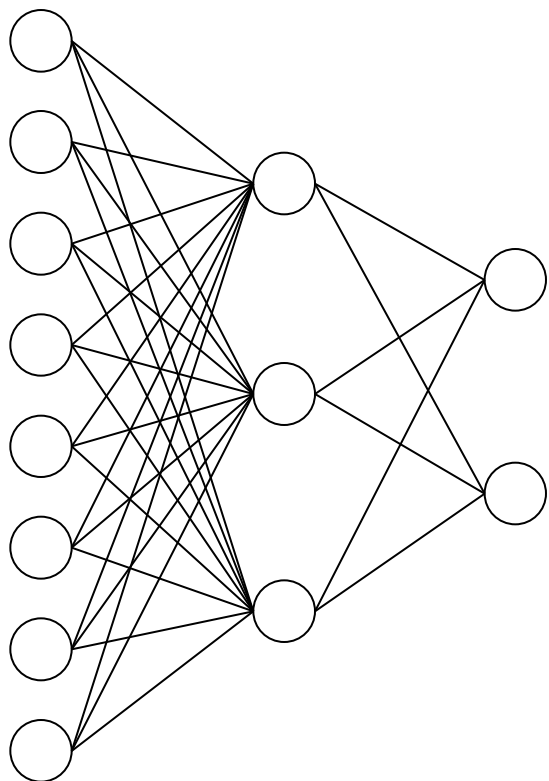


Osnove neuronskih mreža

(Fundamentals of Neural Networks)

Umjetna neuronska mreža (*artificial neural network*) je programski ili sklopovski sustav građen od jednakih, jednostavnih i međusobno povezanih dijelova (čvorova) koji može obavljati raznovrsne i složene zadaće ovisno o parametrima dijelova i veza među njima. Dijelovi (čvorovi) neuronske mreže nazivaju se neuroni (*neurons*).



Građa umjetnih neuronskih mreža je slojevita i nalikuje na (točnije, oponaša) biološka živčana tkiva pa im otuda i naziv.

U tehničkim disciplinama, gdje se “umjetnost” mreže podrazumijeva, najčešće se govori samo *neuronska mreža*.

Ako signali (podatci) u mreži putuju samo u jednom smjeru, od ulaza prema izlazu, radi se o *unaprijednoj* (*feedforward*) mreži. Postoje i složenije strukture u kojima signali putuju u oba smjera i u kojima ima i povratnih veza, ali to su specijalizirane mreže (automatika, obrada signala) kojima se u okviru ovog predmeta nećemo baviti ... 😊

Neuronske mreže mogu se primijeniti u raznorodnim problemima, na primjer:

1. aproksimiranje funkcija: one koje je teško izraziti formulama (to je zapravo preslikavanje domene u kodomenu; *mapping*); regresijska analiza
2. razvrstavanje (*classification*) podataka i prepoznavanje uzoraka, odnosno neobičnosti (izuzetaka)
3. obrada signala i podataka, npr. komplicirano filtriranje ili upravljanje procesima (automatika)

Razlikujemo tri osnovna sloja u mreži: ulazni, skriveni i izlazni sloj, pri čemu skriveni sloj može imati proizvoljno podslojeva (razina).

Broj neurona u pojedinom sloju ovisi o zadaći mreže:

- ulazni sloj ima onoliko neurona koliko je ulaznih varijabli (= dimenzija domene)
- izlazni sloj ima onoliko neurona koliko je željenih izlaznih rezultata (= dimenzija kodomene)
 - na primjer, ako samo treba ustanoviti je li na ulazu neki uzorak ili nije, dovoljan je jedan neuron na izlazu (rezultat DA-NE). Ako, pak, treba filtrirati signal, izlaz treba imati onoliko neurona koliko izlazni signal treba imati točaka.
- ukupni broj neurona u skrivenom sloju mreže s jednim izlazom uglavnom je unutar 2...30 % broja ulaznih varijabli, pri čemu razine u skrivenom sloju redovito imaju sve manje i manje neurona što je razina bliža izlazu

Manji potrební broj neurona u skrivenom sloju intuitivno se može razumjeti ako se neuronska mreža shvati posve apstraktno kao struktura koja dodjeljuje vrijednosti svakoj točki u N -dimenzionalnom prostoru (*hiperprostor*) ulaznih varijabli (pojedininim područjima (djelićima) tog prostora).

Preslikavanje bi se, teorijski, moglo definirati pohranom željenih izlaznih vrijednosti za sve točke tog prostora (npr. u nekakvu *lookup* tablicu), ali čak i kad bismo znali vrijednosti u svim tim točkama, a redovito ih ne znamo, previše ih je za pohranu.

- Na primjer, za samo 10 ulaznih varijabli koje mogu poprimiti samo 10 različitih vrijednosti prostor već ima 10^{10} točaka!

Neuronska mreža “simulira” tabličnu definiciju preslikavanja i izravno izračunava izlaz za svaku moguću kombinaciju vrijednosti ulaznih parametara.

Budući da za “omeđivanje” jednog djelića N -dim. prostora trebamo $2N$ točaka, identifikacija jednog djelića zahtijeva $2N$ prilagodljivih parametara. Ako je ta identifikacija zadaća neuronske mreže, onda je $2N$ ujedno i najmanji potrební broj njezinih prilagodljivih parametara. Ako je, nadalje, mreža troslojna, potpuno povezana i ima N ulaznih varijabli, slijedi da skriveni sloj, načelno, može imati samo dva neurona.

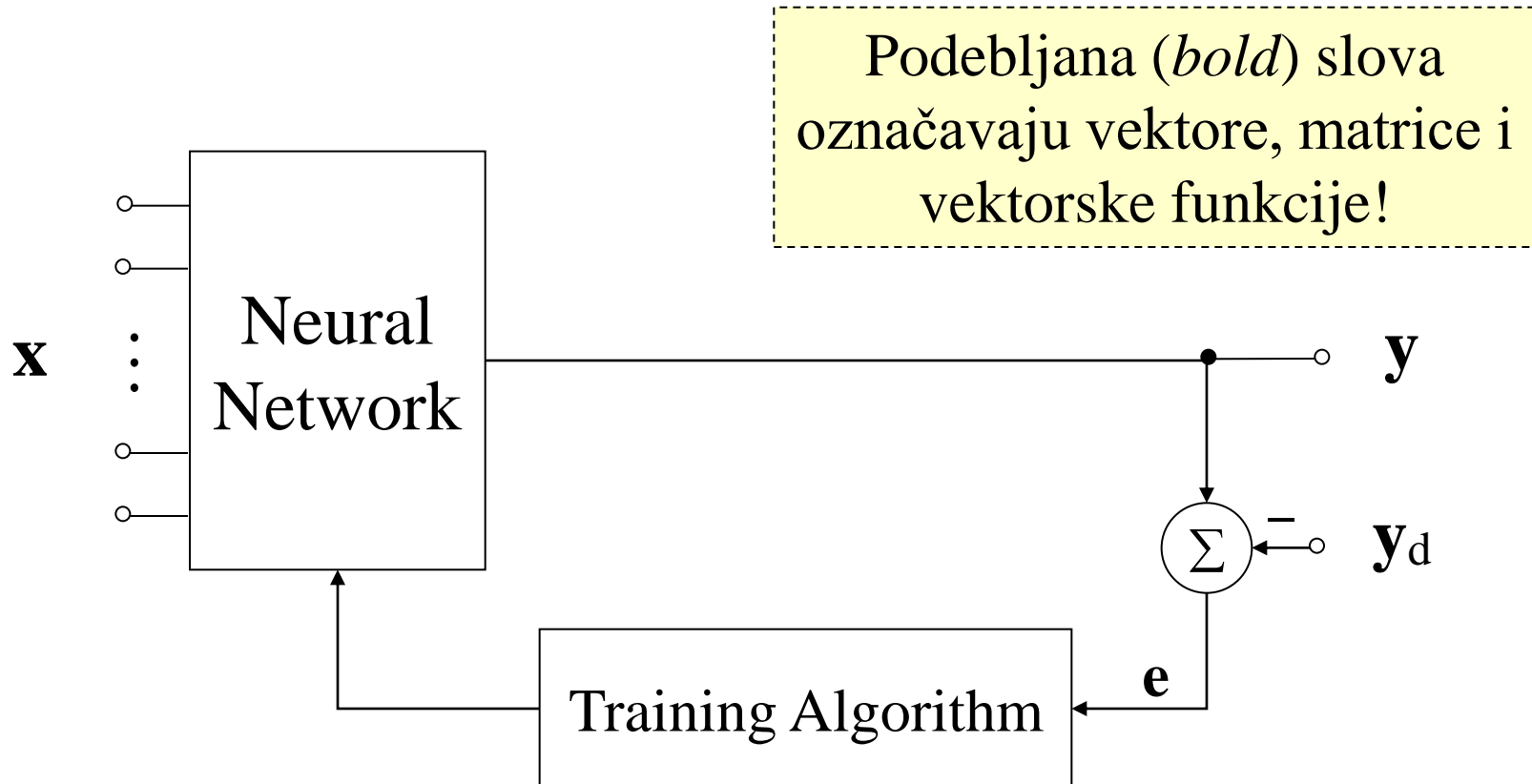
Identifikaciji doprinose i kombinacije neurona u skrivenom sloju te njihove veze s drugim razinama skrivenog sloja i izlaznim slojem pa ovisnost u konačnici nije tako jednostavna.

Broj potrebnih prepoznatljivih djelića prostora ovisi o zadaći mreže, ali prethodno razmatranje intuitivno objašnjava zašto je broj neurona u skrivenom sloju redovito manji od broja ulaznih varijabli ...

Izlaz mreže je odraz (“slika”) ulaza izmijenjena u skladu s definicijom preslikavanja koja je pohranjena u parametrima mreže pa se neuronska mreža može shvatiti kao mapa (“šablona”) po kojoj se na temelju ulaza određuje izlaz. Neuroni u mreži djeluju jedni na druge šaljući svoje izlazne signale na ulaze susjednih neurona, ali svaki neuron u mreži “odlučuje” samostalno, tj. jednom kad primi signale (“mišljenja”) susjednih neurona, svoj izlaz namješta isključivo na temelju svojih parametara. Pri tome svi neuroni istog sloja odluku donose istodobno pa neuronska mreža djeluje kao “paralelni obrađivač” (procesor).

Da bi mreža mogla pravilno preslikavati ulaz u izlaz, moramo joj “objasniti” što želimo, tj. “naučiti ju” kako obaviti željeno preslikavanje. Postoji više načina “podučavanja”, odnosno “uvježbavanja” mreže, a najizravniji je dati joj primjere pravilnog rada.

Neka je $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ preslikavanje (mapa) koju želimo postići. Uvježbavanje mreže je proces tijekom kojeg odgovarajućim algoritmom ugađamo parametre mreže na temelju razlike između njezinog stvarnog izlaznog signala (rezultata) \mathbf{y} i željenog izlaznog signala $\mathbf{y}_d = f(\mathbf{x})$.



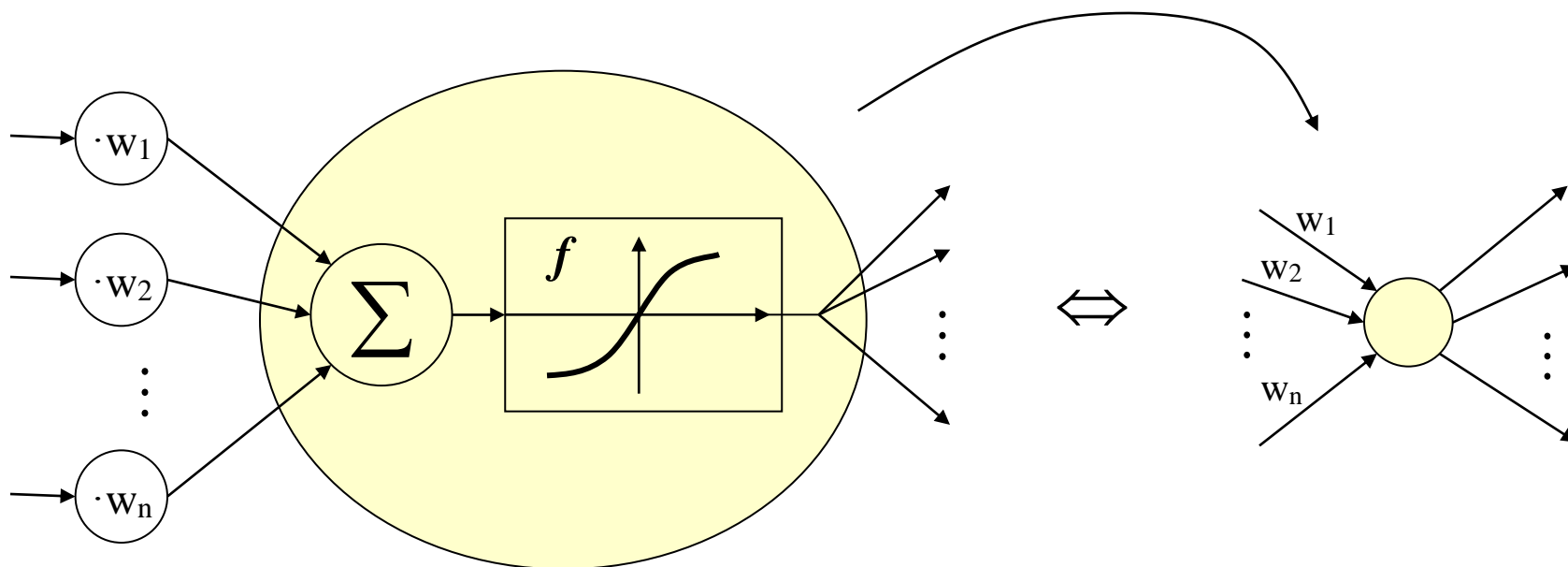
Dakle, mreži dajemo neke od točaka mape f kao primjere željenog preslikavanja, konkretno neki skup od n parova $(x_{d,1}, y_{d,1}), \dots, (x_{d,n}, y_{d,n})$ koji nazivamo *skup za uvježbavanje* (*training set*), i na temelju tog ograničenog skupa primjera određujemo parametre mreže koji će ju osposobiti za predviđene zadaće.

Nakon uvježbavanja očekujemo (!) od mreže pravilno preslikavanje i onih kombinacija ulaznih varijabli koje nisu bile u skupu za uvježbavanje. Prema tome, svrha uvježbavanja mreže, tj. njezina “učenja”, jest njezino “osamostaljivanje” u radu.

Budući da u opisanom načinu učenja mreža ima “učitelja”, taj se pristup naziva *nadzirano učenje* (*supervised learning*). Postojanje tog naziva daje naslutiti da postoji i nenadzirano (*unsupervised*) učenje ... ☺

Uvježbavanje jednog neurona

Neuron je “mapa” preslikavanja strukture kao na slici, dakle više ulaza i samo jedan izlaz, ali koji se može granati. Izlaz neurona je funkcija f zbroja njegovih ulaza, a kako se izlaz često naziva i *aktivnost* (*activity*), izlazna funkcija se naziva *aktivacijska funkcija* (*activation function*).



S obzirom na građu, neuron obavlja preslikavanje $\mathbb{R}^n \rightarrow \mathbb{R}$. Označimo li njegov izlaz s y , a ulazne varijable s x_i , za $f(z) = z$ (tj. izlazne funkcije nema; *Adaline* = *adaptive linear element*) preslikavanje je

$$y = f\left(\sum_{i=1}^n w_i x_i\right) = \sum_{i=1}^n w_i x_i \quad .$$

Uobičajeno je podatke organizirati u stupčane vektore u skladu s prikazom mreže, što znači da će ulazi x_i i težine w_i činiti stupce $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_n]^T$ i $\mathbf{w} = [w_1 \ w_2 \ \dots \ w_n]^T$. Izlaz se tada može izračunati kao skalarni umnožak vektora

$$y = \sum_{i=1}^n w_i x_i = \mathbf{w}^T \mathbf{x} = \mathbf{x}^T \mathbf{w} \quad .$$

Napomena: u statistici je najčešći upravo obrnut zapis, tj. varijable koje čine jedan uzorak (ulaz) se smještaju u isti redak, a uzorci se nižu jedan ispod drugog. No, takav zapis ne odražava intuitivnu sliku neuronske mreže, stoga ćemo koristiti zapis koji smo upravo uveli.

Uvježbavanje takvog neurona s p parova za vježbu $(\mathbf{x}_{d,1}, y_{d,1}), \dots, (\mathbf{x}_{d,p}, y_{d,p})$ je problem oblika

$$\text{minimizirati} \quad \frac{1}{2} \sum_{i=1}^p \left(\mathbf{x}_{d,i}^T \mathbf{w} - y_{d,i} \right)^2$$

pri čemu se minimizacija provodi s obzirom na vektor $\mathbf{w} = [w_1 \dots w_n]^T \in \mathbb{R}^n$.

Izraz koji treba minimizirati je zapravo zbroj kvadrata pogrešaka mreže u svim točkama (parovima) za uvježbavanje, dakle radi se o kriteriju najmanjih kvadrata (LMS = *least mean squares*).

$$\text{minimizirati} \quad \frac{1}{2} \sum_{i=1}^p (y_i - y_{d,i})^2 = \frac{1}{2} \sum_{i=1}^p e_i^2$$

Ako p uzoraka (parova za vježbu) smjestimo u matricu $\mathbf{X}_d = [\mathbf{x}_1 \dots \mathbf{x}_p] \in \mathbb{R}^{n \times p}$ i vektor $\mathbf{y}_d \in \mathbb{R}^{p \times 1}$

$$\mathbf{X}_d = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \quad \text{i} \quad \mathbf{y}_d = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \end{bmatrix},$$

problem se može zapisati kao

$$\text{minimizirati} \quad \frac{1}{2} \|\mathbf{X}_d^T \mathbf{w} - \mathbf{y}_d\|^2 = \frac{1}{2} \|\mathbf{e}\|^2.$$

Izraz koji treba minimizirati je zapravo kvadrat norme vektora pogrešaka $\mathbf{e} = [e_1 \ e_2 \ \dots \ e_p]^T$.

S obzirom na to da minimiziramo kvadrat nekog izraza, jasno je da će teorijski minimum biti nula, dakle cilj je postići da taj izraz, u ovom slučaju norma vektora pogrešaka, postane nula.

Prema tome, optimalne težine w_i jesu rješenje sustava $\mathbf{X}_d^T \cdot \mathbf{w} = \mathbf{y}_d$ od p jednažbi s n nepoznanica ili raspisano

$$x_{11} \cdot w_1 + x_{12} \cdot w_2 + \dots + x_{1n} \cdot w_n = y_1$$

$$x_{21} \cdot w_1 + x_{22} \cdot w_2 + \dots + x_{2n} \cdot w_n = y_2$$

$$\dots \qquad \dots$$

$$x_{p1} \cdot w_1 + x_{p2} \cdot w_2 + \dots + x_{pn} \cdot w_n = y_p \qquad ,$$

tj. želimo savršeno uvježbati neuron tako da bude $\mathbf{y} = \mathbf{y}_d$ u svim točkama za uvježbavanje, ako je to uopće moguće.

Općenito, rješenje sustava $\mathbf{Ax} = \mathbf{b}$ je $\mathbf{x} = \mathbf{A}^- \cdot \mathbf{b}$, gdje je \mathbf{A}^- opći inverz (*generalized inverse*) matrice $\mathbf{A} \in \mathbb{R}^{p \times n}$. Opći inverz je svaka matrica koja zadovoljava relaciju

a) $\mathbf{AA}^- \mathbf{A} = \mathbf{A}$

i takva matrica uvijek postoji, ali nije uvijek jedinstvena. Zato se gotovo uvijek koristi takozvani Moore-Penrose pseudoinverz koji jest jedinstven, a to je matrica koja zadovoljava još sljedeća tri dodatna uvjeta:

b) $\mathbf{A}^- \mathbf{AA}^- = \mathbf{A}^-$

c) $(\mathbf{AA}^-)^T = \mathbf{AA}^-$

d) $(\mathbf{A}^- \mathbf{A})^T = \mathbf{A}^- \mathbf{A}$.

Moore-Penrose pseudoinverz se najčešće označava s \mathbf{A}^+ , ali dosta je zastupljena i oznaka \mathbf{A}^\dagger .

Prilikom rješavanja sustava $\mathbf{A}\mathbf{x} = \mathbf{b}$, $\mathbf{A} \in \mathbb{R}^{p \times n}$ i \mathbf{A} je punog ranga, moguća su tri slučaja:

1. $p = n$: rješenje je jedinstveno i nema optimizacije
2. $p < n$: rješenja ima beskonačno i uobičajeno je odabrati ono najmanje norme $\|\mathbf{x}\|$, a to je upravo rješenje $\mathbf{x} = \mathbf{A}^+ \cdot \mathbf{b}$, pri čemu se pseudoinverz \mathbf{A}^+ računa po formuli (*desni* pseudoinverz)

$$\text{pseudoinverz} = \mathbf{A}^+ = \mathbf{A}^T (\mathbf{A}\mathbf{A}^T)^{-1}$$

što primijenjeno na naš problem ($\mathbf{A} = \mathbf{X}_d^T \in \mathbb{R}^{p \times n}$) daje

$$\mathbf{w}^* = \mathbf{X}_d (\mathbf{X}_d^T \mathbf{X}_d)^{-1} \mathbf{y}_d \quad .$$

Do istog rješenja, ali bez potrebe za invertiranjem matrice $\mathbf{X}_d^T \cdot \mathbf{X}_d$, vodi Kaczmarzov algoritam:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \mu \frac{e^{(k)} \mathbf{x}_{d,R(k)+1}}{\|\mathbf{x}_{d,R(k)+1}\|^2}, \quad k = 0, 1, 2, \dots$$

gdje su:

$$\mathbf{w}^{(0)} = \mathbf{0},$$

$$e^{(k)} = \mathbf{x}_{d,R(k)+1}^T \mathbf{w}^{(k)} - y_{d,R(k)+1},$$

$$0 < \mu < 2$$

i $R(k)$ ostatak dijeljenja k s p , tj. $R(k) = k \bmod p$.

Nedostatak Kaczmarzovog algoritma je taj da je traženi \mathbf{w}^* tek limes za $k \rightarrow \infty$.

3. $p > n$: točno rješenje može, ali i ne mora postojati. Izraz $\mathbf{x} = \mathbf{A}^+ \cdot \mathbf{b}$ je jednak točnom rješenju ako ono postoji, a ako ne, to je optimalno rješenje po kriteriju najmanjih kvadrata odstupanja, odnosno po kriteriju najmanje norme $\|\mathbf{b} - \mathbf{A}\mathbf{x}\|$.

Pseudoinverz \mathbf{A}^+ se računa po formuli (*lijevi pseudoinv.*)

$$\text{pseudoinverz} = \mathbf{A}^+ = \left(\mathbf{A}^T \mathbf{A} \right)^{-1} \mathbf{A}^T$$

pa slijedi (za $\mathbf{A} = \mathbf{X}_d^T \in \mathbb{R}^{p \times n}$)

$$\mathbf{w}^* = \left(\mathbf{X}_d \mathbf{X}_d^T \right)^{-1} \mathbf{X}_d \mathbf{y}_d \quad .$$

Ako neuron ima nelinearnu aktivacijsku funkciju, nema formule za izravno izračunavanje rješenja, nego primjenjujemo razne optimizacijske (iteracijske) algoritme kojima se postupno približavamo optimalnom rješenju.

Na primjer, za $p < n$ i aktivacijsku funkciju f , primjenjujemo Kaczmarzov algoritam koji ostaje jednak, samo se mijenja izraz za pogrešku $e^{(k)}$.

$$e^{(k)} = f\left(\mathbf{x}_{d,R(k)+1}^T \mathbf{w}^{(k)}\right) - y_{d,R(k)+1}$$

Za $p > n$ možemo primijeniti bilo koju optimizacijsku metodu, npr. gradijentnu, Newtonovu, metodu konjugiranih smjerova itd. Najjednostavnija je gradijentna pa ćemo ju i mi koristiti u svim primjerima.

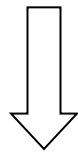
Gradijentna metoda optimizacije

$$f(\mathbf{x}) \quad , \quad f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \cdot \nabla f(\mathbf{x}^{(k)}) \quad ; \quad \alpha_k > 0$$

Ciljna funkcija za *Adaline*:
$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{X}_d^T \mathbf{w} - \mathbf{y}_d\|^2$$

Gradijent ciljne funkcije:
$$\begin{aligned} \nabla f(\mathbf{w}) &= \mathbf{X}_d (\mathbf{X}_d^T \cdot \mathbf{w} - \mathbf{y}_d) \\ &= \mathbf{X}_d \cdot \mathbf{e} \end{aligned}$$



Iteracija:
$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha_k \cdot \mathbf{X}_d \cdot \mathbf{e}^{(k)} \quad ,$$

gdje je
$$\mathbf{e}^{(k)} = \mathbf{X}_d^T \cdot \mathbf{w}^{(k)} - \mathbf{y}_d \quad . \quad (\text{vektor pogreška } k\text{-tog rješenja})$$

Gradijentna metoda zahtijeva izračunavanje gradijenta ciljne funkcije $E(\mathbf{w})$, u našem slučaju

$$\nabla E(\mathbf{w}) = \nabla \frac{1}{2} \sum_{i=1}^p e_i^2 = \frac{1}{2} \sum_{i=1}^p \nabla e_i^2 = \sum_{i=1}^p e_i \cdot \nabla e_i \quad .$$

Gornji izraz pokazuje da se točan gradijent može odrediti tek nakon izračunavanja gradijenta ∇e_i svake pojedine pogreške (u svakoj točki skupa za uvježbavanje).

Prema tome, ispravljanje parametara bi bilo moguće tek nakon što se neuronu “predoči” cijeli skup za uvježbavanje pa se takvo “podučavanje” naziva **skupno uvježbavanje** (*Batch Learning*).

Jedno takvo izračunavanje gradijenta (prolazak svim uzorcima) naziva se *iteracija* ili *epoha*.

No, postoji i druga mogućnost za koju ideja slijedi iz

$$\begin{aligned}\nabla E(\mathbf{w}) &= \nabla \frac{1}{2} \sum_{i=1}^p e_i^2 = \frac{1}{2} \sum_{i=1}^p \nabla e_i^2 = \frac{1}{2} \cdot p \cdot (\nabla e_i^2)_{mean} \\ &= p \cdot (e_i \cdot \nabla e_i)_{mean} \approx p \cdot e_i \cdot \nabla e_i \quad .\end{aligned}$$

Ukupni gradijent kvadrata norme vektora pogrešaka je $p \cdot$ (prosječni gradijent) pa kad bi za svaki par za uvježbavanje umnožak $e_i \cdot \nabla e_i$ bio približno jednak, ispravak parametara nakon “učenja” u jednoj točki

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha'_k \cdot 1/2 \cdot \nabla e_i^2 = \mathbf{w}^{(k)} - \alpha_k \cdot e_i \cdot \nabla e_i$$

bio bi $1/p$ skupnog (ukupnog) ispravka. To, naravno, nije posve točno, ali postupak je djelotvoran i naziva se **koračno uvježbavanje** (*on-line learning*).

Ako se radi o jednom neuronu bez aktivacijske funkcije (*Adaline*), pogreška je $e_i = \mathbf{x}_{d,i}^T \mathbf{w} - y_{d,i}$, njezin gradijent $\nabla e_i = \mathbf{x}_{d,i}$ i formula za ispravak vektora parametara postaje

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha_k \cdot e_i^{(k)} \cdot \mathbf{x}_{d,i} ,$$

odnosno ako se koraci broje kao i točke za uvježbavanje,

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha_k \cdot \mathbf{x}_{d,k} \cdot e_k .$$

Dobiveni izraz opisuje koračno uvježbavanje jednog neurona koje se temelji na minimizaciji kvadrata trenutačne pogreške (u jednoj točki), stoga se algoritam uvježbavanja koji proizlazi iz primjene dobivene formule naziva *algoritam najmanjih kvadrata* (*LMS Algorithm*).

Ovo je bio vrlo pojedonstavnjen izvod! Strogi dokaz ispravnosti (i konvergencije) LMS algoritma polazi od pretpostavke *ergodičnosti* okoline...

Dakle, dva osnovna načina uvježbavanja jesu:

- 1) **skupno** uvježbavanje (*batch learning*) – gradijent se izračunava na temelju svih točaka za uvježbavanje i tek tada se osvježavaju parametri mreže (*epoha po epoha*)
 - memorijski i računski zahtjevnije
 - sigurna i mirnija, ali često sporija konvergencija
- 2) **koračno** uvježbavanje (*on-line learning*) – gradijent se izračunava i parametri mreže osvježavaju u svakoj pojedinoj točki skupa za uvježbavanje (LMS algoritam)
 - konvergencija s oscilacijama (Brownovo gibanje)
 - postupak uvježbavanja bolje prati male promjene podataka (lokalne promjene u prostoru) \Rightarrow lakše izvlačenje iz okoline lokalnog minimuma
 - ponovljeni podatci se bolje iskorištavaju (utjecajniji)

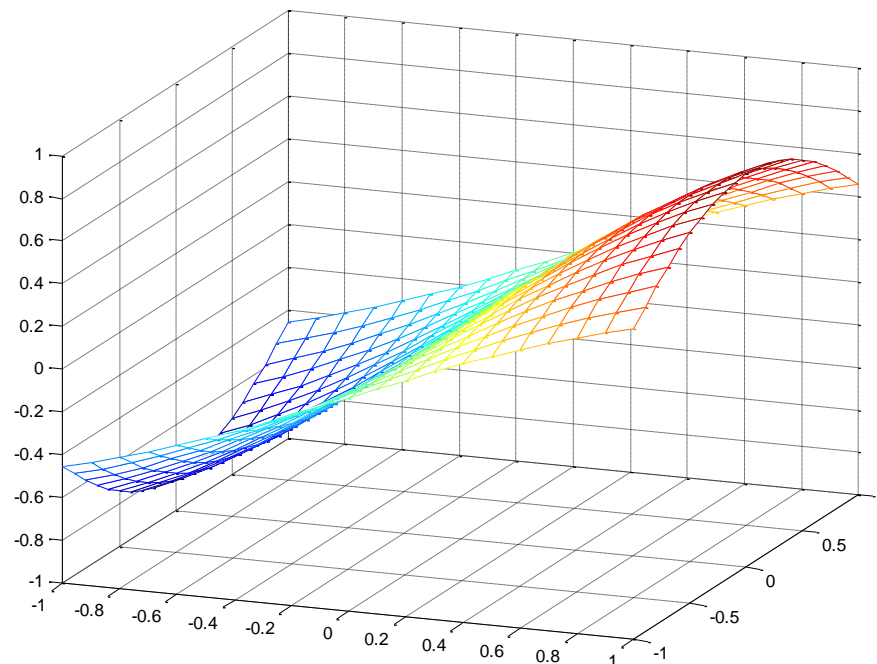
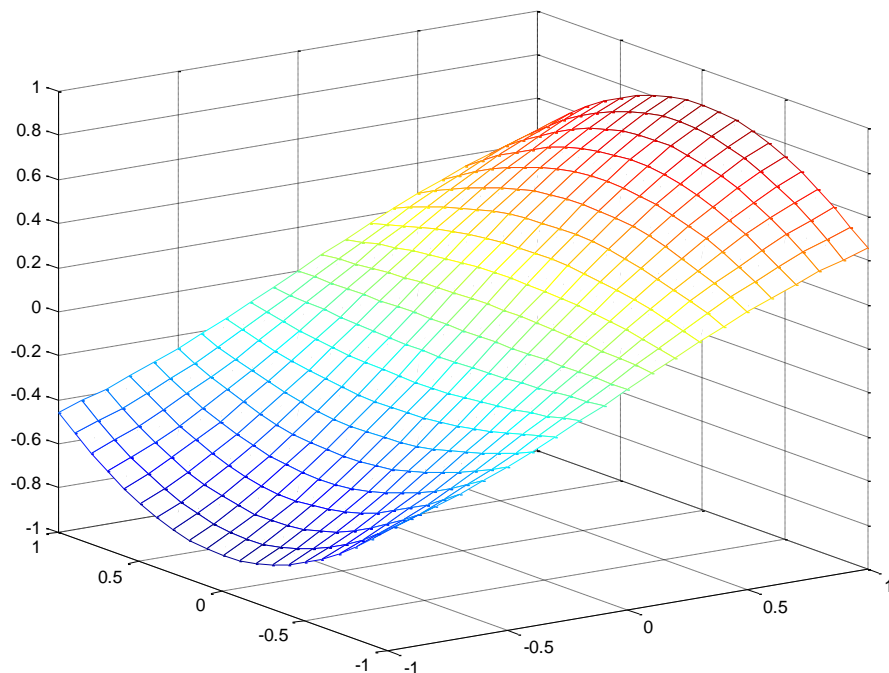
Primjer: uvježbavanje jednog neurona s dva ulaza i bez aktivacijske funkcije (*Adaline*)

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad ; \quad f(x_1, x_2) = \sin(x_1) \cdot \cos(x_2)$$

Domena: $x_1, x_2 \in [-1, 1]$.



SingleNeuron.m

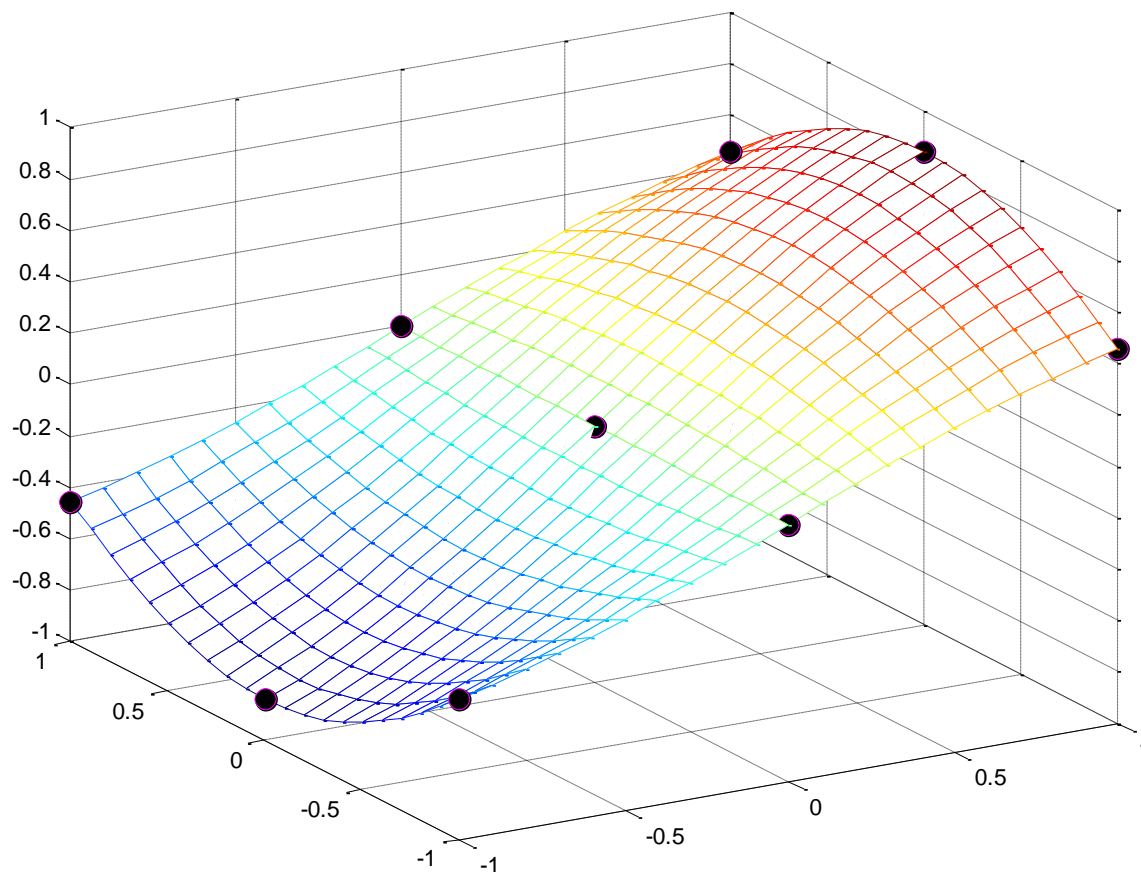


Skup za uvježbavanje: 9 točaka, argumenti funkcije neka su

$$\mathbf{x}_i = [x_{1i} \ x_{2i}]^T \quad ; \quad x_1, x_2 \in \{-1, 0, 1\}.$$

Dakle, $\mathbf{X}_d \in \mathbb{R}^{2 \times 9}$
i $\mathbf{y}_d \in \mathbb{R}^{9 \times 1}$, a cilj
je naći vektor \mathbf{w}
koji će
minimizirati
funkciju

$$\frac{1}{2} \|\mathbf{X}_d^T \mathbf{w} - \mathbf{y}_d\|^2.$$



Tražimo (najbolje moguće) rješenje sustava $\mathbf{X}_d^T \cdot \mathbf{w} = \mathbf{y}_d$. Rješenje je uvijek $(\mathbf{X}_d^T)^- \cdot \mathbf{y}_d$, a kako je broj jednačbi veći od broja nepoznanica, računamo lijevi pseudoinverz, odnosno rješenje je

$$\mathbf{w} = (\mathbf{X}_d \mathbf{X}_d^T)^{-1} \mathbf{X}_d \mathbf{y}_d = \begin{bmatrix} 0,5836 \\ 0 \end{bmatrix} .$$

Do istog rezultata (postupno) dolazimo i gradijentnom metodom. Gradijent ciljne funkcije za *Adaline* je $\mathbf{X}_d(\mathbf{X}_d^T \cdot \mathbf{w} - \mathbf{y}_d)$ pa gradijentna metoda vodi u iteraciju

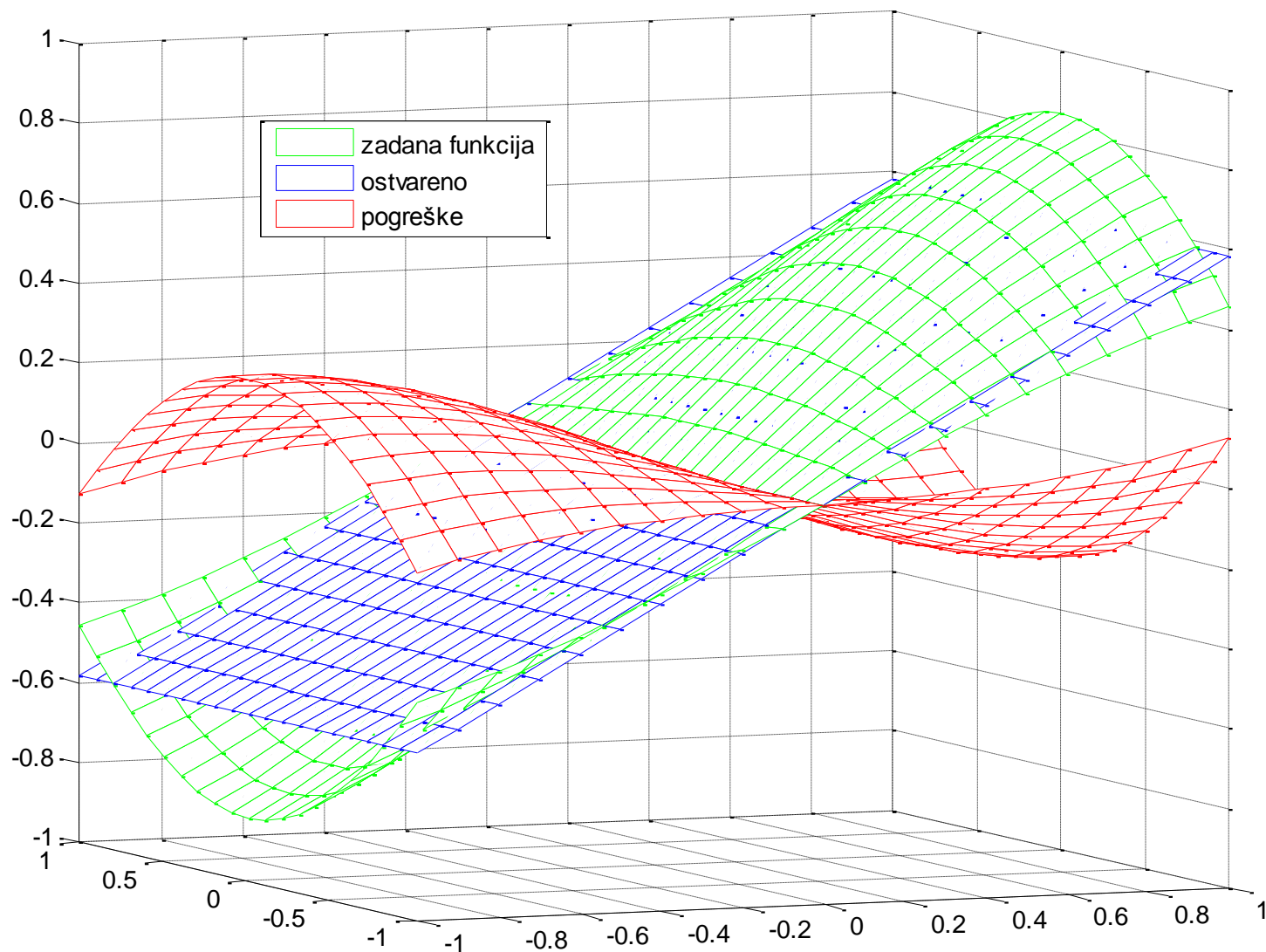
$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha_k \cdot \mathbf{X}_d \cdot \mathbf{e}^{(k)} \quad , \quad (\text{skupno uvježbavanje!})$$

gdje je $\mathbf{e}^{(k)} = \mathbf{X}_d^T \cdot \mathbf{w}^{(k)} - \mathbf{y}_d$ (smisao = pogreška k -tog rješenja).

Konkretno, za $\mathbf{w}^{(0)} = [0 \ 0]^T$ i $\alpha_k = \text{konst.} = 10^{-3}$, nakon 200 iteracija (epoha) imamo $\mathbf{w}^{(200)} = [0,4095 \ 0]^T$.

(do rješenja s pogreškom $< 10^{-6}$ dolazimo nakon 1357 koraka)

Konačno rješenje bilo bi... Nije loše, a riječ je o samo jednom neuronu ...



Prijedlog za vježbu

Prethodna funkcija je bila relativno blago zaobljena (“slična” ravnini) pa je jedan neuron bio dovoljan za njenu aproksimaciju. Međutim, već samo malo složeniju, a vrlo srodnu funkciju

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad ; \quad f(x_1, x_2) = \sin^2(x_1) \cdot \cos^2(x_2), \\ x_1, x_2 \in [-1, 1]$$

jedan neuron ne može pratiti.

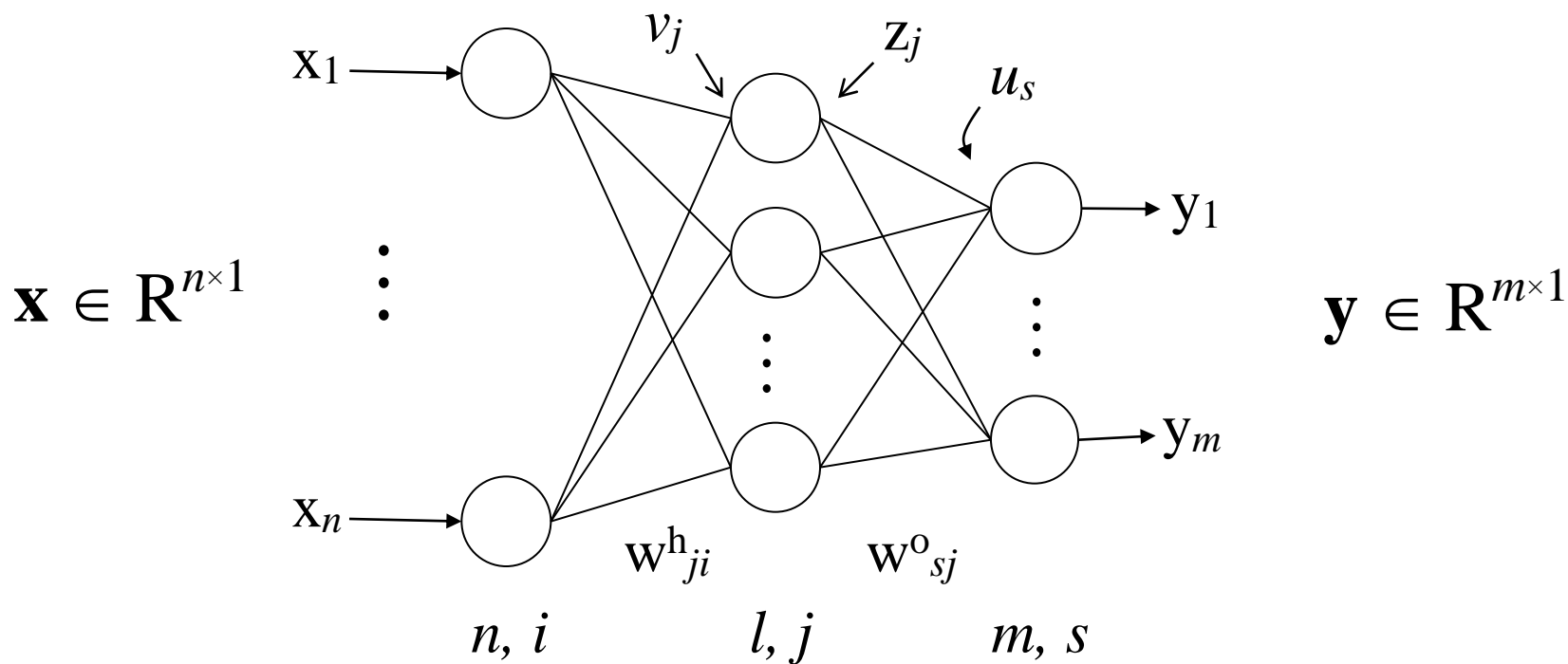
Uvjerite se u tu činjenicu koristeći program kojim je pripremljen prethodni primjer!

Uvježbavanje cijele mreže

Radi preglednosti i razumljivosti izlaganja, ograničit ćemo se na mrežu sa samo jednim skrivenim slojem. Međutim, logika algoritma i građa formula su identični i za složenije mreže pa je analiza koja slijedi valjana općenito.

Imamo n ulaza x_i , $i = 1, \dots, n$. Nadalje, imamo m izlaza y_s , $s = 1, \dots, m$. Skriveni sloj ima l neurona čiji su izlazi z_j , $j = 1, \dots, l$. Aktivacijsku funkciju neurona skrivenog sloja označimo s f^h , a onu izlaznog sloja s f^o . Težine veza iz ulaznog u skriveni sloj su w^h_{ji} , a onih iz skrivenog u izlazni sloj w^o_{sj} .

Takva mreža obavlja (definira) preslikavanje iz \mathbb{R}^n u \mathbb{R}^m .



Odredimo eksplicitnu formulu preslikavanja u jednoj točki.

Ukupni ulaz v_j u j -ti neuron skrivenog sloja (naziva se i *lokalno polje; local field*) i njegov izlaz z_j jesu

$$v_j = \sum_{i=1}^n w_{ji}^h x_i \quad , \quad z_j = f_j^h \left(\sum_{i=1}^n w_{ji}^h x_i \right) = f_j^h(v_j) \quad .$$

Izlaz iz s -tog neurona izlaznog sloja je

$$\begin{aligned} y_s &= f_s^o \left(\sum_{j=1}^l w_{sj}^o z_j \right) = f_s^o \left(\sum_{j=1}^l w_{sj}^o f_j^h(v_j) \right) \\ &= f_s^o \left(\sum_{j=1}^l w_{sj}^o f_j^h \left(\sum_{i=1}^n w_{ji}^h x_i \right) \right) = F_s^o(x_1, x_2, \dots, x_n) \quad , \end{aligned}$$

a iz cijele mreže

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} F_1^o(x_1, \dots, x_n) \\ \vdots \\ F_m^o(x_1, \dots, x_n) \end{bmatrix} = \begin{bmatrix} F_1^o(\mathbf{x}) \\ \vdots \\ F_m^o(\mathbf{x}) \end{bmatrix} = \mathbf{F}^o(\mathbf{x}) \quad .$$

Radi sažetijeg i prikladnijeg zapisa uvodimo

$$\mathbf{W}^h = \begin{bmatrix} w_{11}^h & \cdots & w_{1n}^h \\ \vdots & \ddots & \vdots \\ w_{l1}^h & \cdots & w_{ln}^h \end{bmatrix} \quad \text{i} \quad \mathbf{W}^o = \begin{bmatrix} w_{11}^o & \cdots & w_{1l}^o \\ \vdots & \ddots & \vdots \\ w_{m1}^o & \cdots & w_{ml}^o \end{bmatrix},$$

gdje indeksi elemenata (težinskih faktora w) u matricama imaju smisao “u koji neuron promatranog sloja” (prvi indeks) “iz kojeg neurona prethodnog sloja” (drugi indeks), dakle u redcima su težine svih veza iz prethodnog sloja u jedan te isti neuron promatranog sloja.

Uz uvedene oznake vrijedi $\mathbf{v} = [v_1 \ v_2 \ \dots \ v_l]^T = (\mathbf{W}^h \mathbf{x})$. Ako još uvedemo vektor $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_l]^T$, onda je izlaz s -tog neurona izlaznog sloja jednak $y_s = f^o_s(\mathbf{W}^o_{s*} \mathbf{z})$. Svi izlazi čine vektor $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_m]^T$, a ako su sve $f^o_s = f^o$, proizlazi da je $\mathbf{y} = f^o(\mathbf{W}^o \mathbf{z})$. Zapis se jednako proširuje i za $\mathbf{X} \in \mathbb{R}^{n \times p}$.

Ukratko, izlazi mreže za p točaka na ulazu, smještenih u matricu $\mathbf{X} \in \mathbb{R}^{n \times p}$, čine matricu $\mathbf{Y} \in \mathbb{R}^{m \times p}$ koja se dobiva računanjem redom:

ulazi u skrivene neurone:	$\mathbf{V}^h = \mathbf{W}^h \cdot \mathbf{X}$	$(l \times p)$
izlazi skrivenih neurona:	$\mathbf{Z}^h = f^h(\mathbf{V}^h)$	$(l \times p)$
ulazi u izlazne neurone:	$\mathbf{U}^o = \mathbf{W}^o \cdot \mathbf{Z}^h$	$(m \times p)$
izlazi mreže:	$\mathbf{Y}^o = f^o(\mathbf{V}^o)$	$(m \times p)$

Za p parova za uvježbavanje dobivamo p vektora $\mathbf{e}_i = [e_1, e_2, \dots, e_m]^T$ izlaznih pogrešaka koje smještamo u matricu pogrešaka $\mathbf{e} \in \mathbb{R}^{m \times p}$

$$\mathbf{e} = [\mathbf{e}_1 \quad \mathbf{e}_2 \quad \dots \quad \mathbf{e}_p] = \begin{bmatrix} e_{11} & \dots & e_{1p} \\ \vdots & \ddots & \vdots \\ e_{m1} & \dots & e_{mp} \end{bmatrix}$$

Cilj uvježbavanja je minimizirati zbroj kvadrata normi svih vektora \mathbf{e}_i , odnosno ciljnu funkciju $E(\mathbf{w})$

$$\text{minimizirati } E(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^p \|\mathbf{e}_k\|^2 = \frac{1}{2} \text{trace}(\mathbf{e}^T \cdot \mathbf{e}) \ .$$

Minimizacija se opet provodi gradijentnom metodom, a uz poznati gradijent $\nabla E(\mathbf{w})$, parametri mreže se ispravljaju kao i u svakoj drugoj gradijentnoj optimizaciji.

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha_k \cdot \nabla E(\mathbf{w})$$

Gradijent nalazimo slično kao i u slučaju jednog neurona.

$$\nabla E(\mathbf{w}) = \nabla \frac{1}{2} \sum_{i=1}^p \|\mathbf{e}_i\|^2 = \frac{1}{2} \sum_{i=1}^p \nabla \|\mathbf{e}_i\|^2 = \sum_{i=1}^p \|\mathbf{e}_i\| \cdot \nabla \|\mathbf{e}_i\|$$

Točan gradijent ovisi o pogreškama u svim točkama za uvježbavanje pa bi se mreža trebala uvježbavati po načelu skupnog uvježbavanja.

Skupno uvježbavanje je najpouzdanije jer su svojstva tog postupka (uvjeti i brzina konvergencije) poznati iz opće analize gradijentne optimizacije.

Usprkos tome, u primjeni je češće koračno uvježbavanje (*on-line learning*) zbog već navedenih prednosti.

$$\begin{aligned}\nabla E(\mathbf{w}) &= \nabla \frac{1}{2} \sum_{i=1}^p \|\mathbf{e}_i\|^2 = \frac{1}{2} \sum_{i=1}^p \nabla \|\mathbf{e}_i\|^2 = \frac{1}{2} \cdot p \cdot \left(\nabla \|\mathbf{e}_i\|^2 \right)_{mean} \\ &= p \cdot \left(\|\mathbf{e}_i\| \cdot \nabla \|\mathbf{e}_i\| \right)_{mean} \approx p \cdot \|\mathbf{e}_i\| \cdot \nabla \|\mathbf{e}_i\|\end{aligned}$$

Ukupni gradijent je p puta veći od prosječnog gradijenta pa kad bi u svakom koraku (za svaki par za uvježbavanje) trenutačni umnožak $\|\mathbf{e}_i\| \cdot \nabla \|\mathbf{e}_i\|$ bio približno jednak svojoj prosječnoj vrijednosti, koračni ispravak parametara

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha'_k \cdot 1/2 \cdot \nabla \|\mathbf{e}_i\|^2 = \mathbf{w}^{(k)} - \alpha_k \cdot \|\mathbf{e}_i\| \cdot \nabla \|\mathbf{e}_i\|$$

bio bi $1/p$ skupnog ispravka.

Koračno uvježbavanje ćete, uglavnom, primjenjivati i na ispitima. ☺

Algoritam rasprostiranja (prenošenja) unatrag (backpropagation algorithm)

Za oba postupka uvježbavanja trebamo gradijent kvadrata norme $\nabla \|\mathbf{e}_i\|^2$ izlazne pogreške u jednoj točki. U osnovi, problem je ekvivalentan minimizaciji trenutačne pogreške

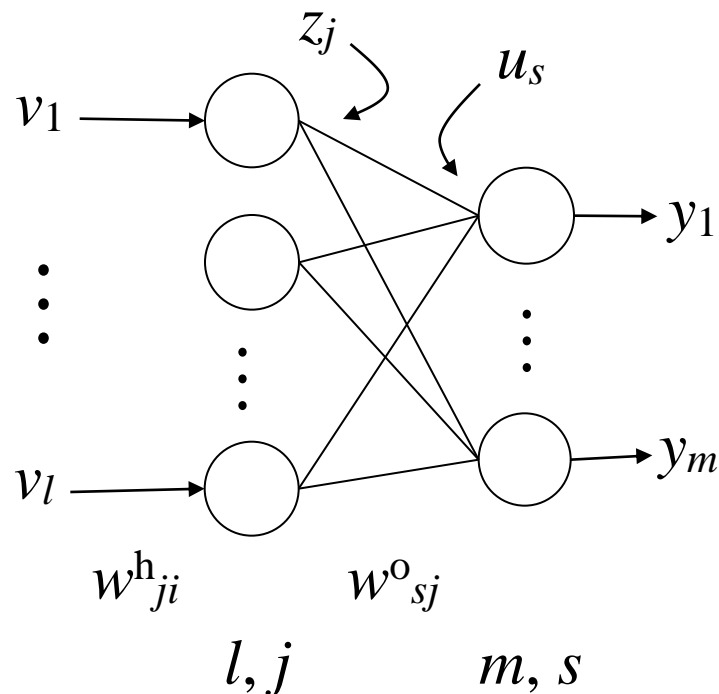
$$E_i(\mathbf{w}) = \frac{1}{2} \|\mathbf{e}_i\|^2 = \frac{1}{2} \|\mathbf{y}_i - \mathbf{y}_{d,i}\|^2 = \frac{1}{2} \sum_{s=1}^m (y_{s,i} - y_{ds,i})^2$$

s obzirom na parametre \mathbf{w} . Izračunavanje gradijenta $\nabla E_i(\mathbf{w})$ iziskuje izračunavanje svih parcijalnih derivacija kvadrata norme $\|\mathbf{e}_i\|^2$ (po svim w_{ij} , za trenutačni \mathbf{w})...

\Rightarrow izračunljivo za mrežu s jednim skrivenim slojem, ali što ako skrivenih slojeva ima više, s desetcima ili stotinama čvorova u svakoj razini?!?

Treba uočiti da je neuronska mreža “prirodno rekurzivna” struktura; iz sloja u sloj građena je i funkcionira jednako. Uspijemo li izračunati parcijalne derivacije po težinama između dva susjedna sloja samo na temelju stanja (aktivnosti) čvorova u ta dva sloja, na isti ćemo način moći izračunati cijeli gradijent!

Promotrimo zadnja dva sloja:



v_j = ukupni ulaz u j -ti neuron prethodnog sloja

z_j = izlaz (*aktivnost*) j -tog neurona prethodnog sloja

u_s = ukupni ulaz u s -ti neuron promatranog sloja

y_s = izlaz (*aktivnost*) s -tog neurona promatranog sloja

1) Derivacija s obzirom na izlaz (aktivnost) izlaznog sloja

$$E_i(\mathbf{w}) = \frac{1}{2} \sum_{s=1}^m (y_{s,i} - y_{ds,i})^2 \Rightarrow EA_s^o = \frac{\partial E_i}{\partial y_s} = y_{s,i} - y_{ds,i} = (\mathbf{e}_i)_s$$

Sve EA^o u i -toj točki, s obzirom na sve izlaze, smještamo u stupčani vektor $\mathbf{EA}^o = \mathbf{y}_i - \mathbf{y}_{di}$ dimenzija $(m \times 1)$.

2) Derivacija s obzirom na uzlaz u izlazni sloj

$$EI_s = \frac{\partial E_i}{\partial u_s} = \frac{\partial E_i}{\partial y_s} \frac{\partial y_s}{\partial u_s} = EA_s^o \cdot f'^o(u_s) = EA_s^o \cdot y_s(1 - y_s)$$

Vektorski:

$$\underset{(m \times 1)}{\mathbf{EI}} = \underset{(m \times 1)}{\mathbf{EA}^o} \cdot \underset{(m \times 1)}{\mathbf{y}} \cdot \underset{(m \times 1)}{(\mathbf{1} - \mathbf{y})}$$

Pozor!

3) Derivacija s obzirom na težinu poveznica s prethodnim slojem

$$EW_{sj} = \frac{\partial E_i}{\partial w_{sj}} = \frac{\partial E_i}{\partial u_s} \frac{\partial u_s}{\partial w_{sj}} = EI_s \cdot z_j$$

Uz $\mathbf{z} = [z_1 \ z_2 \ \dots \ z_l]^T$, sve EW , s obzirom na sve w_{sj} , čine matricu $\mathbf{EW}^o = \mathbf{EI} \cdot \mathbf{z}^T$ dimenzija $(m \times l)$.

4) Derivacija s obzirom na izlaz prethodnog sloja

$$EA_j^h = \frac{\partial E_i}{\partial z_j} = \sum_s \frac{\partial E_i}{\partial u_s} \frac{\partial u_s}{\partial z_j} = \sum_s EI_s w_{sj} = \mathbf{EI}^T \cdot \mathbf{W}_{*j}^o$$

Matrično:

$$\mathbf{EA}_{(l \times 1)}^h = \left(\mathbf{EI}_{(1 \times m)}^T \cdot \mathbf{W}_{(m \times l)}^o \right)^T = \left(\mathbf{W}_{(l \times m)}^o \right)^T \cdot \mathbf{EI}_{(m \times 1)} \quad .$$

Posljednja formula zatvara krug!

Krenuvši od derivacije \mathbf{EA}^o s obzirom na izlaz mreže, izračunali smo derivaciju \mathbf{EA}^h s obzirom na izlaz prethodnog sloja, a “putem” i derivaciju \mathbf{EW}^o s obzirom na težine poveznica, koju smo zapravo i trebali.

Ponavljajući korake 2...4, derivacije s obzirom na izlaze promatranog sloja prevodimo u derivacije s obzirom na izlaze njemu prethodnog sloja i taj se postupak može provesti nad proizvoljnim brojem slojeva. Pritom, prelazeći iz sloja u sloj, u trećem koraku nalazimo potrebne derivacije s obzirom na težine poveznica.

Tijekom opisanog postupka, pogreška izlaza se prenosi u (rasprostire, pronosi kroz) sve prethodne slojeve i to je osnovno obilježje tog algoritma po kojem je i nazvan.

Za one koji ne vole matrice... ☺

Općenito, derivacija kvadrata norme izlazne pogreške po parametru w_{ji} je (j je ovdje indeks u promatranom sloju)

$$\mathbf{E} \mathbf{W}_{ji} = \partial E(\mathbf{w}) / \partial w_{ji} = \delta_j \cdot z_i \quad (z_i \text{ je izlaz } i\text{-tog prethodnika}) .$$

Faktor δ (u izvodu EI) se naziva *lokalni gradijent* i računa:

$$\delta_j = f'(v_j) \cdot (y_j - y_{dj}) \quad ; \quad j = \text{izlazni sloj}$$

ili

$$\delta_j = f'(v_j) \sum_k \delta_k w_{kj} \quad ; \quad j \neq \text{izlazni sloj} ,$$

gdje je k indeks neurona u sljedećem sloju. Za uvijek istu stopu učenja $\alpha_k = \text{konst.} = \eta$, ispravak parametra je

$\Delta w_{ji} = -\eta \cdot \mathbf{E} \mathbf{W}_{ji} = -\eta \cdot \delta_j \cdot z_i$, odnosno osvježeni parametri su

$$(w_{ji})_{\text{novi}} = (w_{ji})_{\text{trenutačni}} + \Delta w_{ji} = (w_{ji})_{\text{trenutačni}} - \eta \cdot \delta_j \cdot z_i .$$

Ta se relacija naziva *delta pravilo*.

Struktura backpropagation algoritma može se dobro uočiti na primjeru troslojne mreže pa ćemo za nju izvesti eksplicitnu formulu gradijenta.

Radi jednostavnosti zapisa, neka je zadan samo jedan par $(\mathbf{x}_d, \mathbf{y}_d)$ za uvježbavanje, $\mathbf{x}_d \in \mathbb{R}^{n \times 1}$ i $\mathbf{y}_d \in \mathbb{R}^{m \times 1}$. Problem određivanja optimalnih parametara mreže je

$$\text{minimizirati} \quad E(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{y}_d\|^2 = \frac{1}{2} \sum_{s=1}^m (y_s - y_{ds})^2 ,$$

a optimizacija se provodi po svim elementima \mathbf{w}^h i \mathbf{w}^o .

Za formulaciju algoritma trebamo gradijent $E(\mathbf{w})$, tj. parcijalne derivacije $E(\mathbf{w})$ po svim w^h_{ji} i w^o_{sj} .

U troslojnoj mreži derivacije po w_{sj}^o su

$$\frac{\partial E(\mathbf{w})}{\partial w_{sj}^o} = (y_s - y_{ds}) f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q \right) z_j = \mathbf{EA}_s^o \cdot f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q \right) z_j ,$$

$$z_q = f_q^h \left(\sum_{i=1}^n w_{qi}^h x_{di} \right) = f_q^h(v_q) \quad .$$

Prepoznamo faktor δ :

$$\delta_s^o = (y_s - y_{ds}) f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q \right) = (y_s - y_{ds}) f_s'^o(u_s) \quad = \mathbf{EI}^o$$

Taj faktor predstavlja skaliranu (pomnoženu nekim faktorom) pogrešku s -tog neurona izlaznog sloja.

Sada se parcijalna derivacija $E(\mathbf{w})$ po w^o_{sj} može kratko pisati

$$\frac{\partial E(\mathbf{w})}{\partial w^o_{sj}} = \delta^o_s z_j \quad . \quad = \mathbf{E} \mathbf{W}^o_{sj}$$

Nadalje, derivacija $E(\mathbf{w})$ po w^h_{ji} je

$$\frac{\partial E(\mathbf{w})}{\partial w^h_{ji}} = \sum_{p=1}^m (y_p - y_{dp}) f_p'^o \left(\sum_{q=1}^l w^o_{pq} z_q \right) w^o_{pj} f_j'^h \left(\sum_{r=1}^n w^h_{jr} x_{dr} \right) x_{di}$$

ili s prethodno uvedenim oznakama

$$\frac{\partial E(\mathbf{w})}{\partial w^h_{ji}} = \underbrace{\left(\sum_{p=1}^m \delta^o_p w^o_{pj} \right)}_{\mathbf{E} \mathbf{A}^h_j} \underbrace{f_j'^h(v_j)}_{f'^o(u_j)} x_{di} = \delta^h_j \cdot x_{di} \quad .$$

Algoritam rasprostiranja unatrag (*backpropagation*) za mrežu s jednim skrivenim slojem:

$$w_{sj}^{o(k+1)} = w_{sj}^{o(k)} - \eta \frac{\partial E(\mathbf{w}^{(k)})}{\partial w_{sj}^o} = w_{sj}^{o(k)} - \eta \delta_s^{o(k)} z_j^{(k)}$$

$$w_{ji}^{h(k+1)} = w_{ji}^{h(k)} - \eta \frac{\partial E(\mathbf{w}^{(k)})}{\partial w_{ji}^h} = w_{ji}^{h(k)} - \eta \delta_j^{h(k)} x_{di}$$

$$= w_{ji}^{h(k)} - \eta \left(\sum_{p=1}^m \delta_p^{o(k)} w_{pj}^{o(k)} \right) f_j'^h(v_j^{(k)}) x_{di} \quad ,$$

pri čemu je η (konstantni) korak u smjeru gradijenta. Sve druge oznake uvedene su prethodno, ali sada ih pregledno navodimo sve zajedno.

Oznake u algoritmu rasprostiranja unatrag (troslojna mreža):

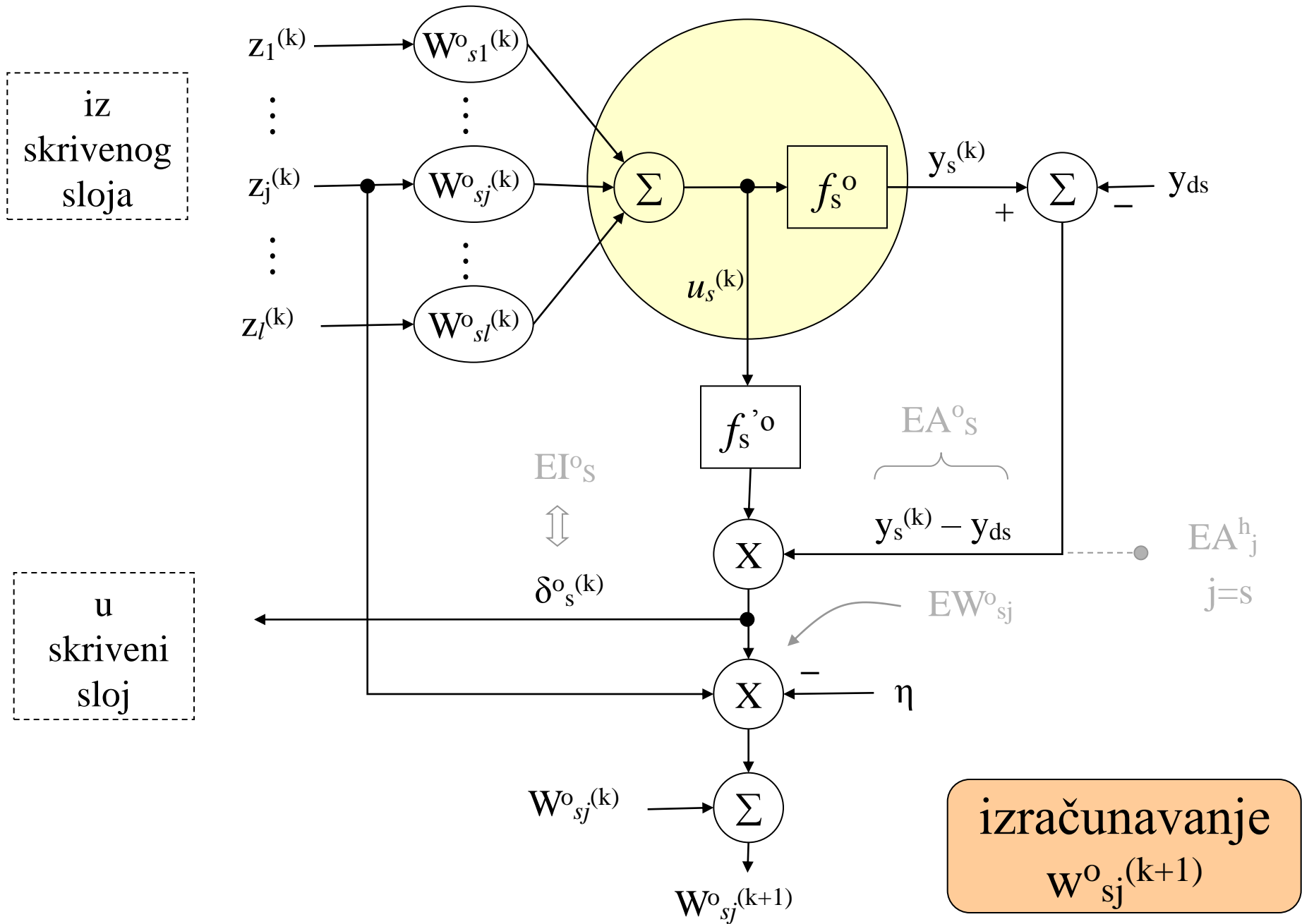
$$v_j^{(k)} = \sum_{i=1}^n w_{ji}^{h(k)} x_{di} \quad , \quad z_j^{(k)} = f_j^h \left(\sum_{i=1}^n w_{ji}^{h(k)} x_{di} \right) = f_j^h \left(v_j^{(k)} \right)$$

$$y_s^{(k)} = f_s^o \left(\sum_{q=1}^l w_{sq}^{o(k)} z_q^{(k)} \right) = f_s^o \left(u_s^{(k)} \right) \quad ,$$

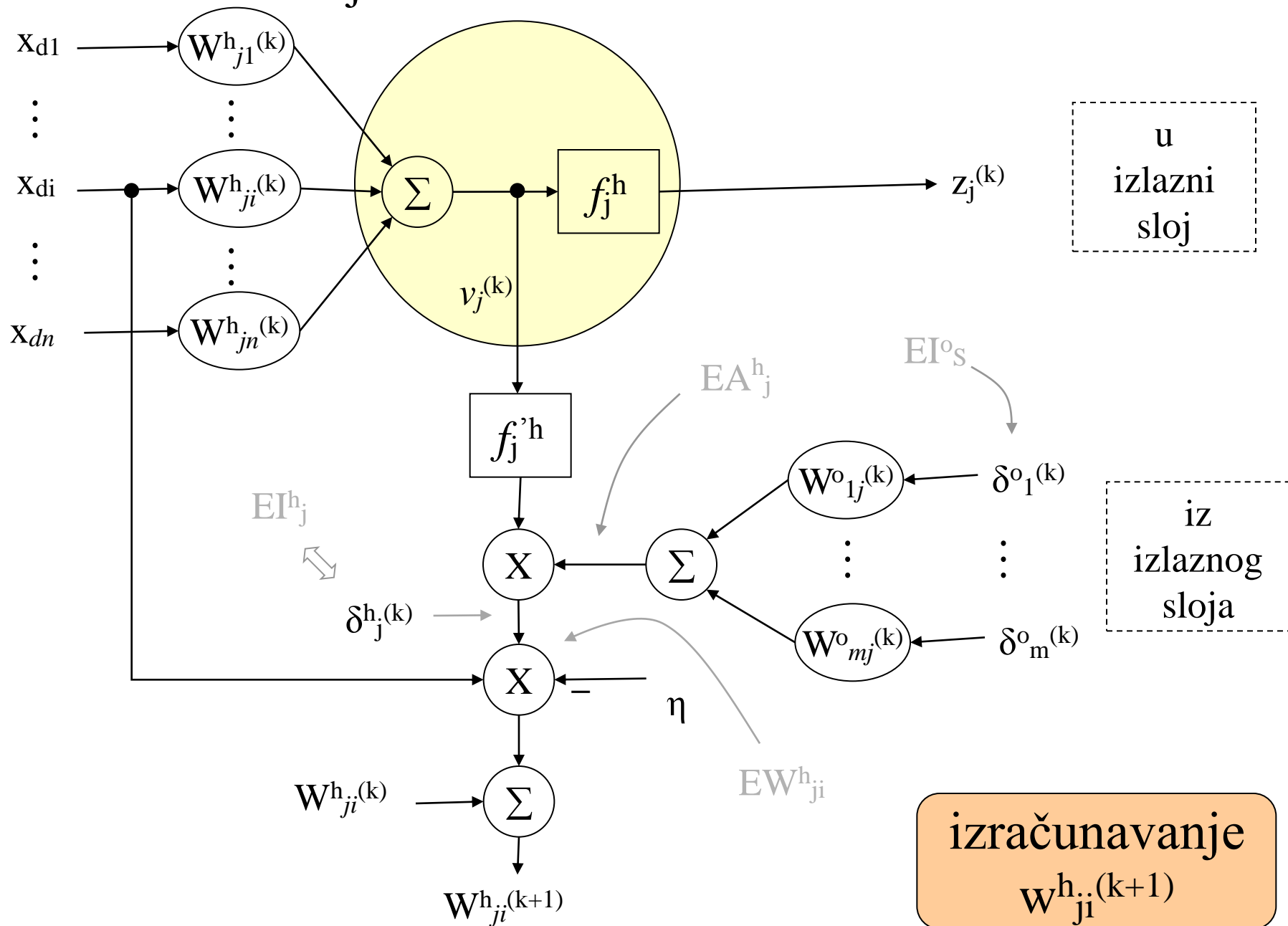
$$\delta_s^{o(k)} = f_s^{'o} \left(u_s^{(k)} \right) \cdot (y_s^{(k)} - y_{ds}) \quad ,$$

$$\delta_j^{h(k)} = f_j^{'h} \left(v_j^{(k)} \right) \cdot \sum_{p=1}^m \delta_p^{o(k)} w_{pj}^{o(k)} \quad .$$

s-ti izlazni neuron



j-ti skriveni neuron



Kao što se vidi iz formula i dijagrama, algoritam rasprostiranja unatrag ima dvije faze:

1. na temelju ulaza x_{di} i trenutanih parametara mreže (težina \mathbf{w}^h i \mathbf{w}^o) izračunamo redom $v_j^{(k)}$, $z_j^{(k)}$, $y_s^{(k)}$ i $\delta_s^{(k)}$. Ova se faza naziva *prolazak prema naprijed* (*forward pass*) jer se utjecaj (stanje) ulaza prenosi sve do izlaza.
2. iz veličina izračunanih tijekom prve faze računamo nove parametre mreže ($w_{sj}^{o(k+1)}$ i $w_{ji}^{h(k+1)}$). Ova faza se naziva *prolazak unazad* ili *obrnuti prolazak* (*backward or reverse pass*) jer se izlazne pogreške $\delta_s^{(k)}$ prenose (rasprostiru) unatrag, od izlaza prema ulazu.

Algoritam ima ista ova svojstva i u složenim mrežama koje imaju višerazinski skriveni sloj.

Samostalni neuron može i ne mora imati aktivacijsku funkciju. A neuroni u mreži?

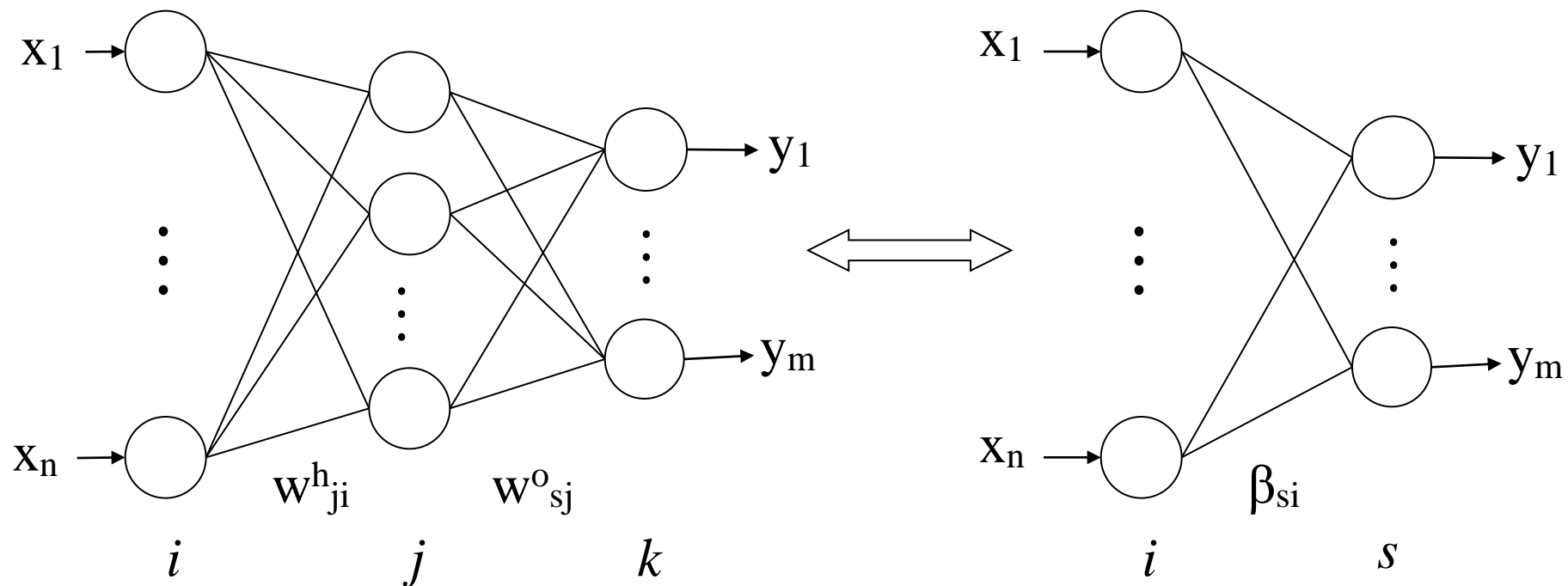
Moraju, i to nelinearnu, inače mreža ne bi bila ono što se čini da je, tj. izgubila bi cijeli skriveni sloj!

$$y_s = \sum_{j=1}^l w_{sj}^o \sum_{i=1}^n w_{ji}^h x_i = \sum_{j=1}^l 1 \cdot \sum_{i=1}^n w_{sj}^o w_{ji}^h x_i = \sum_{j=1}^l 1 \cdot \sum_{i=1}^n \alpha_{sji} x_i$$

To bi bila mreža u kojoj sve poveznice skrivenog i izlaznog sloja imaju težinu =1, a poveznice ulaznog i skrivenog sloja težine α_{sji} . Izlazne težine =1 znače da veze skrivenog i izlaznog sloja nisu “aktivne”, tj. mreža ima samo jedan sloj prilagodljivih poveznica (iz ulaza u prvi skriveni sloj) i ponaša se kao da je dvoslojna: ulaz-izlaz.

Gubitak cijelog skrivenog sloja (svih razina) je još uočljiviji ako prethodni izraz napišemo malo drugačije.

$$y_s = \sum_{i=1}^n \sum_{j=1}^l w_{sj}^o w_{ji}^h x_i = \sum_{i=1}^n x_i \sum_{j=1}^l w_{sj}^o w_{ji}^h = \sum_{i=1}^n x_i \beta_{si}$$

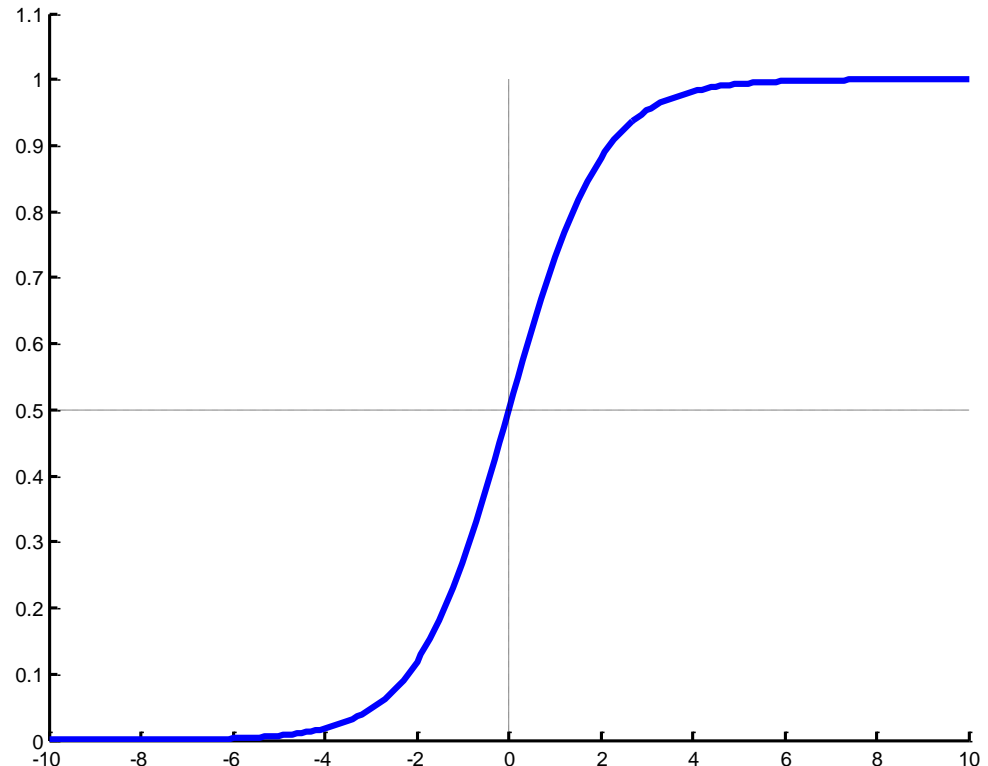


Očito je da aktivacijska funkcija mora biti nelinearna ako se želi izbjeći gubitak skrivenog sloja, a da bi se mogla primijeniti gradijentna optimizacijska metoda, aktivacijska funkcija mora biti i derivabilna.

Osobito prikladna iz više razloga i najčešće primjenjivana funkcija je *sigmoid* čija definicijska formula glasi

$$f(x) = \frac{1}{1 + e^{-x}} \quad ,$$

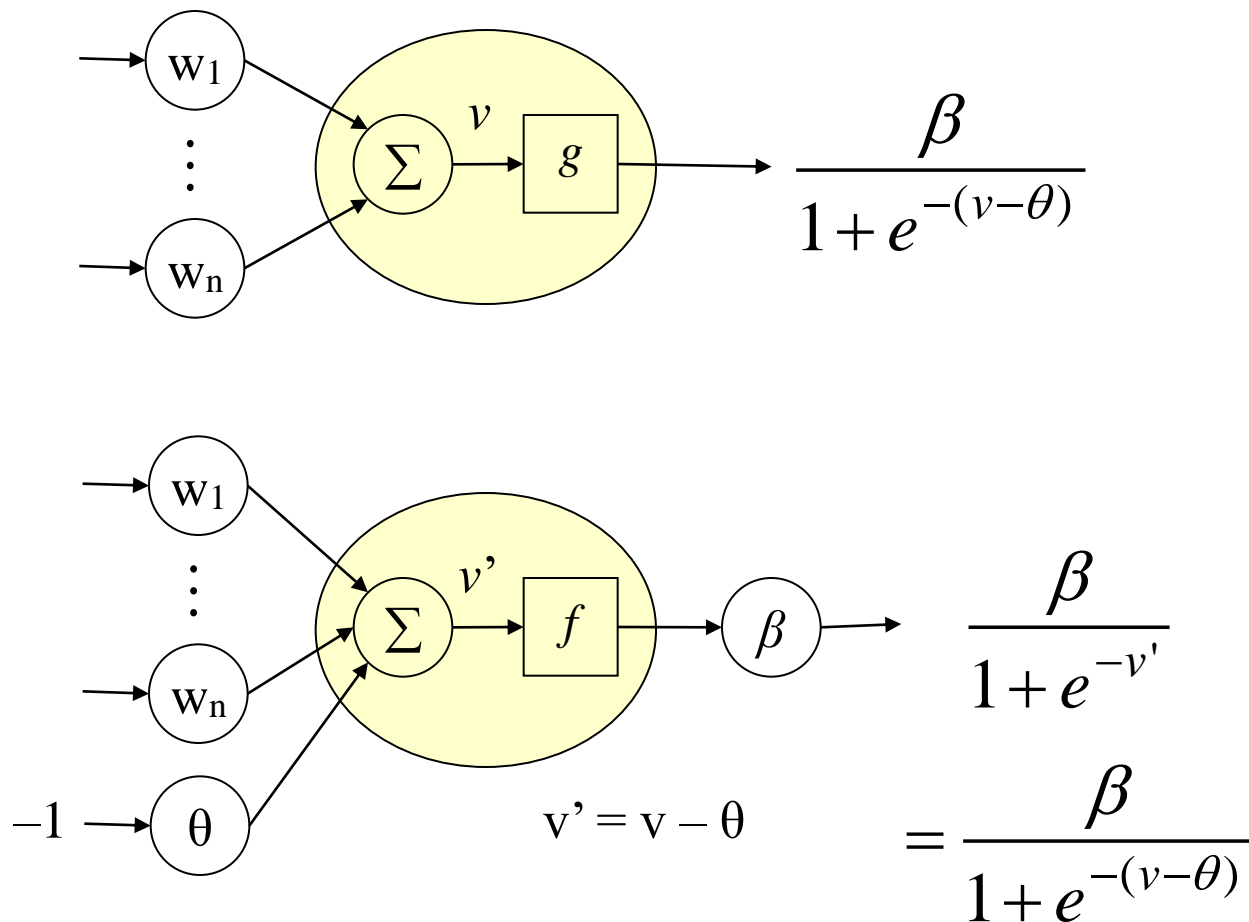
$$\frac{df(x)}{dx} = f(x)[1 - f(x)] \quad .$$



Opći oblik te funkcije je

$$g(x) = \frac{\beta}{1 + e^{-(x-\theta)}} \quad , \quad g'(x) = g(x) \left[1 - \frac{g(x)}{\beta} \right]$$

i ona za svaki neuron može imati drugačije parametre, ali svi se ti parametri lako ugrađuju u algoritam rasprostiranja unatrag.



Parametri β zapravo ne mogu postojati samostalno jer:

1. β neurona i težinu w poveznice sa sljedećim slojem nije moguće razlučiti tijekom optimizacije jer izlaz mreže ovisi o njihovom umnošku, dakle o umnošku dva neovisna faktora, pa rješenja za kombinaciju ta dva faktora ima beskonačno, a samim time beskonačno je i rješenja za cijeli skup parametara
2. izlaz ovisi o umnošku faktora β izlaznog neurona i funkcije svih drugih parametara, dakle opet o dva neovisna faktora pa ih ni u tom umnošku nije moguće razlučiti.

To znači da se mrežu može uvježbati samo za rezultate u opsegu (0,1) pa svaki problem treba svesti na taj opseg. Ili, ukloniti nelinearnost (sigmoide) izlaznim neuronima...

Eksplisitne formule za troslojnu mrežu kada su aktivacijske funkcije sigmoidi s pomakom (tj. sigmoidi s θ , a svi $\beta = 1$):

$$\begin{aligned}
 y_s &= g_s^o \left(\sum_{j=1}^l w_{sj}^o g_j^h \left(\sum_{i=1}^n w_{ji}^h x_i \right) \right) \\
 &= f_s^o \left(\sum_{j=1}^l w_{sj}^o f_j^h \left(\sum_{i=1}^n w_{ji}^h x_i - \theta_j^h \right) - \theta_s^o \right) \\
 &= f_s^o \left(\sum_{j=1}^l w_{sj}^o f_j^h (v_j) - \theta_s^o \right) = f_s^o \left(\sum_{j=1}^l w_{sj}^o z_j - \theta_s^o \right) ,
 \end{aligned}$$

$$v_j = \sum_{i=1}^n w_{ji}^h x_i - \theta_j^h ,$$

$$z_j = f_j^h \left(\sum_{i=1}^n w_{ji}^h x_i - \theta_j^h \right) = f_j^h (v_j) .$$

Parcijalne derivacije su:

$$\frac{\partial E(\mathbf{w})}{\partial w_{sj}^o} = (y_s - y_{ds}) f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q - \theta_s^o \right) z_j = \delta_s^o z_j \quad ,$$

$$\delta_s^o = (y_s - y_{ds}) f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q - \theta_s^o \right) \quad ,$$

$$\frac{\partial E(\mathbf{w})}{\partial w_{ji}^h} = \sum_{p=1}^m (y_p - y_{dp}) f_p'^o \left(\sum_{q=1}^l w_{pq}^o z_q - \theta_p^o \right) w_{pj}^o f_j'^h \left(\sum_{r=1}^n w_{jr}^h x_{dr} - \theta_j^h \right) x_{di}$$

$$= x_{di} f_j'^h(v_j) \sum_{p=1}^m \delta_p^o w_{pj}^o = x_{di} \delta_j^h \quad ,$$

$$\frac{\partial E(\mathbf{w})}{\partial \theta_s^o} = -(y_s - y_{ds}) f_s'^o \left(\sum_{q=1}^l w_{sq}^o z_q - \theta_s^o \right) = -\delta_s^o \quad ,$$

$$\frac{\partial E(\mathbf{w})}{\partial \theta_j^h} = - \sum_{p=1}^m (y_p - y_{dp}) f_p'^o \left(\sum_{q=1}^l w_{pq}^o z_q - \theta_p^o \right) w_{pj}^o f_j'^h \left(\sum_{r=1}^n w_{jr}^h x_{dr} - \theta_j^h \right)$$

$$= -f_j'^h(v_j) \sum_{p=1}^m \delta_p^o w_{pj}^o = -\delta_j^h \quad .$$

Prolazak prema naprijed matrično

$$\underset{(n_h \times 1)}{\mathbf{v}} = \underset{(n_h \times n_i)}{\mathbf{W}^h} \cdot \underset{(n_i \times 1)}{\mathbf{x}} - \underset{(n_h \times 1)}{\boldsymbol{\theta}^h}$$

$$\underset{(n_h \times 1)}{\mathbf{z}} = \textit{sigmoid} \left(\underset{(n_h \times 1)}{\mathbf{v}} \right)$$

$$\underset{(n_o \times 1)}{\mathbf{u}} = \underset{(n_o \times n_h)}{\mathbf{W}^o} \cdot \underset{(n_h \times 1)}{\mathbf{z}} - \underset{(n_o \times 1)}{\boldsymbol{\theta}^o}$$

$$\underset{(n_o \times 1)}{\mathbf{y}} = \textit{sigmoid} \left(\underset{(n_o \times 1)}{\mathbf{u}} \right)$$

Rasprostiranje unatrag matrično

$$\underset{(n_o \times 1)}{\mathbf{E}\mathbf{A}^o} = \underset{(n_o \times 1)}{\mathbf{y}} - \underset{(n_o \times 1)}{\mathbf{y}_d}$$

$$\underset{(n_o \times 1)}{\mathbf{E}\mathbf{I}^o} = \underset{(n_o \times 1)}{\mathbf{E}\mathbf{A}^o} \cdot * \underset{(n_o \times 1)}{\mathbf{y}} \cdot * \left(\underset{(n_o \times 1)}{\mathbf{1}} - \underset{(n_o \times 1)}{\mathbf{y}} \right)$$

$$\underset{(n_o \times 1)}{\boldsymbol{\delta}^o} = \underset{(n_o \times 1)}{\mathbf{E}\mathbf{I}^o}$$

$\cdot *$ = množenje
po parovima,
a ne matrično!

$$\mathbf{E}\mathbf{W}^o = \underset{(n_o \times n_h)}{\boldsymbol{\delta}^o} \cdot \underset{(n_o \times 1)}{\mathbf{z}}^T = \underset{(n_o \times 1)}{\mathbf{E}\mathbf{I}^o} \cdot \underset{(1 \times n_h)}{\mathbf{z}}^T$$

$$\mathbf{E}\mathbf{A}^h = \underset{(n_h \times 1)}{(\mathbf{W}^o)^T} \cdot \underset{(n_h \times n_o)}{\mathbf{E}\mathbf{I}^o}$$

$$\mathbf{E}\mathbf{I}^h = \underset{(n_h \times 1)}{\mathbf{E}\mathbf{A}^h} \cdot \underset{(n_h \times 1)}{*} \mathbf{z} \cdot \underset{(n_h \times 1)}{*} \left(\underset{(n_h \times 1)}{\mathbf{1}} - \underset{(n_h \times 1)}{\mathbf{z}} \right)$$

$$\underset{(n_h \times 1)}{\boldsymbol{\delta}^h} = \underset{(n_h \times 1)}{\mathbf{E}\mathbf{I}^h}$$

$$\mathbf{E}\mathbf{W}^h = \underset{(n_h \times n_i)}{\boldsymbol{\delta}^h} \cdot \underset{(n_h \times 1)}{\mathbf{x}}^T = \underset{(n_h \times 1)}{\mathbf{E}\mathbf{I}^h} \cdot \underset{(1 \times n_i)}{\mathbf{x}}^T$$

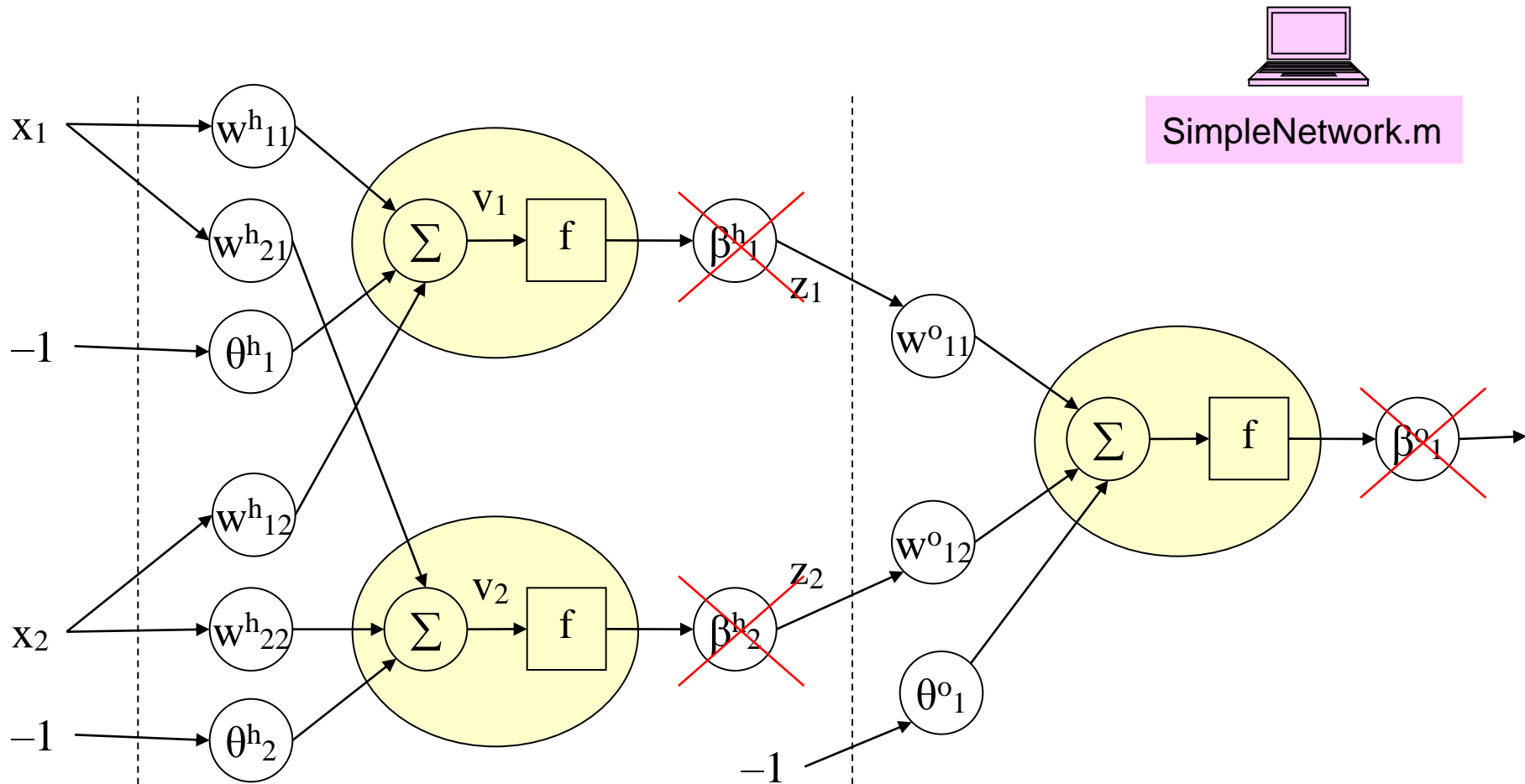
$$\underset{(n_o \times 1)}{\mathbf{E}\boldsymbol{\Theta}^o} = - \underset{(n_o \times 1)}{\mathbf{E}\mathbf{I}^o}$$

$$\underset{(n_h \times 1)}{\mathbf{E}\boldsymbol{\Theta}^h} = - \underset{(n_h \times 1)}{\mathbf{E}\mathbf{I}^h}$$

$\cdot *$ = množenje
po parovima,
a ne matrično!

Primjer: uvježbavanje mreže $2 \times 2 \times 1$

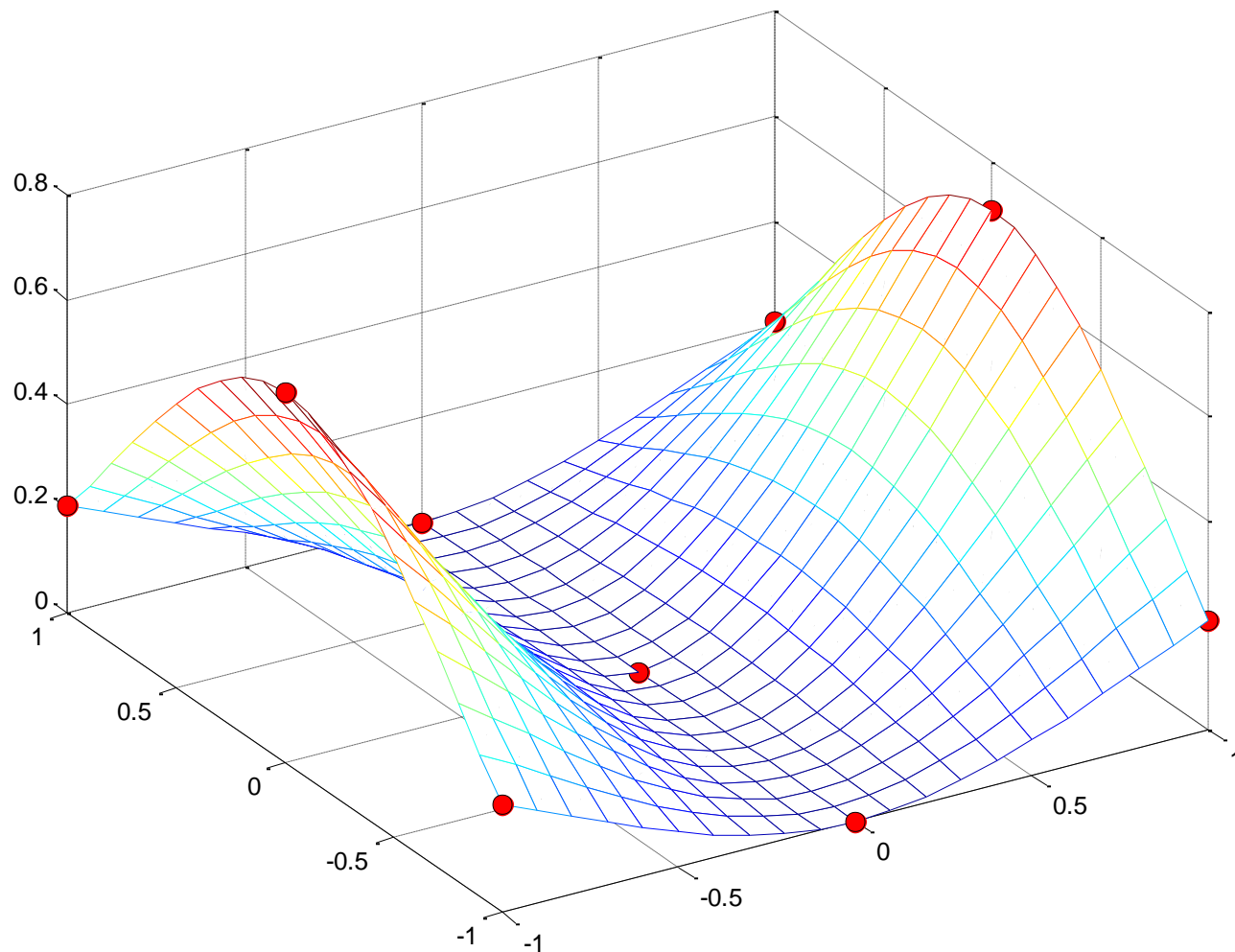
Mreža ima dva ulaza, dva neurona u skrivenom sloju i jedan izlaz, a aktivacijska funkcija je (opći) sigmoid.



Zadaća mreže je aproksimirati funkciju

$$f: \mathbb{R}^2 \rightarrow \mathbb{R} \quad ; \quad f(x_1, x_2) = \sin^2(x_1) \cdot \cos^2(x_2),$$
$$x_1, x_2 \in [-1, 1]$$

koju jedan
neuron ne
može pratiti
(prisjetite se
prijedloga za
vježbu!).



Za polazne parametre $w^h_{11} = -13$, $w^h_{21} = 15$, $w^h_{12} = -5$,
 $w^h_{22} = -6$, $w^o_{11} = 6$, $w^o_{12} = 6$, $\theta^h_1 = 8$, $\theta^h_2 = 9$ i $\theta^o_1 = 4$ te
 $\eta = 1$, nakon 30000 iteracija (epoha) dobivamo

$$w^h_{11} = -16.1107$$

$$w^h_{21} = 17.1725$$

$$w^h_{12} = -6.5361$$

$$w^h_{22} = -6.8322$$

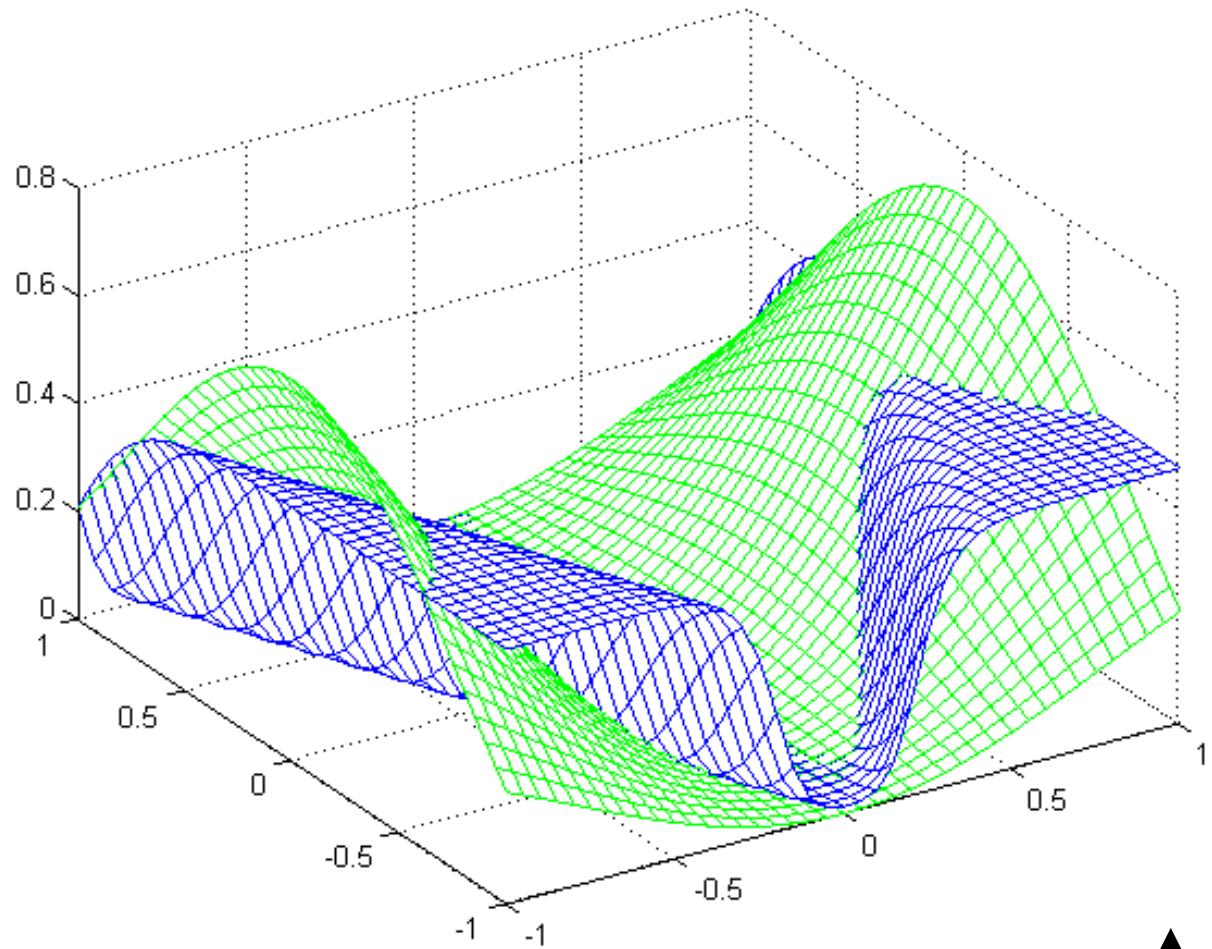
$$w^o_{11} = 6.2848$$

$$w^o_{12} = 6.2848$$

$$\theta^h_1 = 8.1618$$

$$\theta^h_2 = 8.9283$$

$$\theta^o_1 = 6.3961 .$$



Primjer: uvježbavanje mreže $2 \times 2 \times 1$

Mreža ima dva ulaza, dva neurona u skrivenom sloju i jedan izlaz, a aktivacijska funkcija je (opći) sigmoid.

Zadaća mreže je aproksimirati XOR (ekskluzivni ILI) logičku funkciju zadanu tablicom istine.

X1	X2	XOR
0	0	0
0	1	1
1	0	1
1	1	0



SimpleNetwork.m

Tablica neka ujedno bude i skup za uvježbavanje, dakle imamo četiri zadane točke.

Za polazne parametre $w^h_{11}=0.1$, $w^h_{21}=0.3$, $w^h_{12}=0.3$,
 $w^h_{22}= 0.4$, $w^o_{11}= 0.4$, $w^o_{12}= 0.6$, $\theta^h_1= 0.1$, $\theta^h_2=1$ i $\theta^o_1= - 0.3$
 te $\eta=10$, nakon 4000 iteracija (epoha) dobivamo

$$w^h_{11} = -7.8541$$

$$w^h_{21} = -5.7366$$

$$w^h_{12} = -8.443$$

$$w^h_{22} = -5.7857$$

$$w^o_{11} = -11.5318$$

$$w^o_{12} = 11.4628$$

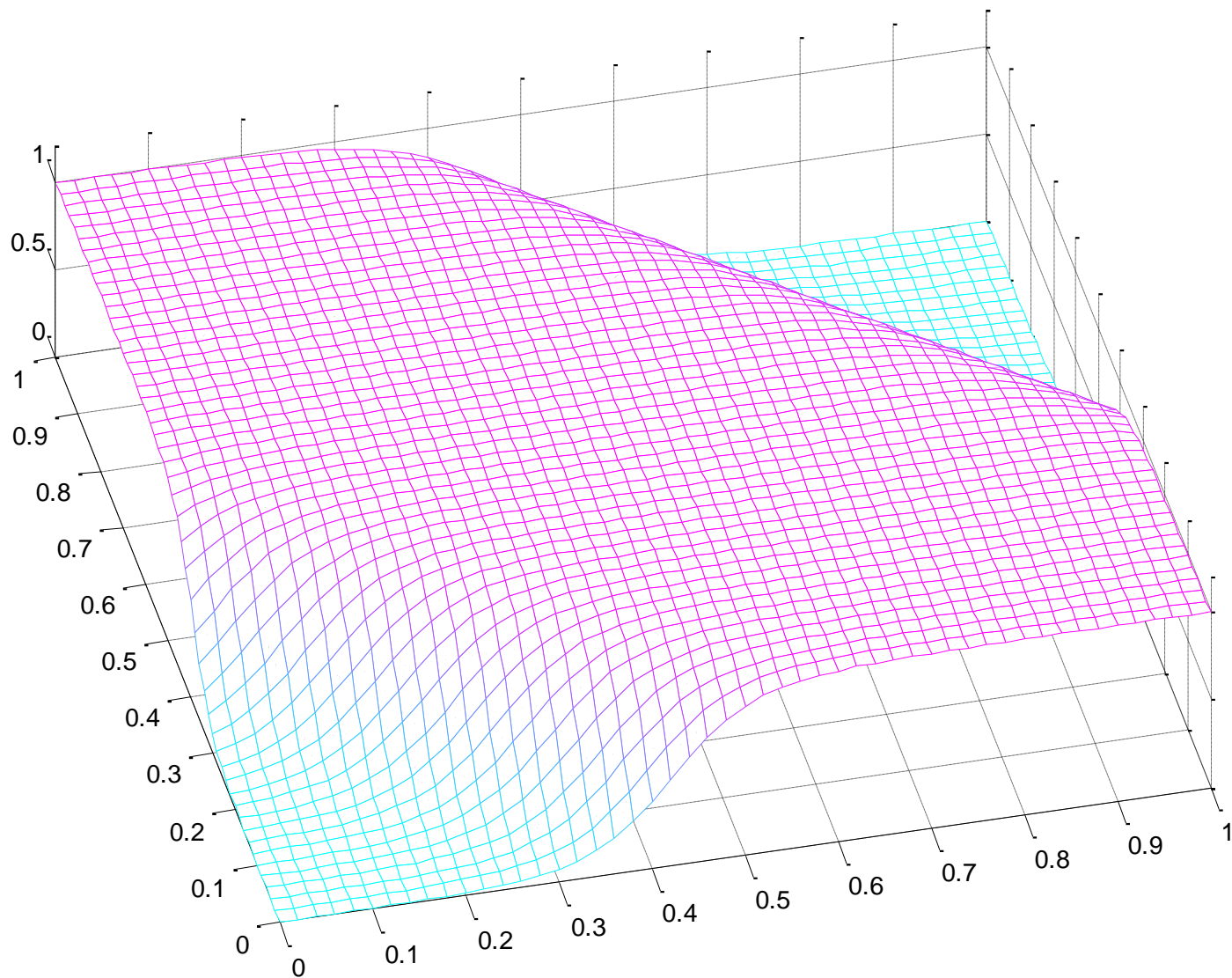
$$\theta^h_1 = -3.3684$$

$$\theta^h_2 = -8.5833$$

$$\theta^o_1 = 5.5385 \text{ .}$$

X1	X2	izlaz mreže
0	0	0.0053
0	1	0.9945
1	0	0.9943
1	1	0.0069

U cijelom prostoru aproksimacija je posve dobra...



Završne napomene i praktični savjeti

1. Početne vrijednosti parametara mreže

- parametri θ svi $=0$, a težine w = nasumični brojevi u opsegu u kojem se očekuju njihove konačne vrijednosti

Proračun w se temelji na:

- opsegu vrijednosti (i varijabilnosti) ulaznih varijabli
- broju ulaza u neurone pojedinog sloja
(= broj neurona u prethodnom sloju, ali samo u potpuno povezanoj mreži)
- aktivnom dijelu domene aktivacijske funkcije;
za sigmoid je to $\approx -5...5$
- izlaznom opsegu (kodomeni) aktivacijske funkcije; za sigmoid je to interval $(0,1)$ (asimptotske granice)

Primjer: troslojna mreža s n ulaza i l neurona u skrivenom sloju, očekivana vrijednost ulaza μ_x i očekivana vrijednost izlaza skrivenih neurona μ_h .

Očekivani zbroj ulaza (*lokalno polje*) u skriveni neuron je $n \cdot \mu_x \cdot w_{\text{mean}}^h$. Da bi neuron bio aktivan, taj zbroj treba biti < 5 . Zbog zalihosti odabiremo < 2 pa slijedi da srednja težina w^h treba biti $w_{\text{mean}}^h < 2/(n \cdot \mu_x)$, tj. težine w^h inicijaliziramo slučajnim brojevima $w^h \in [-w_{\text{mean}}^h, w_{\text{mean}}^h]$.

Očekivani zbroj ulaza (*lokalno polje*) u izlazni neuron jednak je $l \cdot w_{\text{mean}}^o \cdot \mu_h$. Mreža mora ostati aktivna i u ekstremnom slučaju (najveći mogući ulaz u izlazni sloj) pa mora biti $l \cdot w_{\text{mean}}^o \cdot \mu_h = < 5$ i kada je $\mu_h = \max = 1$. Slijedi $w_{\text{mean}}^o < 5/l$ i inicijalno postavljamo $w_{sj}^o \in [-w_{\text{mean}}^o, w_{\text{mean}}^o]$.



2. Stopa učenja

Teorijski uvjet konvergencije LMS algoritma prilikom uvježbavanja jednog neurona jest da je stopa učenja η

$$0 < \eta < 2/\lambda_{\max} \quad ,$$

gdje je λ_{\max} najveća svojstvena vrijednost korelacijske matrice ulaznih vektora $\mathbf{R}_{xx} = E(\mathbf{x} \cdot \mathbf{x}^T)$. U primjeni, \mathbf{R}_{xx} aproksimiramo izrazom $R_{xx} = 1/n \cdot \mathbf{X}^T \cdot \mathbf{X}$, gdje je $\mathbf{X} \in \mathbb{R}^{n \times p}$ matrica ulaznih vektora (vrijedi da je $\mathbf{R}_{xx} = \lim_{n \rightarrow \infty} R_{xx}$ kada $n \rightarrow \infty$). Ako su svi ulazi međusobno statistički nezavisni i imaju očekivanje $E(x_j) = 0$, proizlazi da je $\lambda_{\max} = \max\{\sigma_j^2\}$, odnosno

$$\lambda_{\max} = \max_j \left\{ \frac{1}{n} \sum_{i=1}^n (x_j)_i^2 \right\} \quad ,$$

gdje je j indeks ulaza (ulaznog vektora), a i indeks podataka.

Dvije metode prilagodbe stope učenja najčešće u praksi jesu:

- motriti trenutačnu pogrešku ($y_s - y_{ds}$) i ovisno o njoj mijenjati stopu η tijekom uvježbavanja mreže
 - ako pogreška ($y_s - y_{ds}$) mijenja predznak (oscilira), stopa učenja je prevelika i treba ju smanjiti da bi se postigla brža konvergencija; vrijedi i obrat
 - bilo bi dobro kad bismo u svakom koraku mogli izračunati optimalnu stopu η , tj. primijeniti *steepest descent* algoritam, ali...
- “automatizirati” prilagodbu stope učenja uvođenjem tromosti u izračunavanje ispravaka parametara
 - poopćeno *delta pravilo*

$$\Delta w_{ji}(k) = \gamma \cdot \Delta w_{ji}(k-1) - \eta \cdot \delta_j(k) \cdot z_i(k) \quad ,$$

gdje je γ *tromost* (konstanta inercije) postupka

- unosi povratnu vezu, tj. memoriju u postupak

Poopćeno *delta pravilo* je zapravo jednađžba diferencija (*difference equation*) prvog reda po funkciji $\Delta w_{ji}(k)$.

Rješenje te jednađžbe je

$$\Delta w_{ji}(k) = \eta \sum_{t=0}^k \gamma^{k-t} \delta_j(t) z_i(t) \quad ,$$

što je niz (signal; *time series*) od $k+1$ pribrojnika s eksponencijalno promjenjivim težinama γ^{k-t} . Ako se prisjetimo da je $\delta_j(t) \cdot z_i(t) = \mathbf{E} \mathbf{W}_{ji}(t) = \partial \mathbf{e}(t, \mathbf{w}) / \partial w_{ji}$, slijedi

$$\Delta w_{ji}(k) = \eta \sum_{t=0}^k \gamma^{k-t} \frac{\partial \mathbf{e}(t, \mathbf{w})}{\partial w_{ji}} \quad ,$$

a iz te formule možemo iščitati učinak uvođenja tromosti, kao i uvjete na parametar γ .

- a) Kada su parcijalne derivacije $\partial e(t, \mathbf{w}) / \partial w_{ji}$ izlazne pogreške po parametru w_{ji} “mahom” istog predznaka, dakle kada je promjena u dimenziji nekog parametra “dugo” približno istog smjera, ispravak Δw_{ji} će rasti, odnosno parametar w_{ji} će biti značajnije promijenjen. Vrijedi i obratno pa je jasno da uvođenje tromosti ima stabilizacijski učinak na proces uvježbavanja mreže jer će možebitne oscilacije predznaka ispravaka u dimenziji pojedinog parametra dovesti do smanjenja ispravka.
- b) Ispravci $\Delta w_{ji}(k)$, matematički gledano, jesu red, a iz uvjeta konvergencije redova slijedi da tromost γ mora biti u granicama $0 \leq |\gamma| < 1$. Tromost može biti i pozitivna i negativna, iako negativna baš i nema pretjeranog smisla...

3. Kriterij zaustavljanja

Gradijentna metoda sigurno konvergira ka točki u kojoj gradijent isčezava (jednak je nuli), ali pronalaženje globalnog optimuma nije zajamčeno pa tako nema ni jedinstvenog i matematički neprikosnovenog kriterija za zaustavljanje.

Neki češće primjenjivani kriteriji jesu:

- norma gradijenta manja od zadanog minimuma (upitno prilikom koračnog uvježbavanja)
- ukupna pogreška mreže (za cijeli skup) manja od zadanog minimuma
- promjena ukupne pogreške mreže, u uzastopnim koracima, manja od zadanog minimuma.

4. Ulazni podatci nejednake važnosti (pouzdanosti)

- mrežu uvježbavamo pridjeljujući parovima za uvježbavanje različite težine t_i

minimizirati
$$\frac{1}{2} \sum_{s=1}^m t_i (y_s - y_{ds})^2 \quad ; \quad i=1, 2, \dots, p$$

- za koračno uvježbavanje to znači da cijeli gradijent množimo s t_i u svakoj točki, tj. za svaki par vektora $(\mathbf{x}_d, \mathbf{y}_d)$

5. Promjena redoslijeda ulaznih podataka

- prilikom koračnog uvježbavanja dobro je računati sa slučajno odabranim parovima, tj. skup za uvježbavanje svaki puta (u svakoj iteraciji; epohi) treba “predočiti” mreži drugačijim redoslijedom

6. Normalizacija ulaza

Ako to namjena mreže dozvoljava, poželjno je:

- a) pomaknuti (translatirati) ulazne varijable tako da im srednja vrijednost bude približno nula
- b) ukloniti svaku statističku vezu među ulazima
 - taj se postupak naziva *dekorelacija*, a može se postići primjenom statističke tehnike koja se naziva *analiza glavnih (najvažnijih) sastavnica* (*principal component analysis, PCA*)
- c) skalirati ulazne vrijednosti tako da im kovarijance budu približno jednake

Normalizacija ulaza ujednačava brzinu prilagodbe (učenja) veza (parametra) mreže.