetele datelor de antrenament să fi fost puse eronat, și anume cu probabilitatea ε . Considerând că datele de antrenament sunt distribuite în mod identic și independent unele de altele, scrieți expresia funcției de eroare / cost $E^*(w, \varepsilon)$ care corespunde log-verosimilității cu semn schimbat (engl., the negative log likelihood).

c. Verificați că funcția de eroare $E(w)$ se obține din $E^*(w, \varepsilon)$ pentru cazul particular $\varepsilon = 0$.

d. Explicați de ce anume funcția de eroare $E^*(w, \varepsilon)$ va face ca modelul obținut să fie [mai] robust în raport cu date incorect etichetate, spre deosebire de funcția uzuală $E(w)$.

a. Conform enunțului, $y_i \stackrel{def.}{=} P(t_i = 1|x_i; w)$. Așadar, putem scrie

$$P(t_i|x_i; w) = \begin{cases} y_i & \text{dacă } t_i = 1, \\ 1 - y_i & \text{dacă } t_i = 0. \end{cases}$$

Se verifică imediat că, din punct de vedere matematic, relația de mai sus se poate exprima în mod echivalent sub forma

$$P(t_i|x_i; w) = (y_i)^{t_i} (1 - y_i)^{1-t_i}.$$

Această expresie, în ciuda faptului că este mai puțin intuitivă, este mult mai convenabilă pentru calculele pe care trebuie să le facem în vederea aplicării principiului verosimilității maxime:

$$\begin{aligned} w_{MLE} &\stackrel{def.}{=} \operatorname{argmax}_w P(t_1, \dots, t_n | x_1, \dots, x_n; w) \stackrel{i.i.d.}{=} \operatorname{argmax}_w \prod_{i=1}^N P(t_i | x_i; w) \\ &= \operatorname{argmax}_w \sum_{i=1}^N \ln P(t_i | x_i; w) = \operatorname{argmax}_w \sum_{i=1}^N \ln (y_i)^{t_i} (1 - y_i)^{1-t_i} \\ &= \operatorname{argmax}_w \sum_{i=1}^N (t_i \ln y_i + (1 - t_i) \ln(1 - y_i)) = \operatorname{argmax}_w (-E(w)) = \operatorname{argmin}_w E(w). \end{aligned}$$

b. În cazul în care etichetarea datelor de antrenament s-a făcut în mod eronat, cu probabilitatea ε , rezultă imediat că

$$x_i \text{ a fost etichetat cu } \begin{cases} 1 & \text{cu probabilitatea } 1 - \varepsilon, \text{ dacă } t_i = 1, \\ 0 & \text{cu probabilitatea } \varepsilon, \text{ dacă } t_i = 1, \\ 0 & \text{cu probabilitatea } 1 - \varepsilon, \text{ dacă } t_i = 0, \\ 1 & \text{cu probabilitatea } \varepsilon, \text{ dacă } t_i = 0. \end{cases}$$

Prin urmare, putem scrie în manieră condensată

$$P(t_i|x_i; w) = \begin{cases} y_i(1 - \varepsilon) + (1 - y_i)\varepsilon & \text{dacă } t_i = 1, \\ (1 - y_i)(1 - \varepsilon) + y_i\varepsilon & \text{dacă } t_i = 0. \end{cases}$$

Ba chiar și mai mult, această expresie este echivalentă cu următoarea:

$$P(t_i|x_i; w) = [y_i(1 - \varepsilon) + (1 - y_i)\varepsilon]^{t_i} [(1 - y_i)(1 - \varepsilon) + y_i\varepsilon]^{1-t_i}.$$

Făcând un calcul similar cu cel de la punctul a, va rezulta:^a

$$E^*(w, \varepsilon) = - \sum_{i=1}^N \{t_i \ln[y_i(1 - \varepsilon) + (1 - y_i)\varepsilon] + (1 - t_i) \ln[(1 - y_i)(1 - \varepsilon) + y_i\varepsilon]\}.$$

^aSe verifică imediat că $[y_i(1 - \varepsilon) + (1 - y_i)\varepsilon] + [(1 - y_i)(1 - \varepsilon) + y_i\varepsilon] = 1$.

c. Înlocuind parametrul ε cu valoarea 0 în expresia pe care tocmai am obținut-o mai sus, vom avea:

$$\begin{aligned} E^*(w, 0) &= - \sum_{i=1}^N \{ [t_i \ln y_i + (1 - y_i) \times 0] + (1 - t_i) \ln[(1 - y_i) + y_i \times 0] \} \\ &= - \sum_{i=1}^N [t_i \ln y_i + (1 - t_i) \ln(1 - y_i)] = E(w) \end{aligned}$$

d. Vom considera ca *exemplu* cazul extrem (însă instructiv) în care toate etichetele au fost puse incorect. Dacă am folosi funcția de eroare $E(w)$, atunci modelul produs de algoritmul de retro-propagare ar fi complet eronat. Dacă în schimb vom folosi funcția de eroare $E^*(w, \varepsilon)$ cu $\varepsilon = 1$, atunci algoritmul de retro-propagare va încerca să învețe un model care clasifică datele exact invers față de etichetele originale, ceea ce corespunde obiectivului nostru.

Important Remarks (in Romanian)

1.

Analizând cu atenție demonstrația pentru deducerea regulilor de actualizarea ponderilor w pentru algoritmul de retro-propagare atunci când funcția obiectiv este prezenta cross-entropie (vezi CMU, 2008 fall, E. Xing, HW2, pr. 2.2), și făcând comparația cu demonstrația din cartea lui Tom Mitchell (alternativ, vezi CMU, 2008(?) spring, HW2, pr. 2.2-4) — cazul utilizării semi-sumei pătratelor erorilor ca funcție obiectiv —, se constată că singura diferență [între cele două seturi de reguli de actualizare] apare în expresia lui δ_j pentru j indice de unitate de pe nivelul ascuns: în acel caz, în componenta lui δ_j este inclus factorul $f'(net_j)$, însă aici (folosind σ , funcția sigmoidală) acest factor lipsește. În rest, totul este la fel!

2.

Ar fi interesant de văzut care este forma regulilor de actualizare atunci când în locul prezentului criteriu de optimizat se folosește $E(w, \varepsilon)$ definit [în acest set de slide-uri] la punctele b și c .

**Deep neural networks:
gradient vanishing / explosion**

CMU, 2015 fall, Eric Xing, Zib Bar-Joseph, HW3, pr. 1.3

In this problem we will study the difficulty of back-propagation in training deep neural networks. For simplicity, we consider the simplest deep neural network: one with just a single neuron in each layer, where the output of the neuron in the j th layer is $z_j \stackrel{not.}{=} f(net_j) = f(w_j z_{j-1} + b_j)$. Here f is some activation function whose derivative on x is $f'(x)$. Let m be the number of layers in the neural network, L the training loss function.

a. Derive the derivative of L w.r.t. b_1 (the bias of the neuron in the first layer).

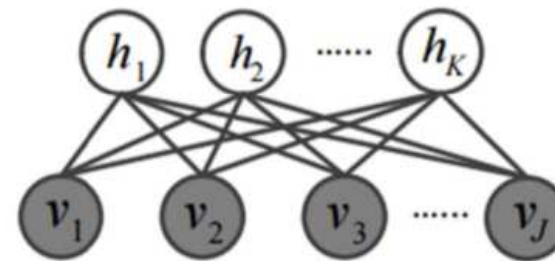
b. Assume the activation function is the usual sigmoid function $\sigma(x) = 1/(1 + \exp(-x))$. The weights \bar{w} are initialized to be $|w_j| < 1$ ($j = 1, \dots, m$).

Explain why the above gradient ($\partial L / \partial b_1$) tends to vanish ($\rightarrow 0$) when m is large.

Even if $|w|$ is large, the above gradient would also tend to vanish, rather than explode ($\rightarrow \infty$). Explain why.

c. One of the approaches to (partially) address the gradient vanishing/explosion problem is to use the *rectified linear* (ReLU) activation function instead of the sigmoid. The ReLU activation function is $\sigma(x) = \max\{0, x\}$. Explain why ReLU can alleviate the gradient vanishing problem as faced by sigmoid.

d. A second approach to (partially) address the gradient vanishing/explosion problem is *layer-wise pre-training*. *Restricted Boltzmann machine* (RBM) is one of the widely-used models for layer-wise pre-training. The nearby figure shows an example of RBM which includes K hidden units \bar{h} , and J input units \bar{v} .



Let us define the *joint distribution* over \bar{v} and \bar{h} as having the following general form:

$$P(\bar{v}, \bar{h}) = \frac{1}{Z} \exp \left(\sum_i \theta_i \phi_i(\bar{v}, \bar{h}) \right),$$

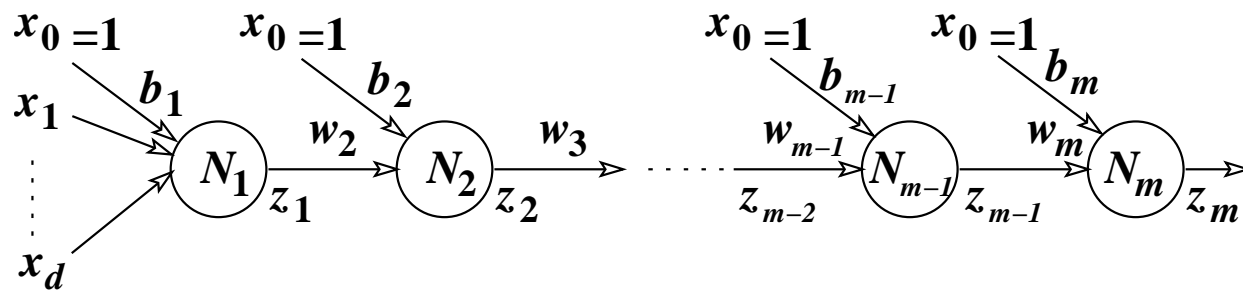
where $Z = \sum_{\bar{v}', \bar{h}'} \exp \left(\sum_i \theta_i \phi_i(\bar{v}', \bar{h}') \right)$ is the normalization term; $\phi_i(\bar{v}, \bar{h})$ are some features; θ_i are the parameters corresponding to the weights in the RBM. Consider the simplest learning algorithm, gradient descent.

Show that

$$\frac{\partial \log P(\bar{v})}{\partial \theta_i} = \sum_{\bar{h}} \phi_i(\bar{v}, \bar{h}) P(\bar{h}|\bar{v}) - \sum_{\bar{v}', \bar{h}'} \phi_i(\bar{v}', \bar{h}') P(\bar{v}', \bar{h}').$$

Răspuns

a. Ilustrăm rețeaua neuronală profundă definită în enunț:



$$\Rightarrow z_1 = f(\text{net}_1) = f(b_1 + w_{11}x_1 + \dots + w_{1d}x_d)$$

$$z_2 = f(\text{net}_2) = f(b_2 + w_2z_1)$$

...

$$o = z_m = f(\text{net}_m) = f(b_m + w_mz_{m-1})$$

Funcția de pierdere (engl., loss function) L se va exprima astfel:

$$\begin{aligned}
 & L(w_{11}, \dots, w_{1d}, w_2, \dots, w_m, b_1, \dots, b_m) \\
 & \stackrel{not.}{=} L(f(\underbrace{\quad}_{b_m + w_m z_{m-1}}^{net_m})) \\
 & = L(f(b_m + w_m f(\underbrace{\quad}_{b_{m-1} + w_{m-1} z_{m-2}}^{net_{m-1}}))) \\
 & \quad \dots \\
 & = L(f(b_m + w_m f(b_{m-1} + w_{m-1} f(b_{m-2} + \dots + w_3 f(\underbrace{\quad}_{b_2 + w_2 z_1}}^{net_2}) \dots)))) \\
 & = L(f(b_m + w_m f(b_{m-1} + w_{m-1} f(b_{m-2} + \dots + w_3 f(b_2 + w_2 z_1) \dots)))) \\
 & = L(f(b_m + w_m f(b_{m-1} + w_{m-1} f(b_{m-2} + \dots + w_3 f(b_2 + w_2 f(\underbrace{\quad}_{b_1 + w_{11} x_1 + \dots + w_{1d} x_d}}^{net_1}) \dots))))
 \end{aligned}$$

Pentru a deriva funcția L în raport cu b_1 , ne vom aminti mai întâi regula de derivare a unei compuneri multiple de funcții $f_1(f_2(\dots f_n(x)\dots))$. Pe lângă forma clasică pe care o știm din liceu,

$$f'_1(f_2(\dots f_n(x)\dots)) \cdot f'_2(\dots f_n(x)\dots) \cdot f'_n(x),$$

ea se poate scrie sub forma așa-numitei *reguli de înlănțuire*:

$$\frac{\partial f_1}{\partial f_2} \cdot \frac{\partial f_2}{\partial f_3} \cdot \dots \cdot \frac{\partial f_2}{\partial f_1} \cdot \frac{\partial f_1}{\partial x}.$$

Adoptând / aplicând această regulă pentru a calcula derivata parțială cerută în enunț, vom avea:

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= \frac{\partial L}{\partial net_m} \cdot \frac{\partial net_m}{\partial net_{m-1}} \cdot \dots \cdot \frac{\partial net_2}{\partial net_1} \cdot \frac{\partial net_1}{\partial b_1} \\ &= L'(f(net_m)) \cdot f'(net_m) \cdot \\ &\quad w_m \cdot f'(net_{m-1}) \cdot \\ &\quad w_{m-1} \cdot f'(net_{m-2}) \cdot \\ &\quad \dots \\ &\quad w_2 \cdot f'(net_1) \cdot \\ &\quad 1 \\ &= L'(f(net_m)) \cdot f'(net_1) \cdot \prod_{k=2}^m (f'(net_k) \cdot w_k) \end{aligned}$$

b. Știm că $\sigma(x) \in (0, 1)$ și $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ pentru orice $x \in \mathbb{R}$, iar maximul funcției $z(1 - z)$ este $1/4$ și se atinge în punctul $z = 1/2$. Prin urmare,

$$|\prod_{k=2}^m \sigma'(net_k)| \leq \left(\frac{1}{4}\right)^{m-1}.$$

Întrucât $|w_k| < 1$ pentru orice k (conform ipotezei din enunț), rezultă imediat că $\left|\frac{\partial L}{\partial b_1}\right| = |L'(\sigma(net_m)) \cdot \sigma'(net_1)| \cdot \prod_{k=2}^m |\sigma'(net_k) \cdot w_k|$ tinde la 0 pentru $k \rightarrow +\infty$ și pentru orice input \bar{x} fixat.

Pentru a contracara acest fenomen de „dispariție” a gradientului, ar trebui să impunem condiții de forma $|w_k \sigma'(net_k)| > 1$. Trebuie însă să remarcăm faptul că $\sigma(net_k)$ depinde de w_k : $\sigma(net_k) = \sigma(b_k + w_k z_{k-1})$.

În consecință, dacă îi vom da lui w_k posibilitatea să ia valori mari în modul, atunci $z_k \stackrel{not.}{=} b_k + w_k z_{k-1}$ va lua de asemenea valori mari în modul, iar $\sigma(z_k)$ va tinde fie la 0 fie la 1, deci $\sigma'(z_k) = \sigma(z_k)(1 - \sigma(z_k))$ va tinde la 0.

În sfârșit, se poate verifica în mod riguros că $\lim_{z \rightarrow \pm\infty} z\sigma'(z) = 0$ (lăsăm această demonstrație ca exercițiu).

c. Dacă f este funcția de activare ReL, atunci $f'(x) = 0$ pentru $x < 0$ și $f'(x) = 1$ pentru $x > 0$. Prin urmare, folosind această funcție putem împiedica „dispariția” gradientului.

d. Pornind de la expresia distribuției probabiliste $P(\bar{v}, \bar{h})$ care a fost dată în enunț, explicitând mai întâi constanta de normalizare Z , vom putea scrie apoi expresia distribuției marginale $P(\bar{v})$:

$$P(\bar{v}, \bar{h}) = \frac{1}{\sum_{\bar{v}', \bar{h}'} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}'))} \exp\left(\sum_k \theta_k \phi_k(\bar{v}, \bar{h})\right)$$

$$\Rightarrow P(\bar{v}) = \frac{1}{\sum_{\bar{v}', \bar{h}'} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}'))} \sum_{\bar{h}} \exp\left(\sum_k \theta_k \phi_k(\bar{v}, \bar{h})\right)$$

Prin urmare,

$$\ln P(\bar{v}) = \ln \sum_{\bar{h}} \exp\left(\sum_k \theta_k \phi_k(\bar{v}, \bar{h})\right) - \ln \sum_{\bar{v}', \bar{h}'} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}')).$$

În consecință, derivata parțială a funcției $\ln P(\bar{v})$ în raport cu parametrul θ_i se poate calcula astfel:

$$\begin{aligned}
\frac{\partial \ln P(\bar{v})}{\partial \theta_i} &= \\
&= \frac{\sum_{\bar{h}} \exp(\sum_k \theta_k \phi_k(\bar{v}, \bar{h})) \cdot \phi_i(\bar{v}, \bar{h})}{\sum_{\bar{h}} \exp(\sum_k \theta_k \phi_k(\bar{v}, \bar{h}))} - \frac{\sum_{\bar{v}', \bar{h}'} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}')) \cdot \phi_i(\bar{v}, \bar{h})}{\underbrace{\sum_{\bar{v}', \bar{h}'} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}'))}_Z} \\
&= \frac{\sum_{\bar{h}} \frac{1}{Z} \exp(\sum_k \theta_k \phi_k(\bar{v}, \bar{h})) \cdot \phi_i(\bar{v}, \bar{h})}{\sum_{\bar{h}} \frac{1}{Z} \exp(\sum_k \theta_k \phi_k(\bar{v}, \bar{h}))} - \sum_{\bar{v}', \bar{h}'} \frac{1}{Z} \exp(\sum_k \theta_k \phi_k(\bar{v}', \bar{h}')) \cdot \phi_i(\bar{v}, \bar{h}) \\
&= \frac{\sum_{\bar{h}} P(\bar{v}, \bar{h}) \cdot \phi_i(\bar{v}, \bar{h})}{P(\bar{v})} - \sum_{\bar{v}', \bar{h}'} P(\bar{v}', \bar{h}') \phi_i(\bar{v}', \bar{h}') \\
&= \sum_{\bar{h}} \frac{P(\bar{v}, \bar{h})}{P(\bar{v})} \cdot \phi_i(\bar{v}, \bar{h}) - \sum_{\bar{v}', \bar{h}'} P(\bar{v}', \bar{h}') \phi_i(\bar{v}', \bar{h}') \\
&= \sum_{\bar{h}} P(\bar{h}|\bar{v}) \cdot \phi_i(\bar{v}, \bar{h}) - \sum_{\bar{v}', \bar{h}'} P(\bar{v}', \bar{h}') \phi_i(\bar{v}', \bar{h}')
\end{aligned}$$

The kernelized *Perceptron*

Liviu Ciortuz, following

CMU, 2013 spring, A. Smola, B. Póczos, HW2, pr. 2.d

Stanford, 2016 fall, A. Ng, J. Duchi, HW2, pr. 2

MIT, 2006 fall, Tommy Jaakkola, HW2, pr. 3

Majoritatea clasificatorilor liniari pot fi kernel-izați, în vederea clasificării de date care sunt neseperabile liniar. Aici vom elabora varianta kernel-izată [duală] a algoritmului *Perceptron*.

Pseudo-codul variantei simple, ne-kernel-izate a algoritmului *Perceptron* este prezentat alăturat. (Operatorul \cdot reprezintă produsul scalar.)

```

initialize  $\bar{w} \leftarrow \bar{0}$ 
for  $i = 1, \dots, n$  do
  if  $y_i (\bar{w} \cdot \bar{x}_i) \leq 0$  then
     $\bar{w} \leftarrow \bar{w} + y_i \bar{x}_i$ 
  end if
end for

```

Ca input, *Perceptron*-ul kernel-izat va lua secvența de instanțe etichetate $\{x_i, y_i\}_{i=1}^n$ cu $x_i \in \mathbb{R}^d$ și $y_i \in \{-1, 1\}$, precum și funcția-nucleu $K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^m$, cu proprietatea că există $m > d$ și o funcție („mapare“) $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^m$ astfel încât $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)$ pentru orice x_i, x_j . (Din punct de vedere practic, este necesar ca funcția K să poată fi calculată în mod eficient.)

Perceptron-ul kernel-izat va lucra nu cu instanțele x_i , ci cu imaginile lor, $\phi(x_i)$, și va încerca să găsească pentru acestea un separator liniar în spațiul \mathbb{R}^m .

Arătați cum anume se scrie regula de actualizare a coeficienților $\alpha_1, \dots, \alpha_n$ la procesarea unui nou exemplu (x_i, y_i) de către algoritmul *Perceptron* kernelizat, și cum se face predicția ($y = 1$ sau -1) pentru o instanță nouă x .

Răspuns

În spațiul \mathbb{R}^m (care se mai numește, în contextul kernel-izării, *spațiul de trăsături*), regula de actualizare a ponderilor w se scrie astfel:

$$w^{(i)} \leftarrow w^{(i-1)} + y_i w^{(i-1)} \cdot \phi(x_i) \text{ dacă } y_i w^{(i-1)} \cdot \phi(x_i) \leq 0,$$

unde operatorul \cdot desemnează produsul scalar.

În consecință, întrucât facem inițializarea $w^{(0)} = \bar{0}$, vectorul $w \in \mathbb{R}^m$ va fi întotdeauna o combinație liniară de vectori de trăsături, $\phi(x_i)$. Aceasta înseamnă că există coeficienții $\alpha_i \in \mathbb{R}$ astfel încât $w^{(i)} = \sum_{l=1}^i \alpha_l \phi(x_l)$ după procesarea primelor i instanțe de antrenament. Așadar, vectorul $w^{(i)}$ poate fi reprezentat în mod compact prin coeficienții α_l (cu $l = 1, \dots, i$) din această combinație liniară.

În particular, valoarea inițială $w^{(0)}$ corespunde cazului când suma nu conține niciun termen (adică avem o listă vidă de coeficienți α_l).

Arătăm acum că putem actualiza în mod eficient coeficienții α_i .

La iterația i trebuie să calculăm produsul scalar $w^{(i-1)} \cdot \phi(x_i)$. Întrucât produsul scalar în spațiul de trăsături \mathbb{R}^m este o operație costisitoare atunci când m este mare, vom ține cont că

$$w^{(i-1)} \cdot x_i = \left(\sum_{l=1}^{i-1} \alpha_l \phi(x_l) \right) \cdot \phi(x_i) = \sum_{l=1}^{i-1} \alpha_l (\phi(x_l) \cdot \phi(x_i)) = \sum_{l=1}^{i-1} \alpha_l K(x_l, x_i),$$

ceea ce înseamnă că acești coeficienți se pot calcula într-adevăr în mod eficient.

În mod similar, se poate face în mod eficient *predicția* clasei/etichetei pentru o instanță nouă (de test) $x \in \mathbb{R}^d$:

$$w^{(i)} \cdot \phi(x) = \sum_{l=1}^i \beta_l \phi(x_l) \cdot \phi(x) = \sum_{l=1}^i \beta_l K(x_l, x).$$

Sumarizând, algoritmul pentru antrenarea *Perceptron*-ului kernel-izat se poate scrie în pseudo-cod astfel:

```

initialize  $\alpha_i = 0$  for  $i = 1, \dots, n$ ;
for  $i = 1, \dots, n$  do
  if  $y_i \sum_{l=1}^{i-1} \alpha_l K(x_l, x_i) \leq 0$  then
     $\alpha_i \leftarrow y_i$ 
  end if
end for

```

Observație importantă: Se poate arăta relativ ușor că raționamentele (și rezultatele) din acest exercițiu se pot extinde și la cazul perceptronului care folosește funcție de activare de tip prag (în particular, funcția *sign*), rată de învățare oarecare $\eta > 0$ — dar face inițializarea vectorului de ponderi $w \in \mathbb{R}^d$ cu $\bar{0}$, ceea ce este echivalent cu $\alpha_i = 0$, pentru $i = 1, \dots, n$ — și care eventual parcurge de mai multe ori setul de date de antrenament. În acest caz, regula de actualizare a perceptronului kernel-izat este de forma

$$\alpha_i \leftarrow \alpha_i + \eta(y_i - \text{sign}(w^{(i-1)} \cdot \phi(x_i))) \text{ dacă } y_i \text{sign}(w^{(i-1)} \cdot \phi(x_i)) \leq 0$$

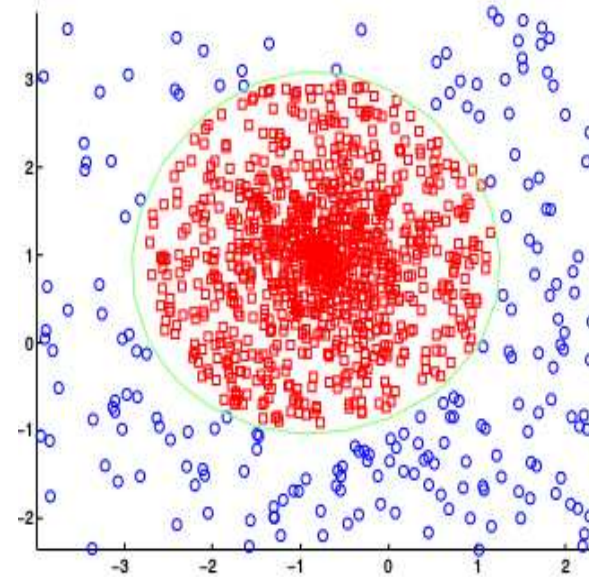
unde, exact ca și mai sus, $w^{(i-1)} \cdot \phi(x_i) = \left(\sum_{l=1}^{i-1} \alpha_l \phi(x_l) \right) \cdot \phi(x_i) = \sum_{l=1}^i \alpha_l K(x_l, x_i)$.

```
kernel = 'exp(-transpose(x_i - x_j)(x_i - x_j)/(2s^2))'; s = 3;
function  $\alpha$ =TRAIN_KERNEL_PERCEPTRON(X, y, kernel)
```

```
    n, d = size(X);
    K = [ ];
    for i = 1 : n
        x_i = X(i, :)' ;
        for j = 1 : n
            x_j = X(j, :)' ;
            K(i, j) = eval(kernel);
        end
    end
```

```
     $\alpha$  = zeros(n, 1);
    mistakes = 1;
    while mistakes > 0
        mistakes = 0;
        for i = 1 : n
            if  $\alpha' K(:, i) y(i) \leq 0$ 
                 $\alpha(i) = \alpha(i) + y(i)$ ;
                mistakes = mistakes + 1;
            end
        end
    end
```

Source: MIT, 2006 fall, Tommy Jaakkola, HW2,
pr. 3



```
function f=DISCRIMINANT_FUNCTION( $\alpha$ , X, kernel, X_test)
```

```
    n, d = size(X);
    K = [ ];
    for i = 1 : n
        x_i = X(i, :)' ;
        x_j = X_test;
        K(i) = eval(kernel);
    end
    f =  $\alpha' K$ ;
```