

```
def feed_forward(self, inputs):
```

```
    # Convert the list of inputs into a 2D array
```

```
    I = np.array(inputs, ndmin=2).T
```

```
    # x_hidden = W_1h * I
```

```
    X_hidden = np.dot(self.W_1h, I)
```

```
    # O_hidden = sigmoid(X_hidden): the output of the hidden layer is its input squashed using the sigmoid function
```

```
    O_hidden = self.activation_function(X_hidden)
```

```
    # x_final = W_2h * O_hidden: the signals into the final output layer
```

```
    X_final = np.dot(self.W_2h, O_hidden)
```

```
    # O_final = sigmoid(X_final)
```

```
    O_final = self.activation_function(X_final)
```

```
    return O_final, O_hidden, I
```

```
def train(self, inputs, targets):
```

```
    # Feed forward
```

```
    O_final, O_hidden, I = self.feed_forward(inputs)
```

```
    # Backpropagation
```

```
    # Convert the list of target values into a 2D array
```

```
    T = np.array(targets, ndmin=2).T
```

```
    # Compute the error: E_output = T - O_final
```

```
    E_output = T - O_final
```

```
    # Compute the backpropagated errors for the hidden layer nodes
```

```
    # E_hidden = transposed(W_2h) * E_output
```

```
    E_hidden = np.dot(self.W_2h.T, E_output)
```

```
    # update the weights between the hidden and output layers using the formula obtained in the first part of the workshop
```

```
    # delta_w_2h = lr * E_h * sigmoid(O_h) * (1 - sigmoid(O_h)) * transposed(O_j)
```

```
    self.W_2h += self.lr * np.dot((E_output * O_final * (1 - O_final)), np.transpose(O_hidden))
```

```
    # update the weights between the hidden and the input layer
```

```
    self.W_1h += self.lr * np.dot((E_hidden * O_hidden * (1 - O_hidden)), np.transpose(I))
```

2NClub #1

2N Club

My First Neural Network



Part I

The Math behind a Neural Network

1.1. Modele liniare

Andrei Nicolicioiu

Machine Learning
Researcher

@ *Bitdefender*



Clasificare - Detecție - Segmentare

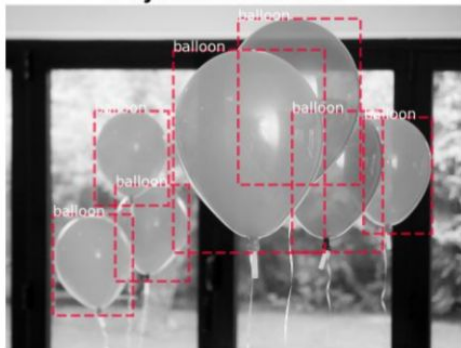
Classification



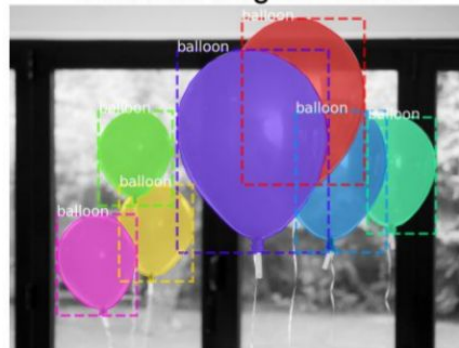
Semantic Segmentation



Object Detection



Instance Segmentation



Generare Immagini



2014



2015



2016



2017

(Brundage et al, 2018)

Adaptare



- CycleGAN



Captioning



Generated:

0.650. a group of people are sitting in a line with a tiger

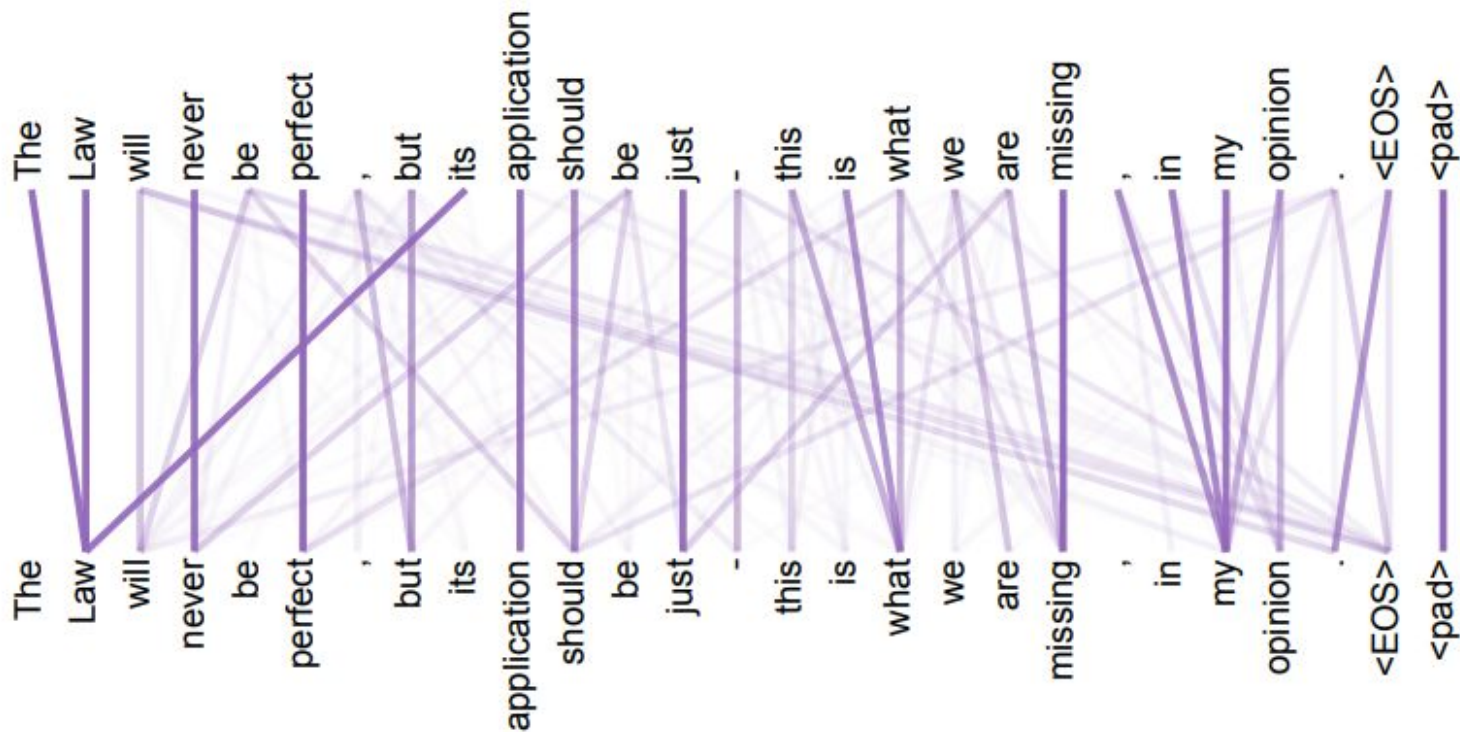
0.599. a man is sitting in a chair with a tiger

0.588. a man and a woman are eating a tiger in a bowl

0.547. a man is talking about a tiger

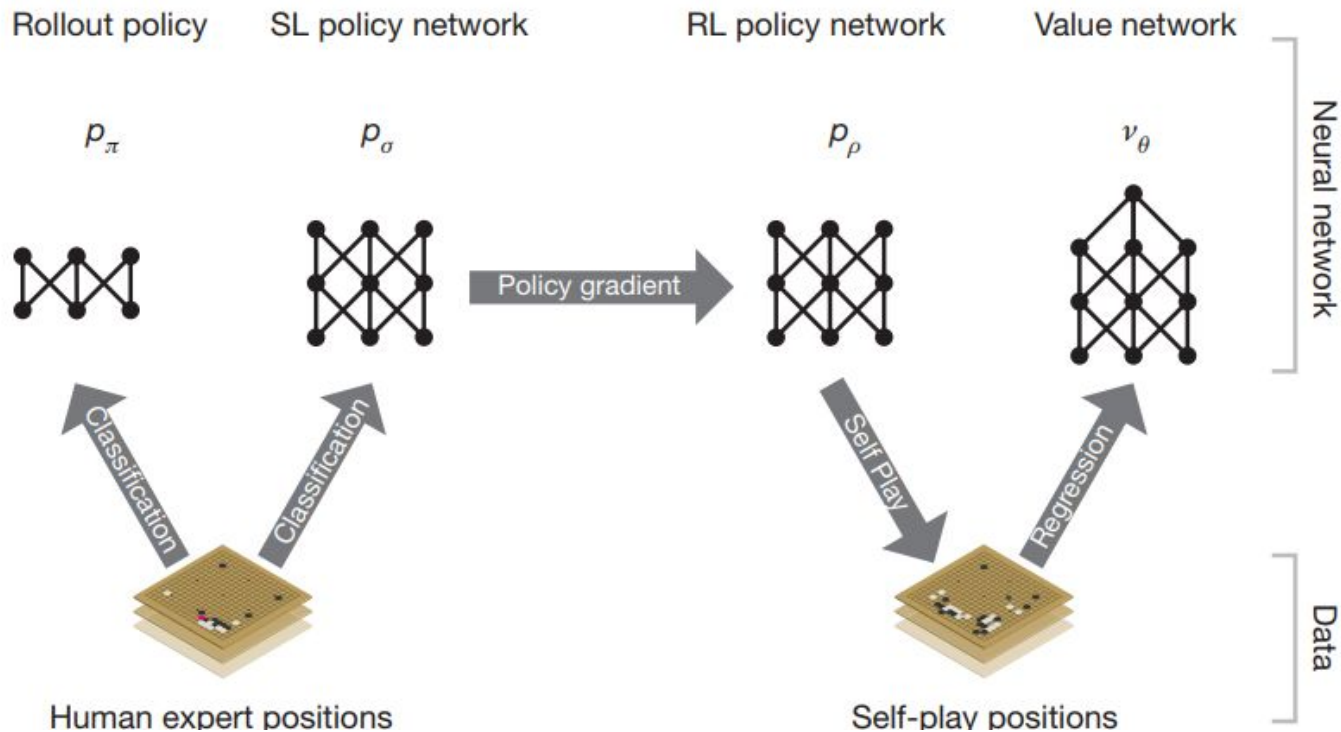
0.180. a man and a woman are sitting in a table

Analiză text - Traduceri



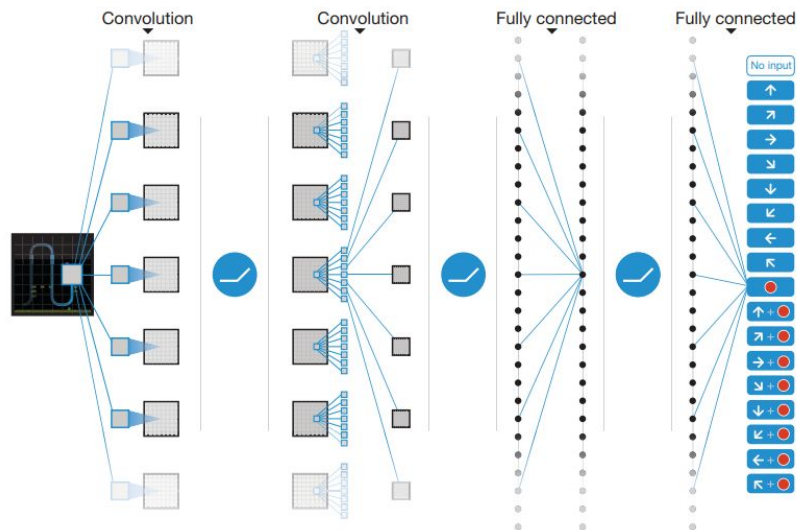
- Vaswani et al - Attention is all you need [2017]

Alpha Go



- Alpha Go - 2016

Reinforcement Learning



V. Mnih et. al - DQN - 2015



De ce Machine Learning?

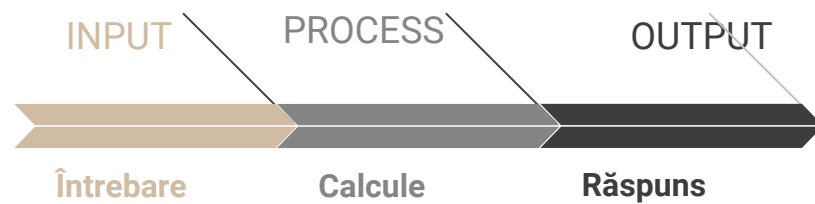
- Problemă clasică: Recunoașterea cifrelor scrise de mână

De ce Machine Learning?

- Vrem să recunoaștem cifre scrise de mână. Cum procedăm?
- Găsim reguli pentru fiecare cifră. Exemplu:
 - pentru cifra 8 găsim 2 cercuri de aceeași dimensiune așezate unul peste altul
 - cunoaștem ecuația cercului: $x^2 + y^2 = R$
 - găsim toate cercurile, le păstrăm doar pe cele așezate unul peste altul
- Probleme:
 - foarte multe reguli greu de identificat
 - varietate foarte mare în date
 - nu pot găsi reguli care să acopere toate posibilitățile

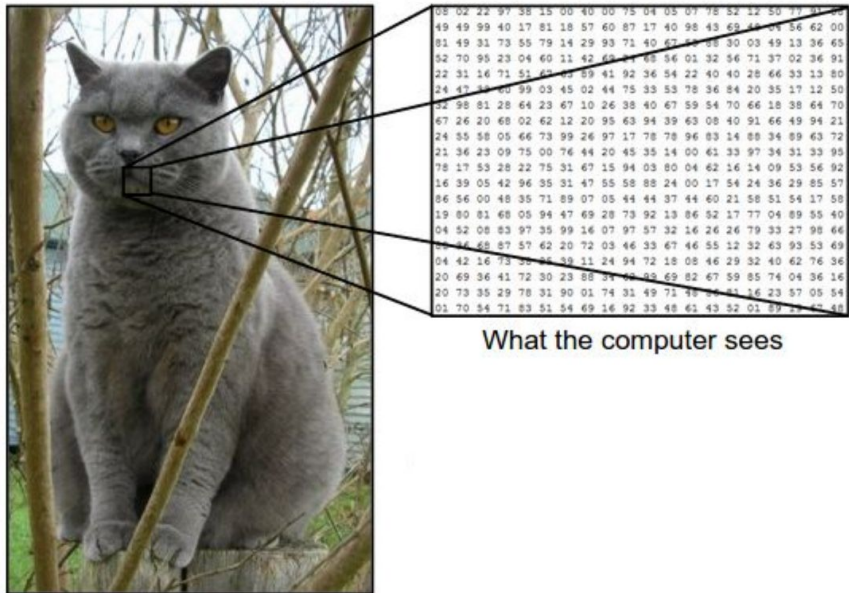
Machine Learning

- Scop:
 - Crearea unui program capabil să învețe automat
- Mijloace:
 - **Dataset:** perechi **input - target**
 - **Procesare** input
 - La ieșire vom avea un **rezultat**
- Urmărim 3 lucruri:
 - Ce primește programul la input: **reprezentare**
 - Care sunt operațiile de procesare: **model**
 - Cum învață modelul : **optimizare**



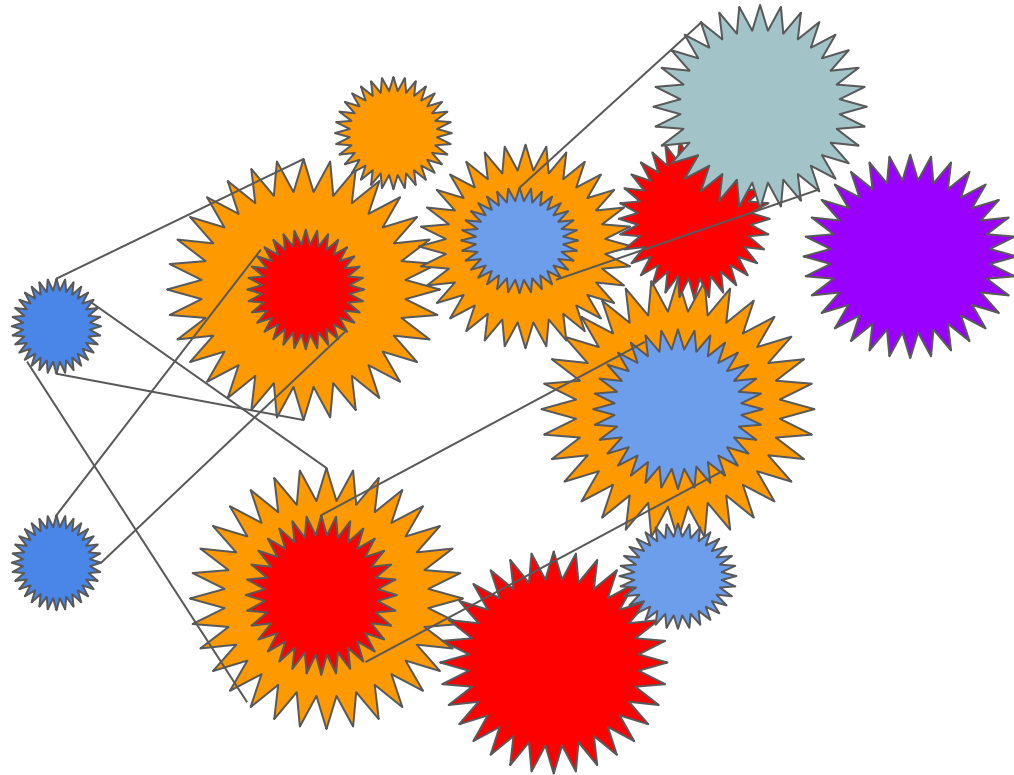
Human vs Computer

- Procesare de date
 - => **reprezentare**
- Din date neprelucrate (pixeli) extragem niște reprezentări, caracteristici - **features**
- Putem folosi direct datele neprelucrate și vom **învăța** automat reprezentări



What the computer sees

Model - Analogie



- Un model are nevoie de:
 - **Arhitectura** (cum sunt conectate roțile)
 - **Parametrii**: (dimensiunea roților, cu cât amplifică)
-
- Pentru o arhitectura aleasa vrem un mode capabil să **învețe** singur parametrii

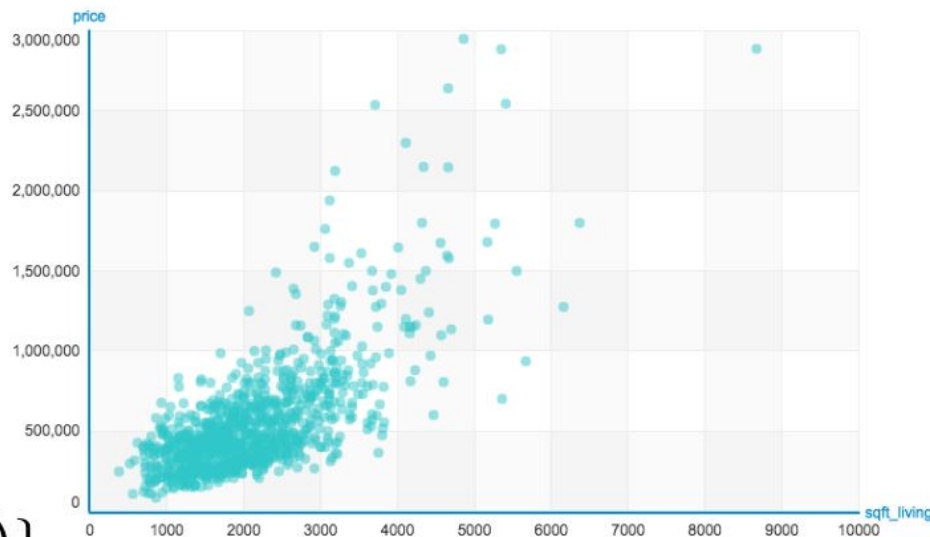
Regresie

- Problemă:
 - Prezicerea prețului de vânzare a unei locuințe în funcție de suprafață

- Set de date despre locuințe:

$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_M, t_M)\}$$

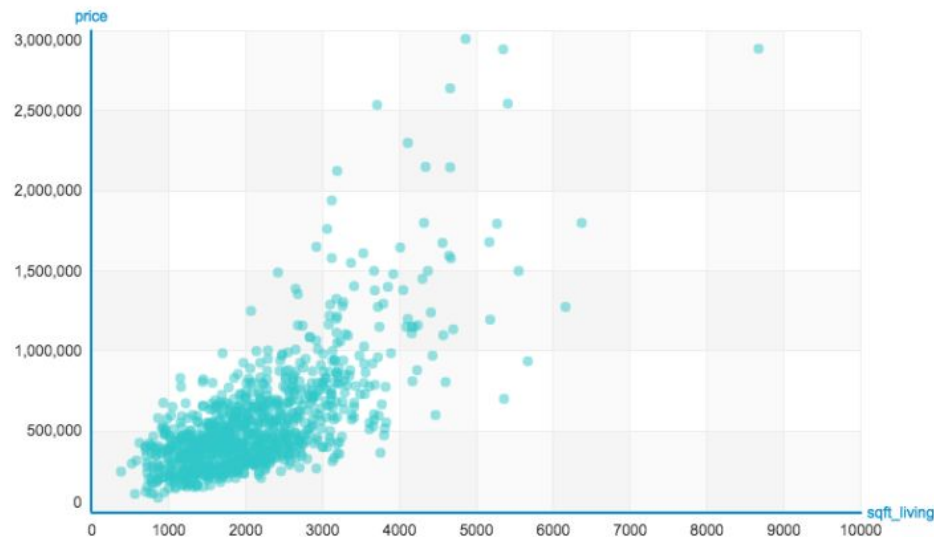
- **Reprezentare:** suprafața x
 - **Target:** prețul adevărat al apartamentului
- Scop:
 - Folosirea unui **model** adecvat pentru a putea prezice prețul
 - **Optimizarea** (învățarea) modelului cu ajutorul setului de date



Regresie

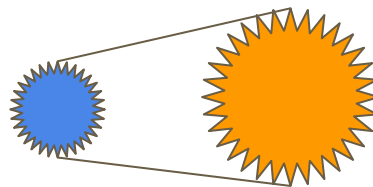
Model Liniar

- Reprezentare: $x \in \mathbb{R}$
- Cel mai simplu model:
$$y = w \cdot x + b$$
- Analogie cu o roată zimțată
- Găsește parametrul (raza roții)
a.î. prețul prezis pentru orice
casă sa fie cât mai apropiat de
cel real



x

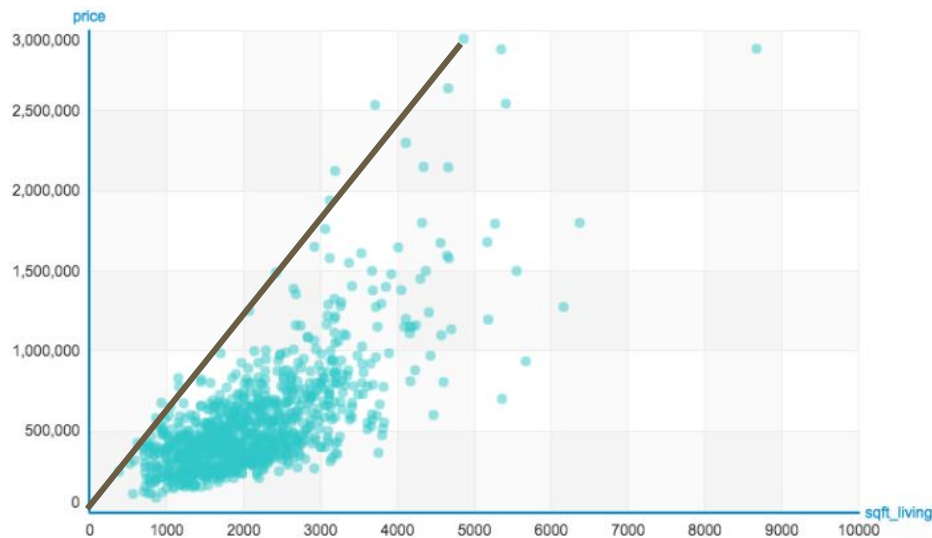
$$y = w \cdot x$$



Regresie

Model Liniar

- Reprezentare: $x \in \mathbb{R}$
- Cel mai simplu model:
$$y = w \cdot x + b$$
- Analogie cu o roată zimțată
- Găsește parametrul (raza roții) a.î. prețul prezis pentru orice casă sa fie cât mai apropiat de cel real



x

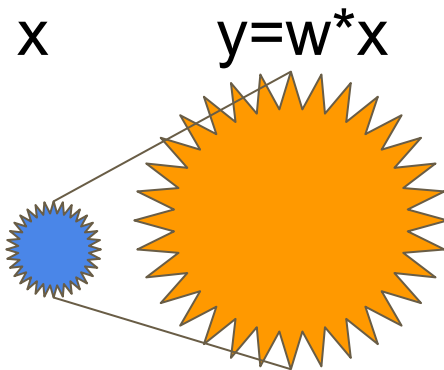
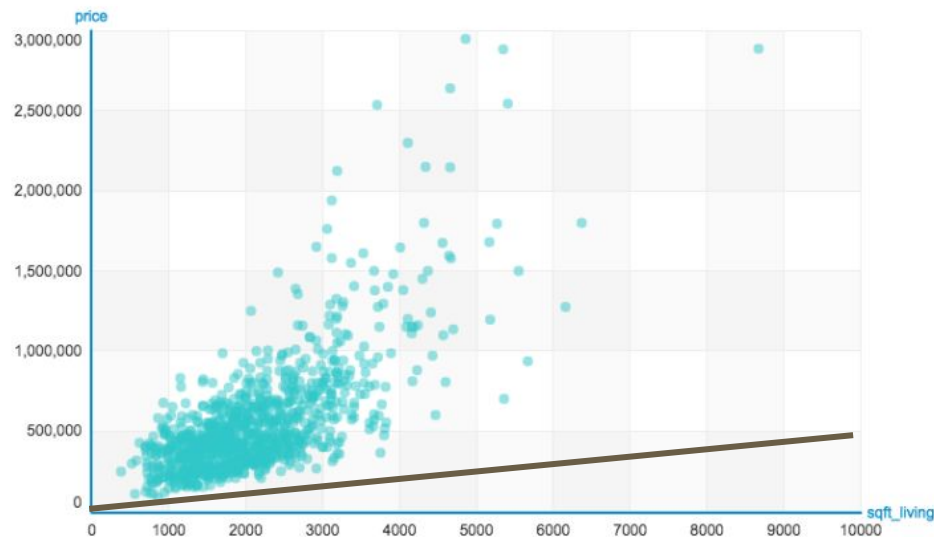
$$y = w * x$$



Regresie

Model Liniar

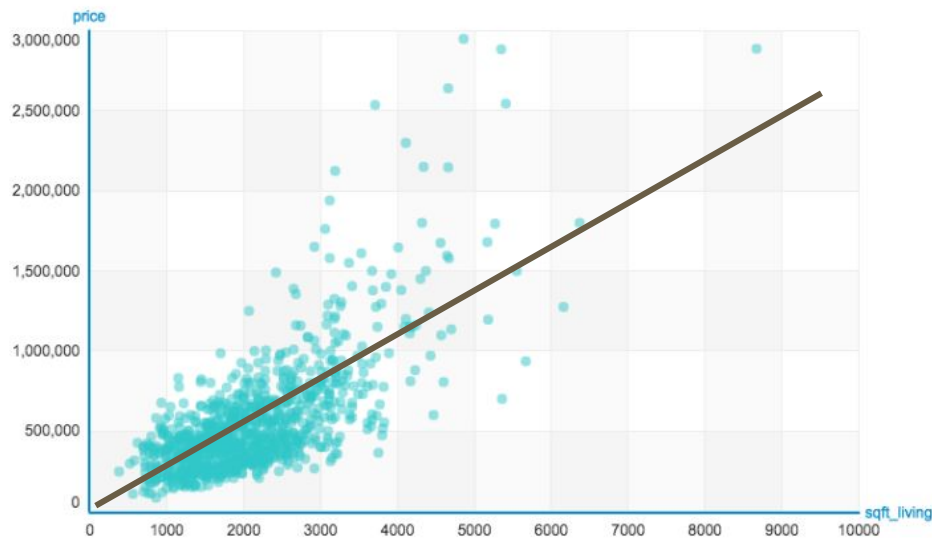
- Reprezentare: $x \in \mathbb{R}$
- Cel mai simplu model:
$$y = w \cdot x + b$$
- Analogie cu o roată zimțată
- Găsește parametrul (raza roții) a.î. prețul prezis pentru orice casă sa fie cât mai apropiat de cel real



Regresie

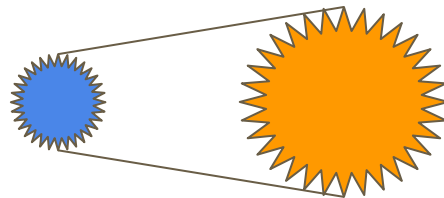
Model Liniar

- Reprezentare: $x \in \mathbb{R}$
- Cel mai simplu model:
$$y = w \cdot x + b$$
- Analogie cu o roată zimțată
- Găsește parametrul (raza roții) a.î. prețul prezis pentru orice casă sa fie cât mai apropiat de cel real



x

$$y = w * x$$



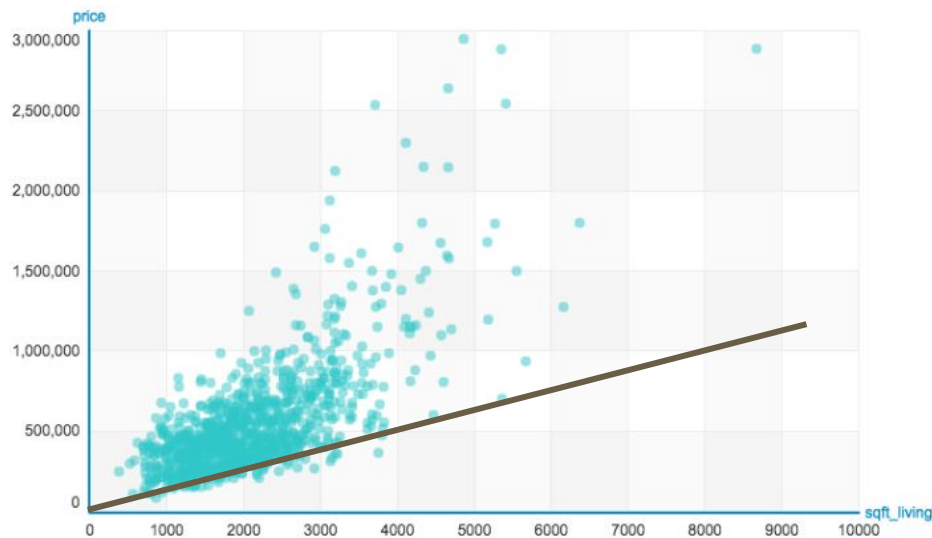
Regresie

Optimizare

- Ponderi aleatoare
- Funcție de eroare(cost), cât de prost prezice modelul nostru baza de date:

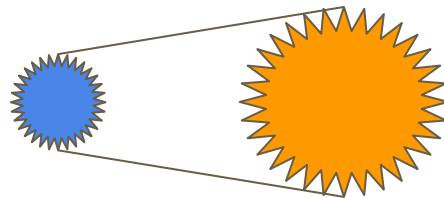
$$E = \sum_i |t_i - y_i|$$

- Ajustăm parametrii
- Păstrăm parametrii care minimizează eroarea



x

$$y=w*x$$

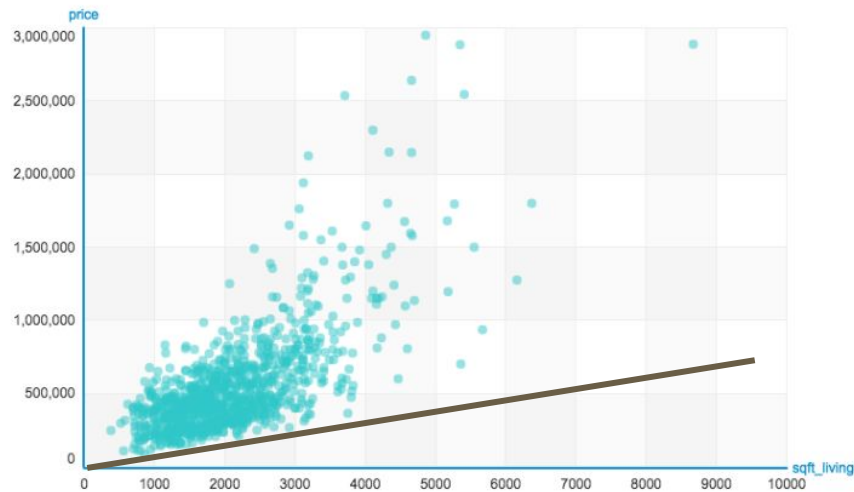


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente

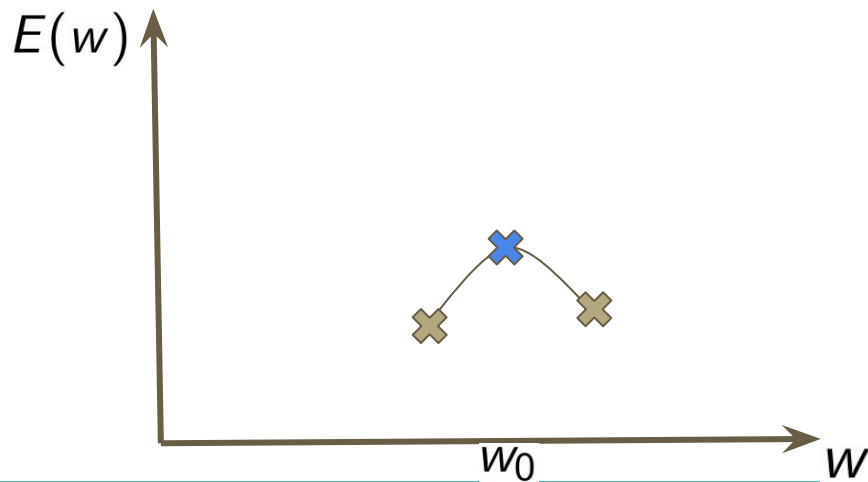
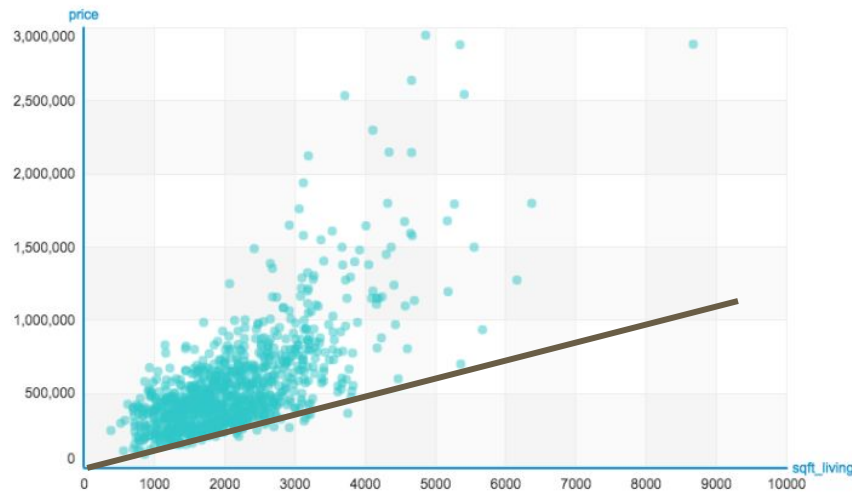


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente

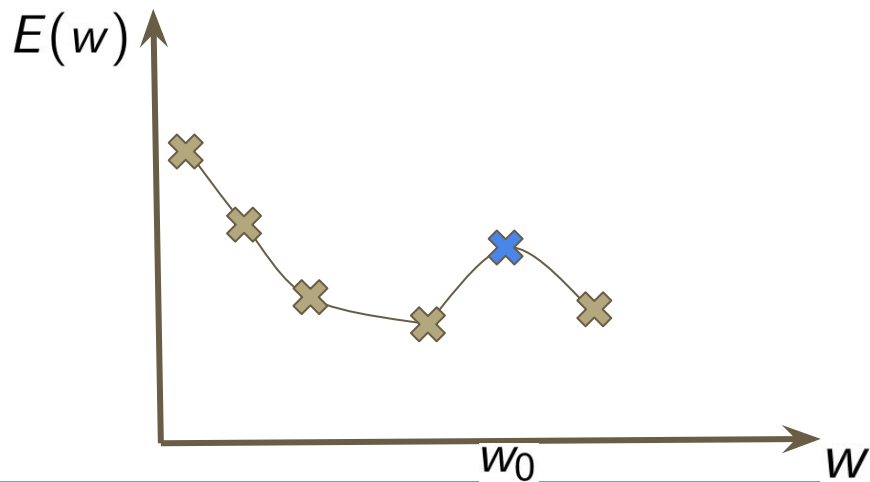
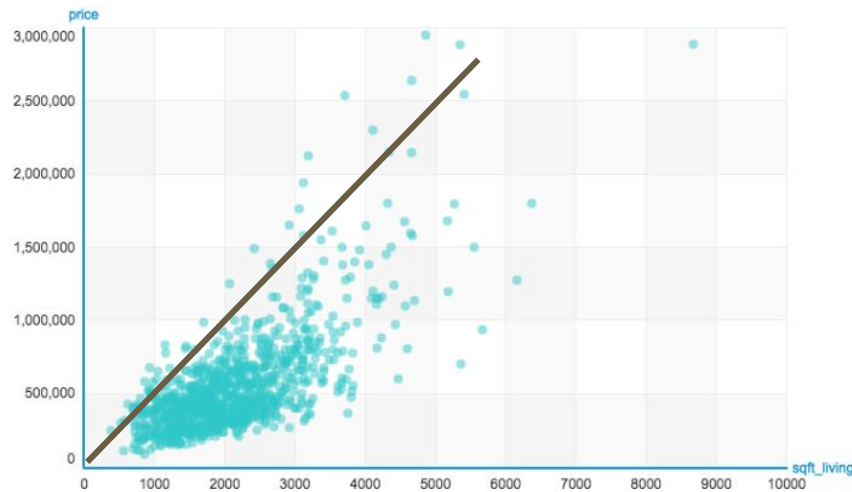


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente

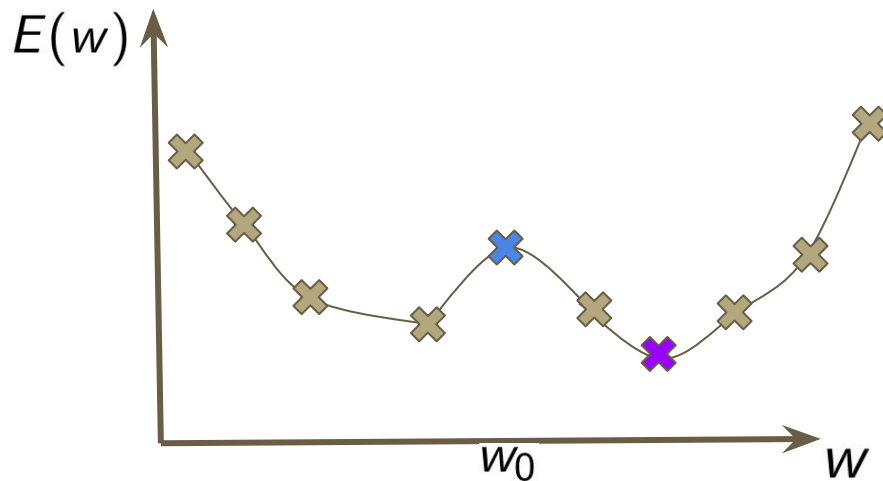
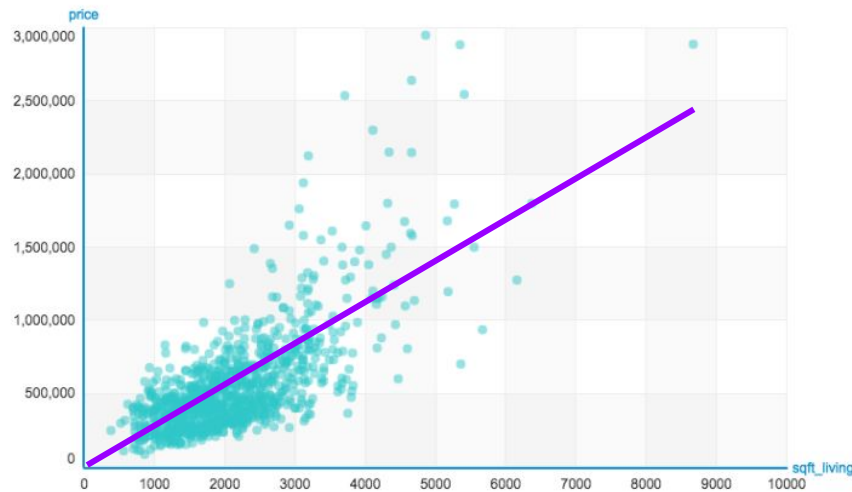


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente

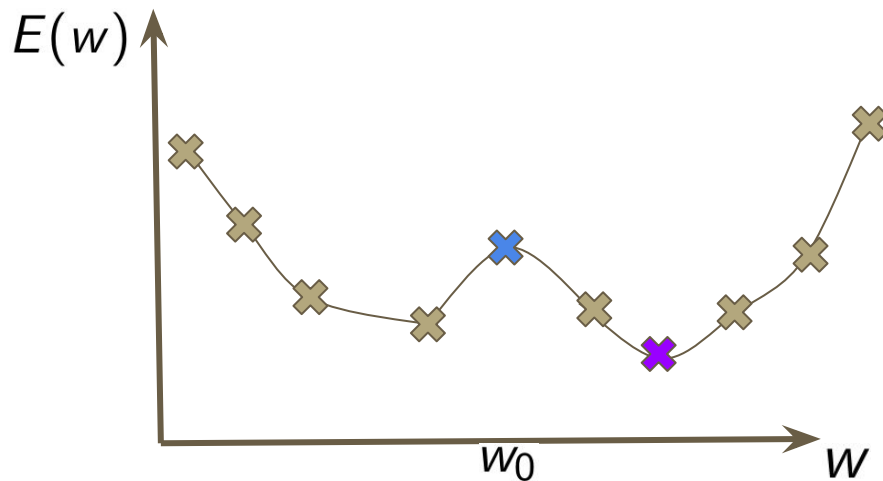
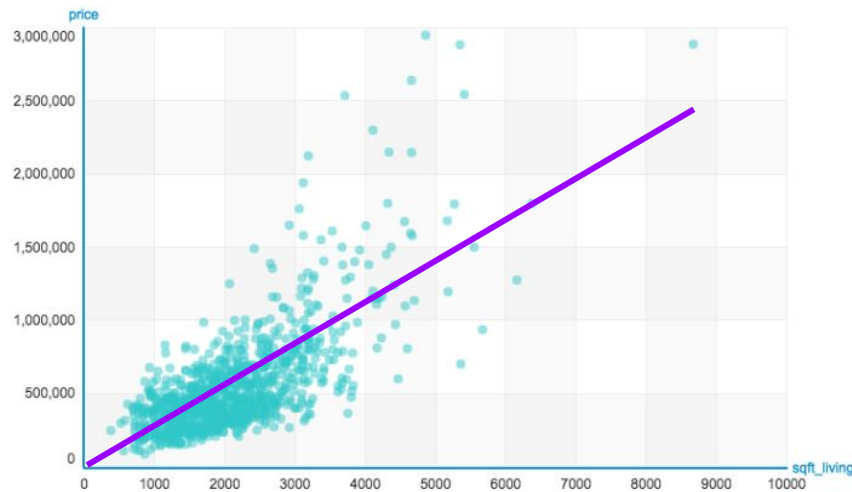


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente
- Ineficient!
 - Calcule pentru fiecare variantă

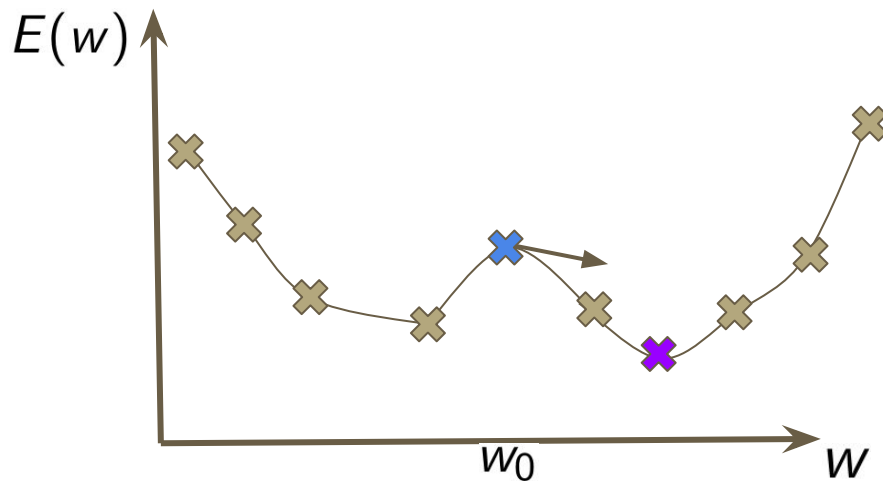
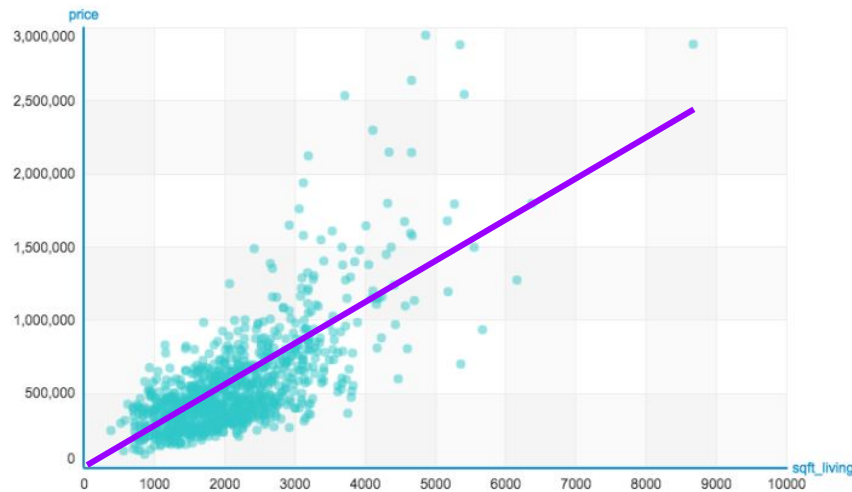


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente
- Ineficient!
 - Calcule pentru fiecare variantă
- Soluție: urmăm panta funcției

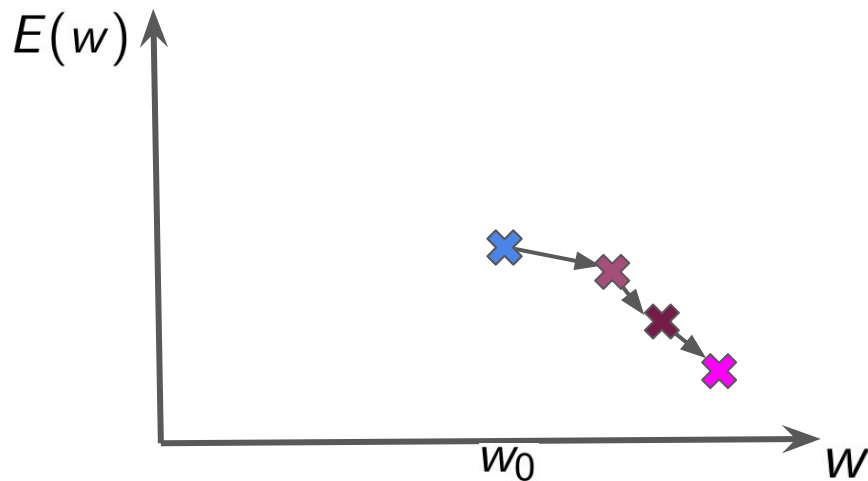
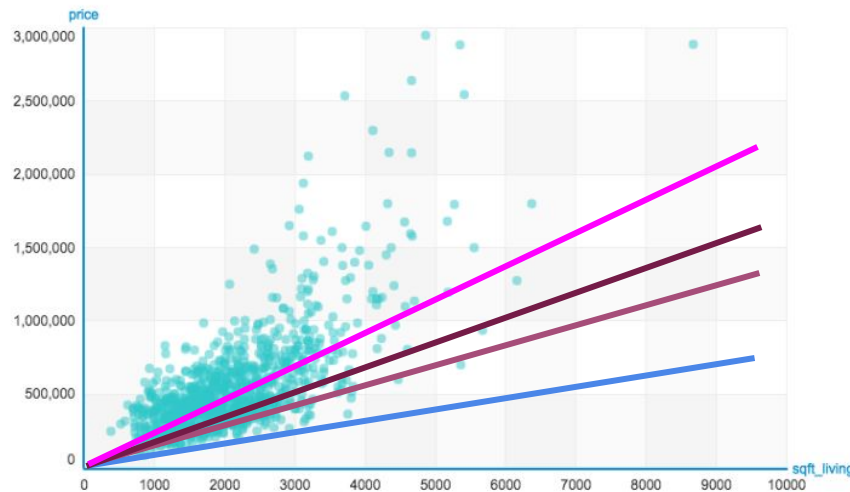


Regresie

Optimizare

Varianta 1:

- Încercăm multe valori pentru parametrii w_0 în jurul valorii curente
- Ineficient!
 - Calcule pentru fiecare variantă
- Soluție: urmăm panta funcției
 - În punctul curent w_0 calculăm panta funcției de eroare
 - Găsim parametri noi urmând această direcție
 - Necesită câțiva pași



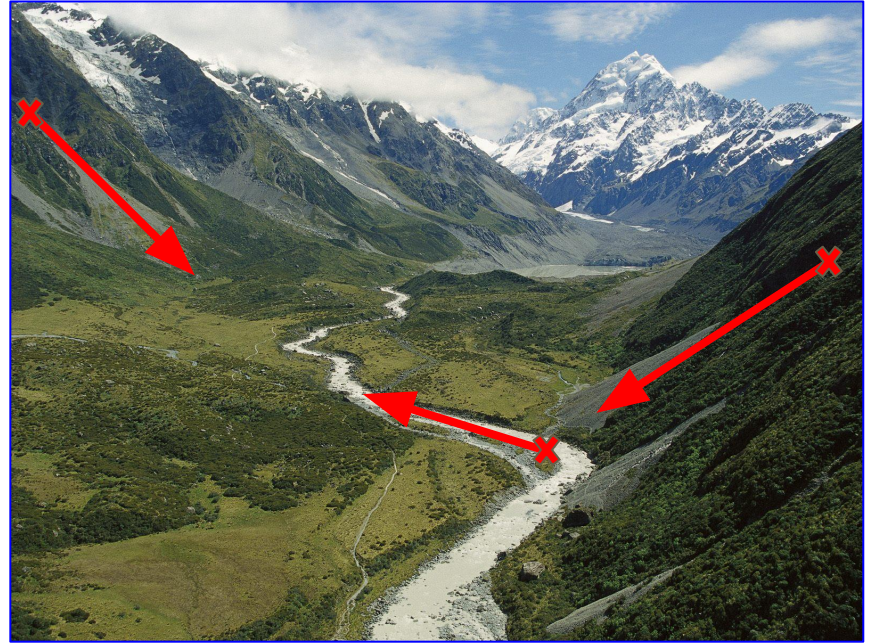
Gradient descent

Scop:

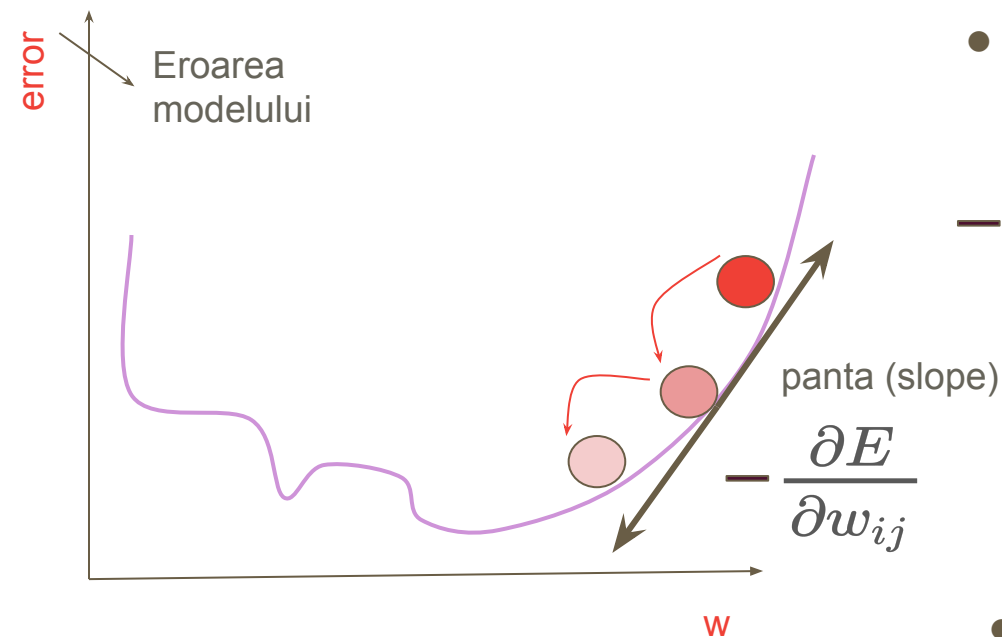
- găsirea minimului unei funcții.

Soluție:

- urmarea “pantei”
- Panta = direcția în care funcția scade cel mai brusc



Gradient descent

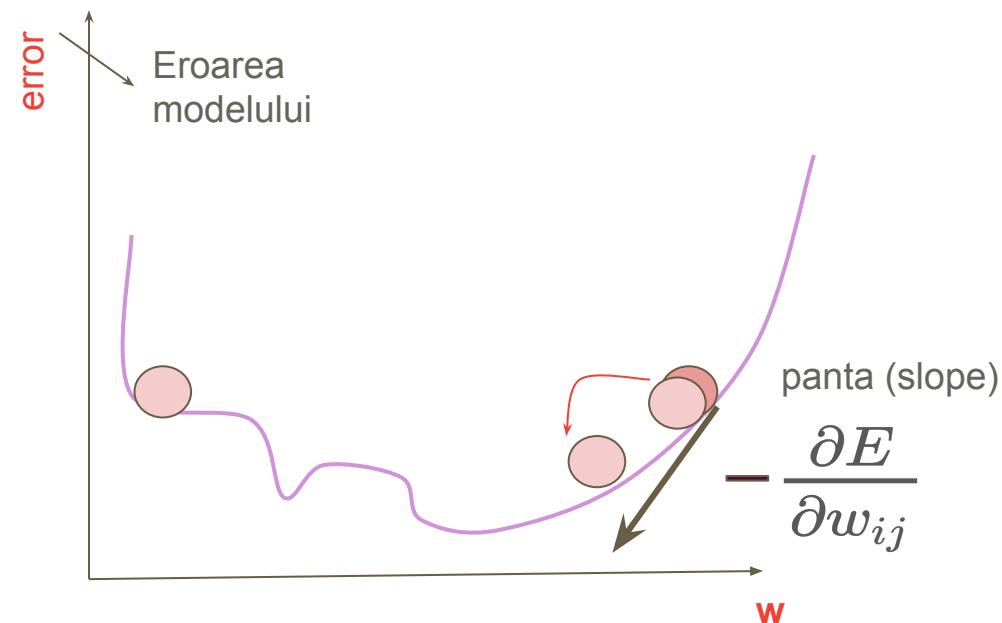


- Căutăm direcția în care E scade

$$E(w + \Delta w) - E(w)$$

- Panta = Derivată = Gradient

Gradient descent



$$-\lim_{\Delta w \rightarrow 0} \frac{E(w + \Delta w) - E(w)}{\Delta w} = \frac{dE}{dw}$$

$$w = w - \alpha \frac{dE}{dw}$$

Learning rate

- Alpha mic -> învățare prea lentă
- Alpha mare -> risc să sar peste puncte de minim

Regresie liniară

- Avem dat un dataset

$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_M, t_M)\}$$

- Stabilim o funcție de cost / eroare E
- Căutăm parametrii unei funcții f care să aibă eroare cât mai mică

$$y = f(x) = w \cdot x + b$$

Regresie liniară

1. Inițializăm w random

2. Calculăm ieșirile:

$$y_i = w \cdot x_1 + b$$

3. Calculăm eroarea pe întreg datasetul:

$$E = \frac{1}{2} \sum_i^M (t_i - y_i)^2$$

4. Actualizăm parametrii:

$$w = w - \alpha \frac{dE}{dw}$$

5. Reluăm pașii 2-3 cât timp eroarea scade

Clasificare

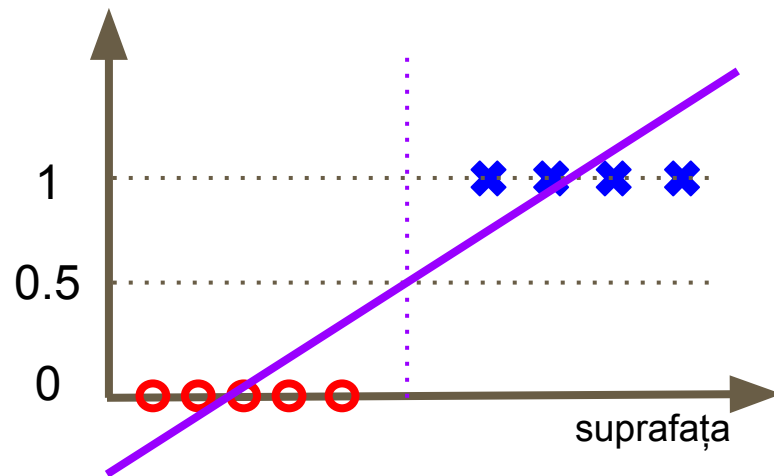
- Avem de ales între două variante.
- De exemplu: știm doar suprafața unui apartament. Cum decidem dacă are 2 sau 3 camere?

Clasificare

- Avem de ales între două variante.
- De exemplu: știm doar suprafața unui apartament. Cum decidem dacă are 2 sau 3 camere?

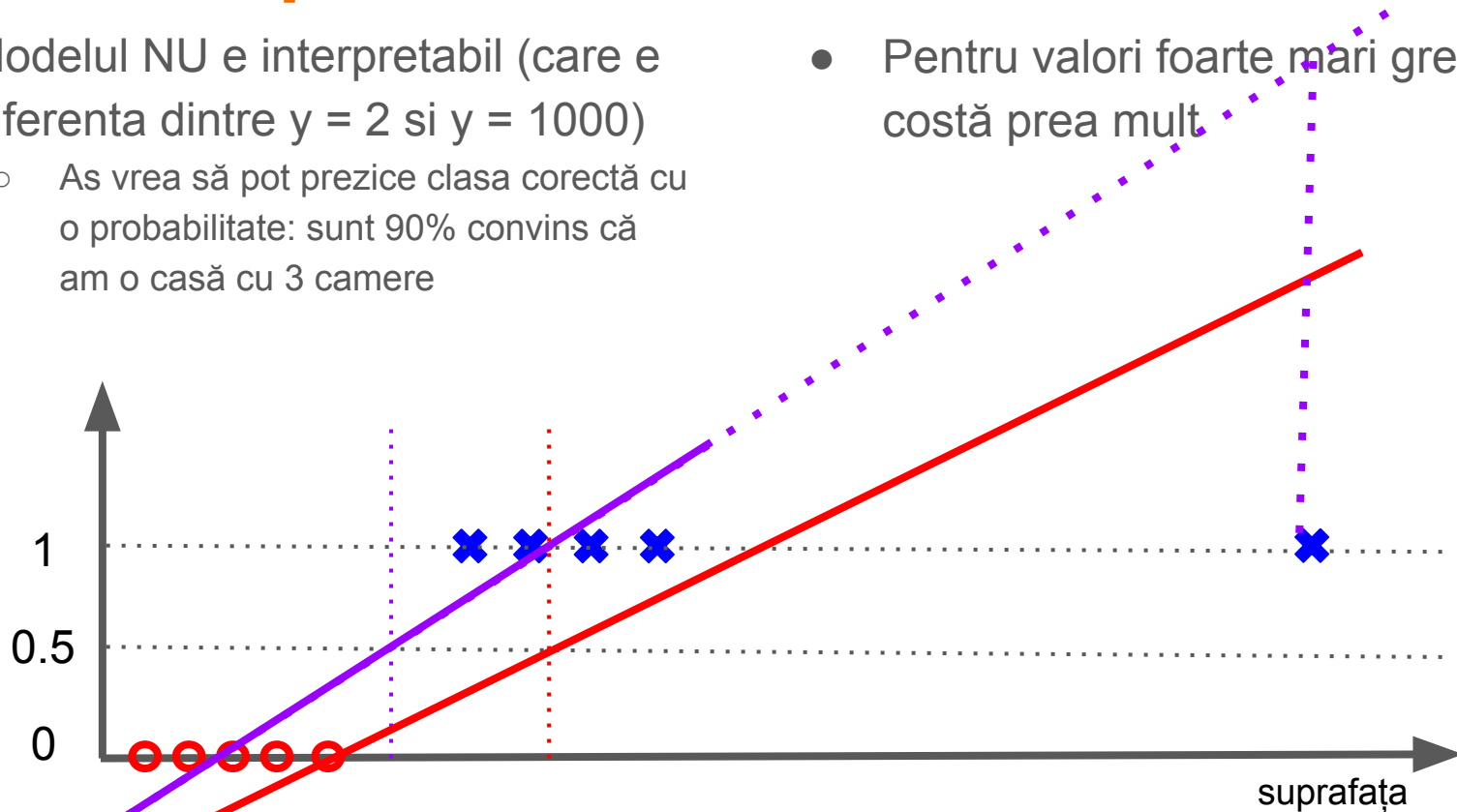
Varianta 1:

1. Colectăm date:
 - a. $y = 0 \Rightarrow$ două camere
 - b. $y = 1 \Rightarrow$ trei camere
2. Aplicăm regresie liniară
 - a. Toate valorile $y > 0.5$ vor fi clasificate ca exemplu pozitiv (3 camere).



Clasificare - probleme

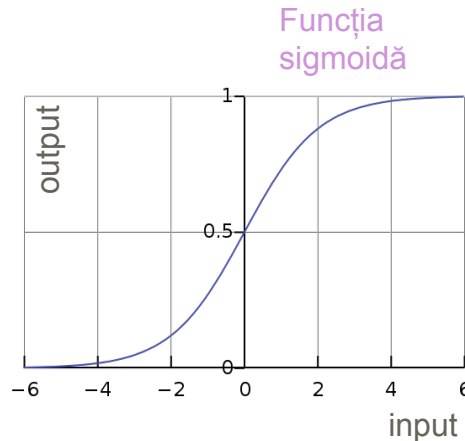
- Modelul NU e interpretabil (care e diferența dintre $y = 2$ și $y = 1000$)
 - As vrea să pot prezice clasa corectă cu o probabilitate: sunt 90% convins că am o casă cu 3 camere
- Pentru valori foarte mari greșelile costă prea mult.



Clasificare - Soluție

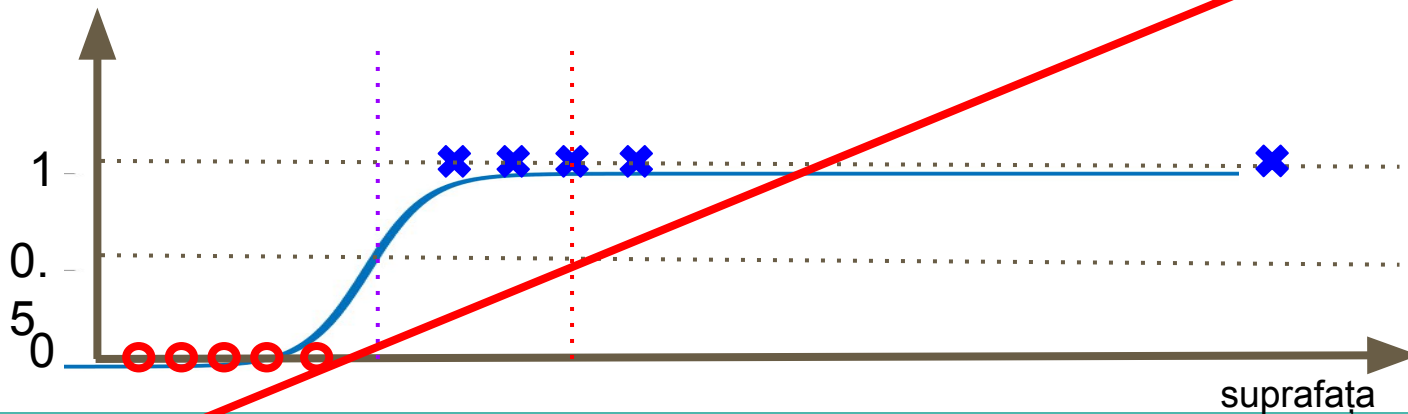
- Vom folosi funcția sigmoid
- Modelul nostru devine

$$a = \text{sigmoid}(w * x + b)$$



$$y = \frac{1}{1 + e^{-x}}$$

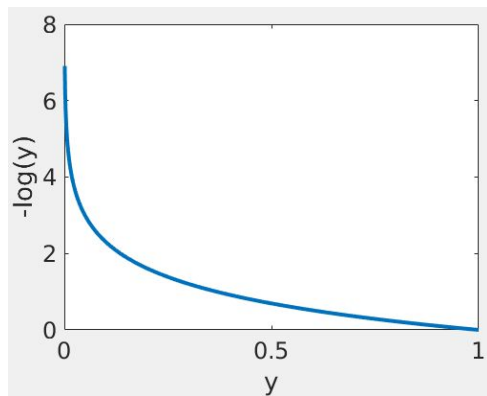
- Toate valorile mari se comportă la fel: sunt apropiate de 1
- Modelul nostru ia valori între 0 (2 camere) și 1 (3 camere)



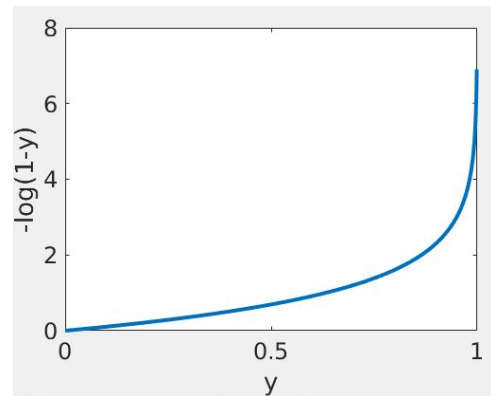
Clasificare - Logistic Regression

- Aceeasi pași ca la regresie liniară:
 - Inițializez parametrii
 - Calculăm ieșirile
 - Calculăm eroarea
 - Updatăm parametrii
 - Repetăm
- Ce funcție de eroare folosim?
 - În cazul regresiei liniare am folosit loss-ul L2:
- Pentru clasificare folosim cross-entropy:

$$E = -\frac{1}{M} \sum_i^M (t_i \log a_i + (1 - t_i) \log (1 - a_i))$$

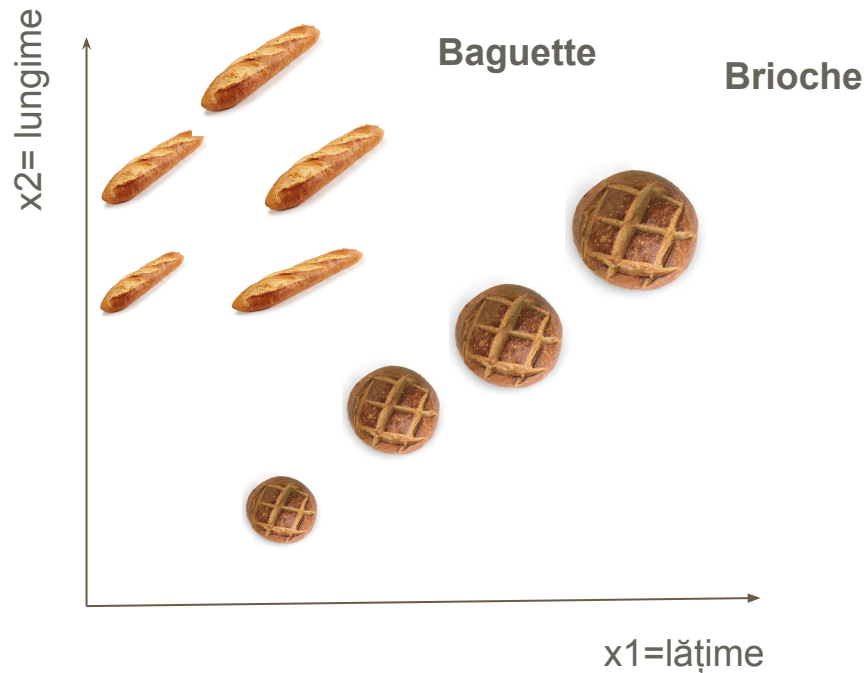


$t_i = 1$



$t_i = 0$

Clasificare - Exemplu

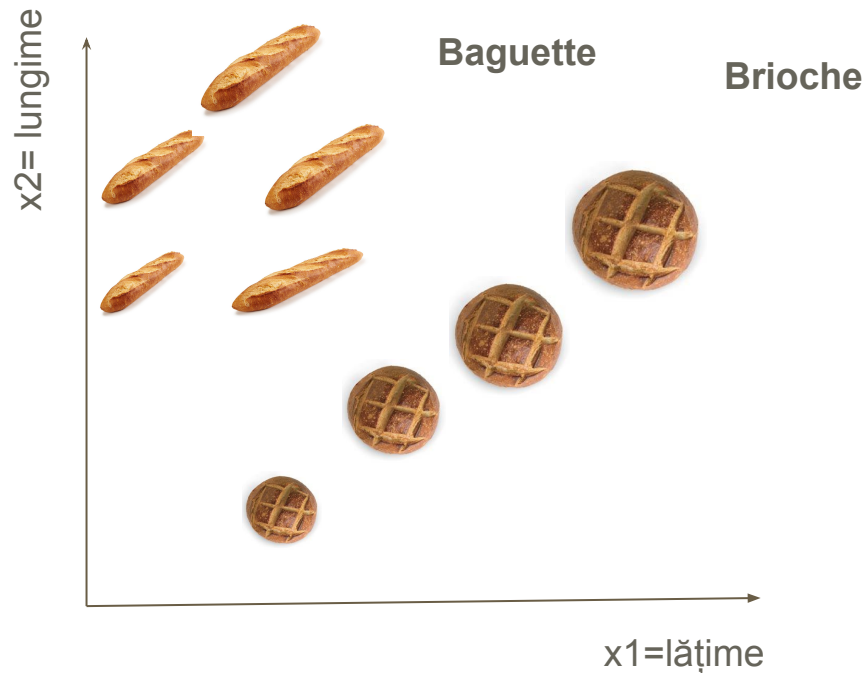


- Ce se întâmplă dacă inputul e caracterizat de mai multe dimensiuni?
- Inputul are dimensiune mai mare: \mathbb{R}^2
- Trebuie să învățăm o funcție cu mai mulți parametri

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \in \mathbb{R}^2 \quad W = [w_1 \ w_2]$$

$$\begin{aligned} y(x) &= w_1 \cdot x_1 + w_2 \cdot x_2 + b \\ &= [w_1 \ w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + b \\ &= W x + b \end{aligned}$$

Clasificare - Exemplu

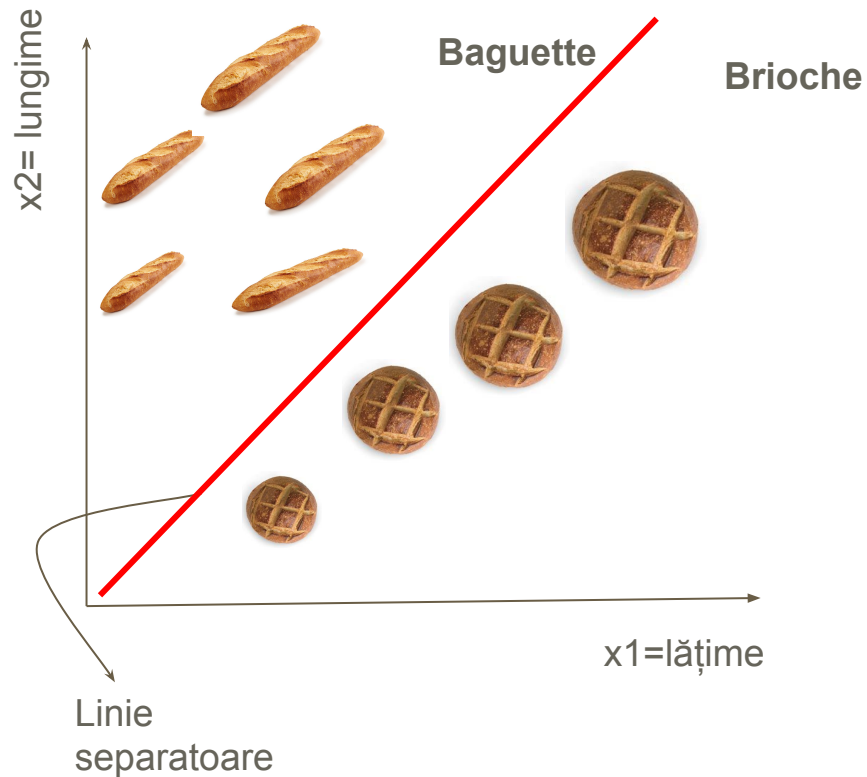


- Ce se întâmplă dacă inputul e caracterizat de mai multe dimensiuni?
- Inputul are dimensiune mai mare: \mathbb{R}^N
- Trebuie să învățăm o funcție cu mai mulți parametri

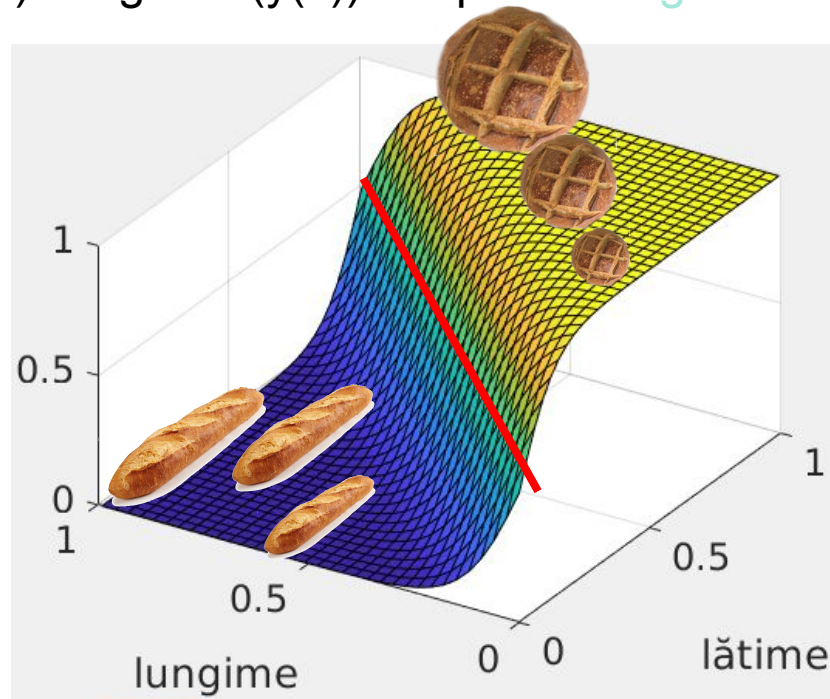
$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \in \mathbb{R}^N \quad W = [w_1 w_2 \dots w_N]$$

$$\begin{aligned} y(x) &= w_1 \cdot x_1 + w_2 \cdot x_2 + \dots + w_N \cdot x_N + b \\ &= [w_1 w_2 \dots w_N] \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} + b \\ &= W x + b \end{aligned}$$

Clasificare - Exemplu

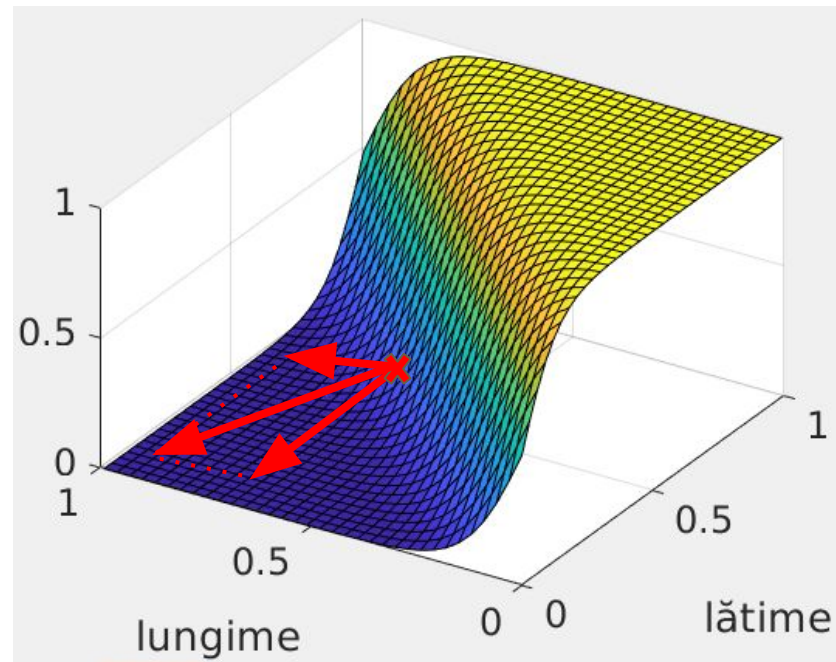


- Vrem să învățăm W astfel încât
 $a(x) = \text{sigmoid}(y(x)) = 1$ pentru brioche
 $a(x) = \text{sigmoid}(y(x)) = 0$ pentru baguette



Gradientul în două dimensiuni

- Avem doi parametri de modificat
- Putem calcula pentru fiecare parametru în parte direcția pentru care Eroarea scade cel mai mult
- Direcția finală e formată din compunerea acestor direcții

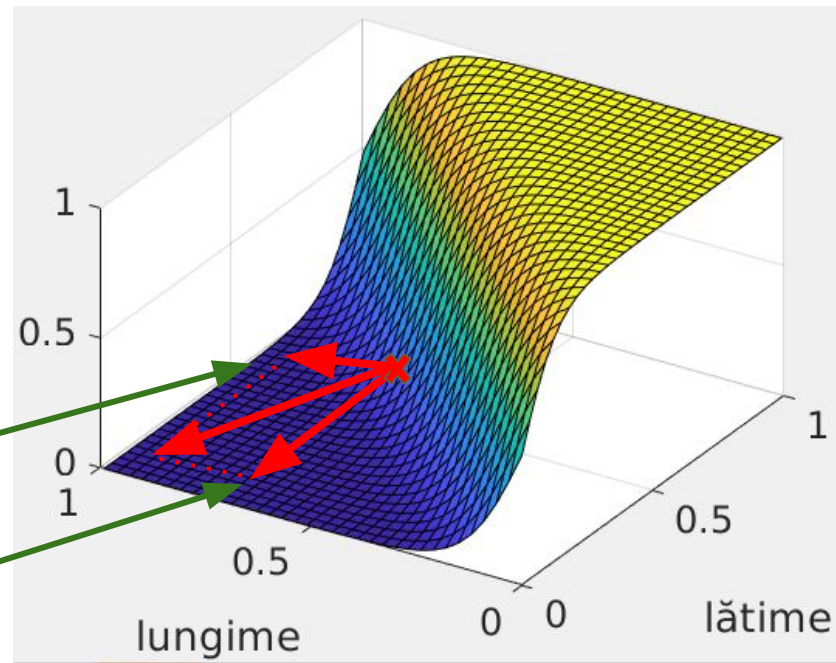


Gradientul în două dimensiuni

$$y = \sigma(w_1 * \text{latime} + w_2 * \text{lungime})$$

- Calculăm **derivatele parțiale**
- Formăm **gradientul**

$$\nabla E = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \frac{\partial E}{\partial w_2} \end{bmatrix}$$



Cum calculăm $\frac{dE}{dw}$ pentru un singur exemplu

Regresie liniară

Clasificare: Regresie logistică

Cum calculăm $\frac{dE}{dw}$ pentru un singur exemplu

Regresie liniară

Clasificare: Regresie logistică

$$y = w \cdot x + b$$

$$E = \frac{1}{2}(t - y)^2$$

$$\frac{dE}{dy} = \frac{1}{2} \cdot 2(t - y) \cdot (-1) = y - t$$

$$\frac{dE}{dw} = \frac{dE}{dy} \frac{dy}{dw} = (y - t) \cdot x$$

Cum calculăm $\frac{dE}{dw}$ pentru un singur exemplu

Regresie liniară

$$y = w \cdot x + b$$

$$E = \frac{1}{2}(t - y)^2$$

$$\frac{dE}{dy} = \frac{1}{2} \cdot 2(t - y) \cdot (-1) = y - t$$

$$\frac{dE}{dw} = \frac{dE}{dy} \frac{dy}{dw} = (y - t) \cdot x$$

Clasificare: Regresie logistică

$$y = w \cdot x + b$$

$$a = \text{sigmoid}(y) = \frac{1}{1 + e^{-y}}$$

$$E = -t \log a - (1 - t) \log(1 - a)$$

$$\frac{dE}{da} = -\frac{t}{a} + \frac{1 - t}{1 - a}$$

$$\frac{da}{dy} = \dots = a(1 - a)$$

$$\frac{dE}{dy} = \frac{dE}{da} \frac{da}{dy} = \left(-\frac{t}{a} + \frac{1 - t}{1 - a}\right) a(1 - a) = a - t$$

$$\frac{dE}{dw} = \frac{dE}{dy} \frac{dy}{dw} = (a - t)x$$

1.2. Rețele complet conectate - Intro

Andrei Roman

Software Developer

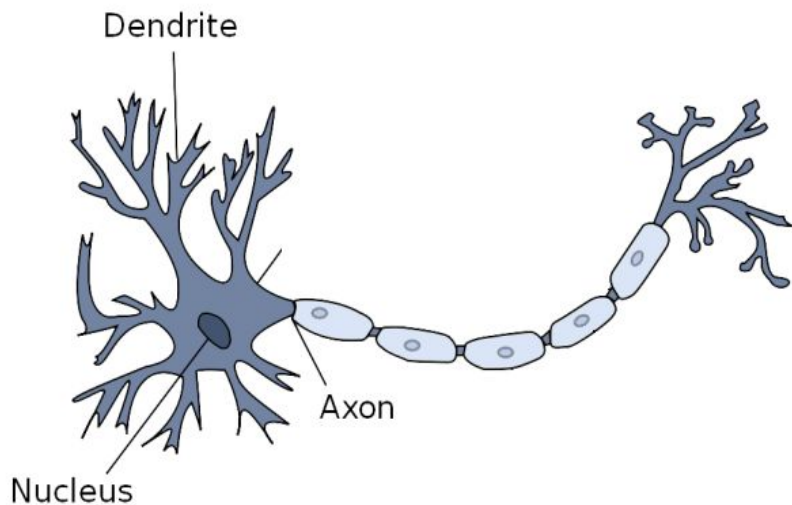
@Teamnet




Neuroni

- Mașini de calcul din natură
- Cum funcționează neuronii?
 - Transmit **semnale electrice** de la un capăt la celălalt
 - (sinapse -> dendrite -> nucleu -> axon -> sinapse)
 - Apoi la următorul neuron, etc.
- Creierul unui animal vs computer:
 - **Fuziness** vs **calcul secvențial**
- Creierul uman:
 - 100 de miliarde (10^{11}) de neuroni
 - **Conștiința**:
 - un mister so far
 - DAR: știm suficient cât să-i punem la treabă **programatic** în rezolvarea unor probleme.

Neuron

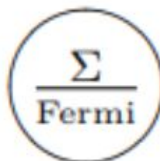


Neuroni	Nume	Imagine
10^4	Furnică	
$5 \cdot 10^7$	Liliac	
$1.6 \cdot 10^8$	Câine	
$3 \cdot 10^8$	Pisică	
$2 \cdot 10^{11}$	Elefant	

Neuroni (2)

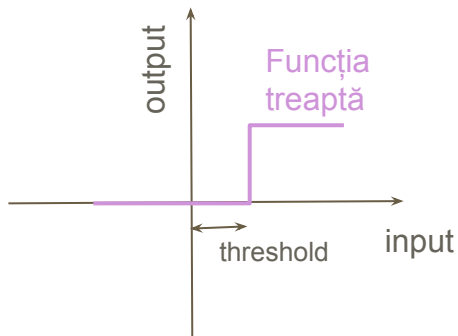
- **Reprezentare**
 - Nu putem folosi funcții liniare ($\text{out} = c1 * \text{in} + c2$)
- **Considerente**
 - Nu acționează instantaneu
 - Nu dau mai departe **input-ul** până nu ating un **threshold...**
 - ... apoi produc **output-ul**
 - **Analogie:**
 - Precum o cană care nu se varsă până nu este plină.
 - **Altfel spus...**
 - ...fără zgomot.

Neuron (2)



Neuroni	Nume	Task
$6 \cdot 10^5$	LeNet	MNIST
$6 \cdot 10^7$	AlexNet	ImageNet 15.3%
$1.3 \cdot 10^8$	VGG-16	ImageNet 7.3%
$4 \cdot 10^6$	GoogleNet	ImageNet 6.67%
$2.5 \cdot 10^7$	ResNet-50	ImageNet 3.6%

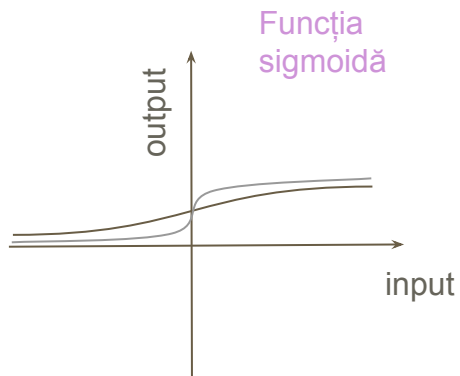
Neuroni (3)



Funcția treaptă

Valori de intrare mici
=> output = 0

Trecem de prag (threshold)
=> neuronii se activează (fire)
=> avem output



Funcția logistică (sigmoidă)

Funcția de activare utilizată mai departe.
De ce?

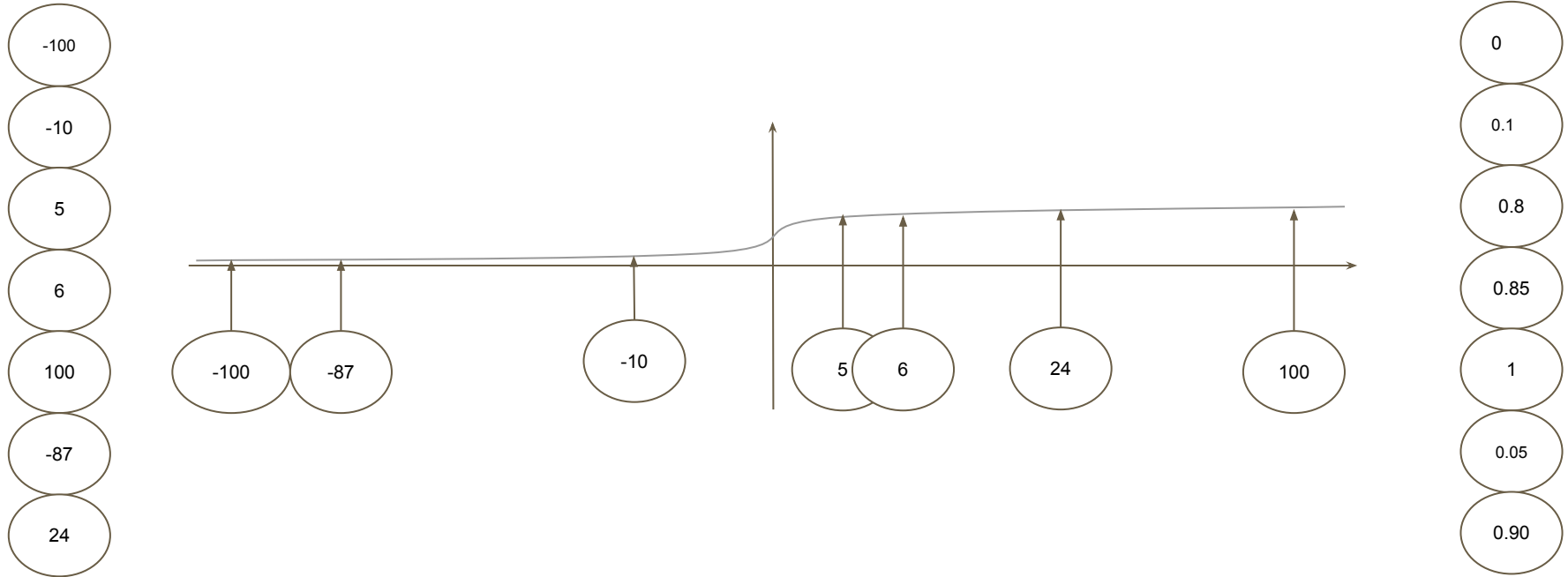
- Mai smooth decât funcția treaptă
- Simplitate
- Scurtături pentru calcule

Funcția sigmoidă de activare

$$y = \frac{1}{1+e^{-x}}$$

Taie axa la $y = 1/2$:
 $x = 0 \Rightarrow y = 1/2$

Neuron (3)



Neuroni (4)

Inputs

a

b

c

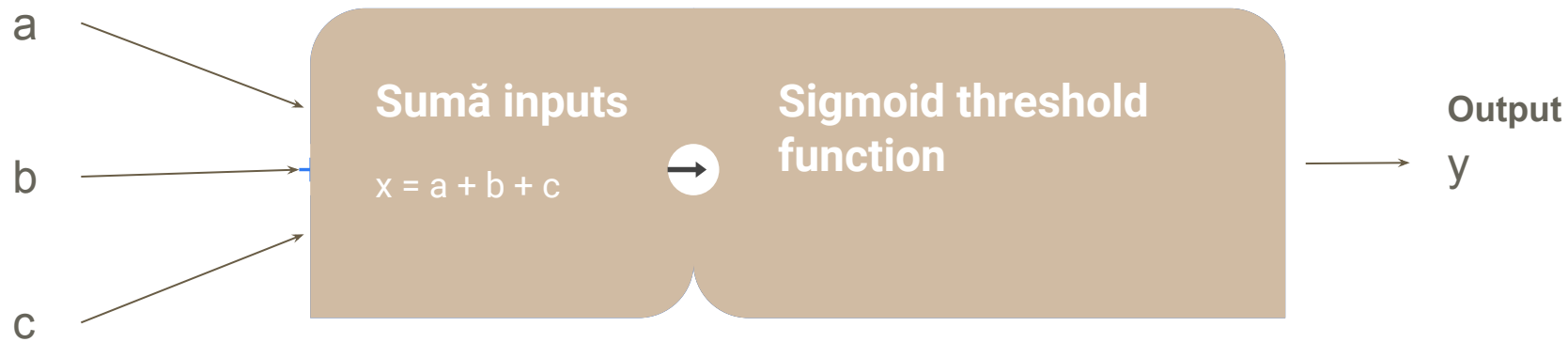
Sumă inputs

$$x = a + b + c$$

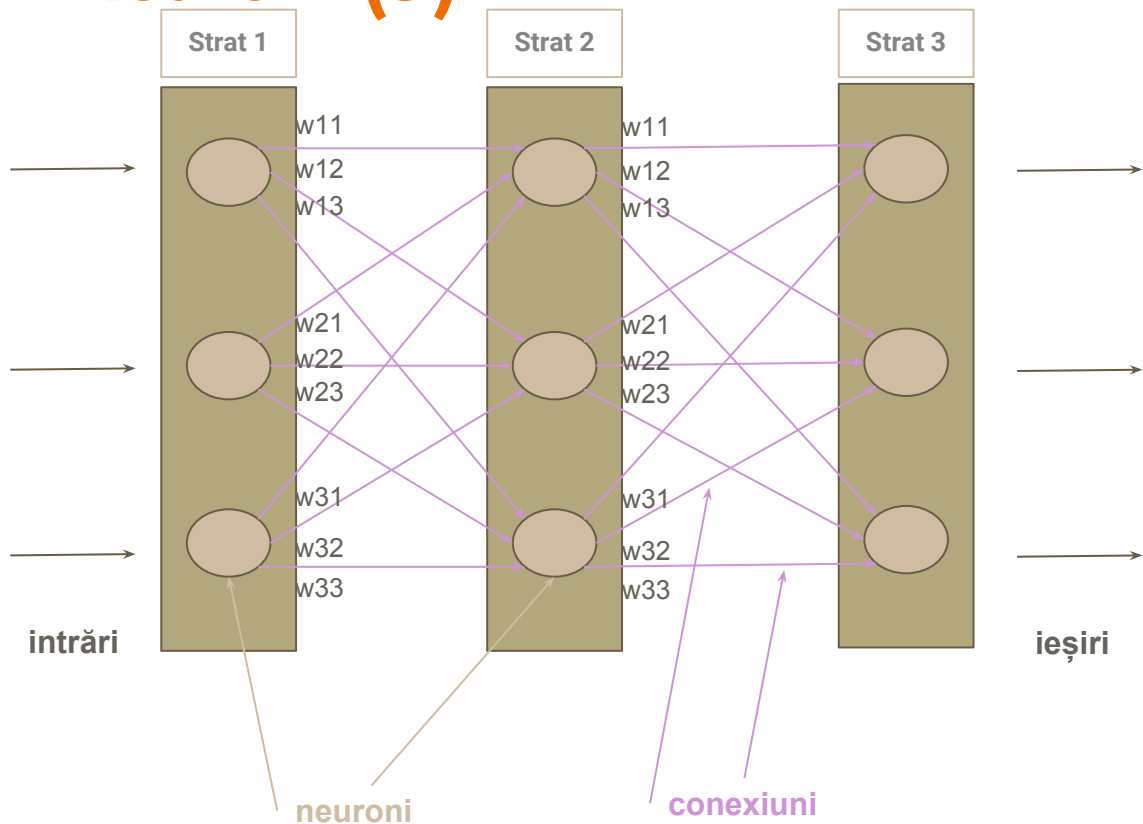


Sigmoid threshold
function

Output
y



Neuroni (5)

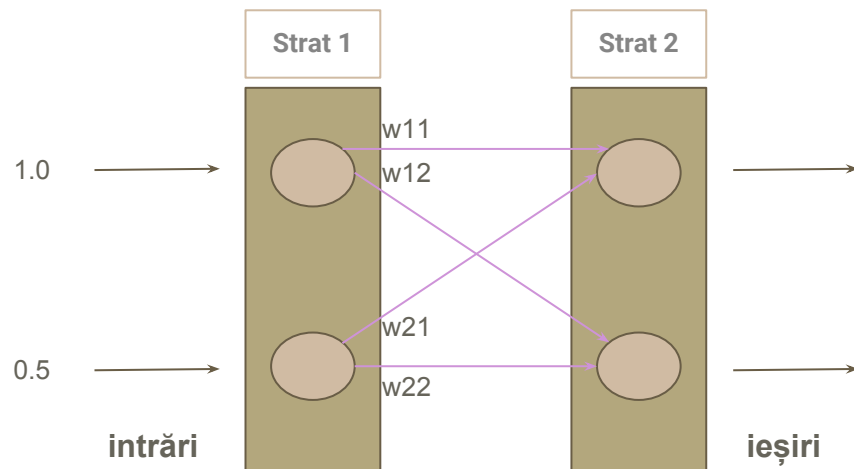


- Exemplu de rețea neuronală **fully connected** cu 3 straturi a câte 3 neuroni fiecare strat.

Învățarea:

- Ajustăm weight-urile.
- **Weight:**
 - Slabe: diminuează semnalul
 - Puternice: amplifică semnalul (au mai mult aport la decizie).

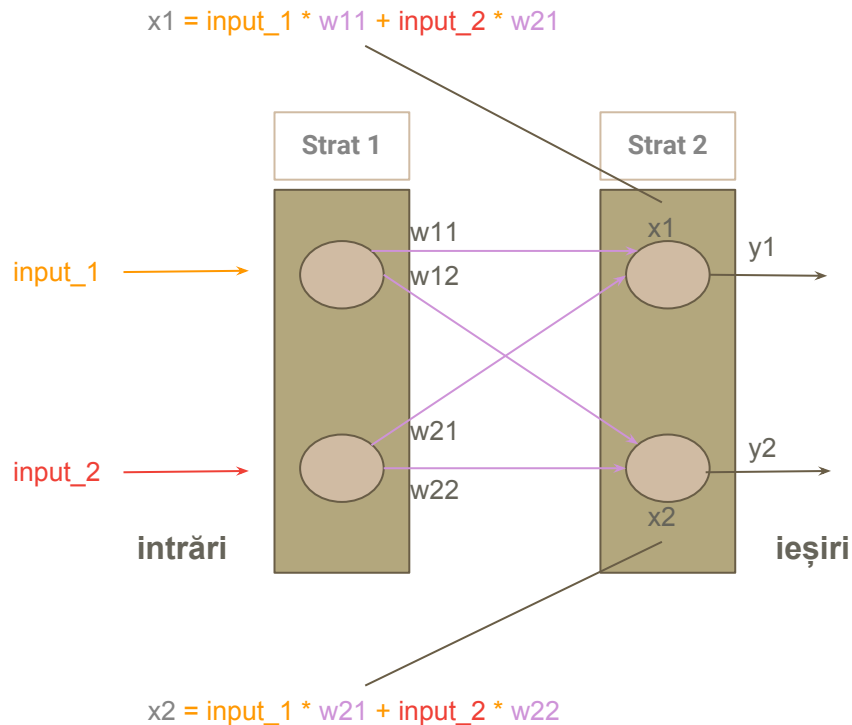
Propagarea semnalelor prin rețea



- Weight-uri aleatorii:
 - $w11 = 0.9, w12 = 0.2$
 - $w21 = 0.3, w22 = 0.8$
- Funcția de activare:

$$y = \frac{1}{1+e^{-x}}$$

Propagarea semnalelor prin rețea (2)



Layer 1: Reprezentare intrări

Layer 2:

- Input-ul moderat al neuronului 1 de pe layer 2:

$$x_1 = 1.0 * 0.9 + 0.5 * 0.3 = 1.05$$

- Output-ul neuronului 1 de pe layer 2:

$$y_1 = \frac{1}{1+0.3499} = 0.7408$$

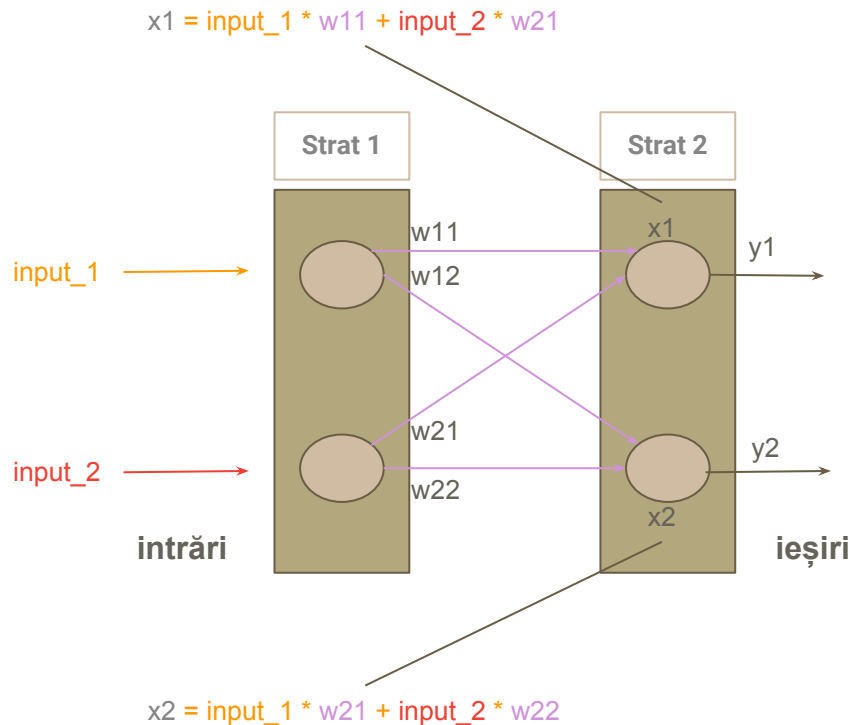
- Input-ul moderat al neuronului 2 de pe layer 2:

$$x_2 = 1.0 * 0.2 + 0.5 * 0.8 = 0.6$$

- Output-ul neuronului 2 de pe layer 2:

$$y_2 = \frac{1}{1+0.5488} = 0.6457$$

Propagarea semnalelor prin rețea (2)



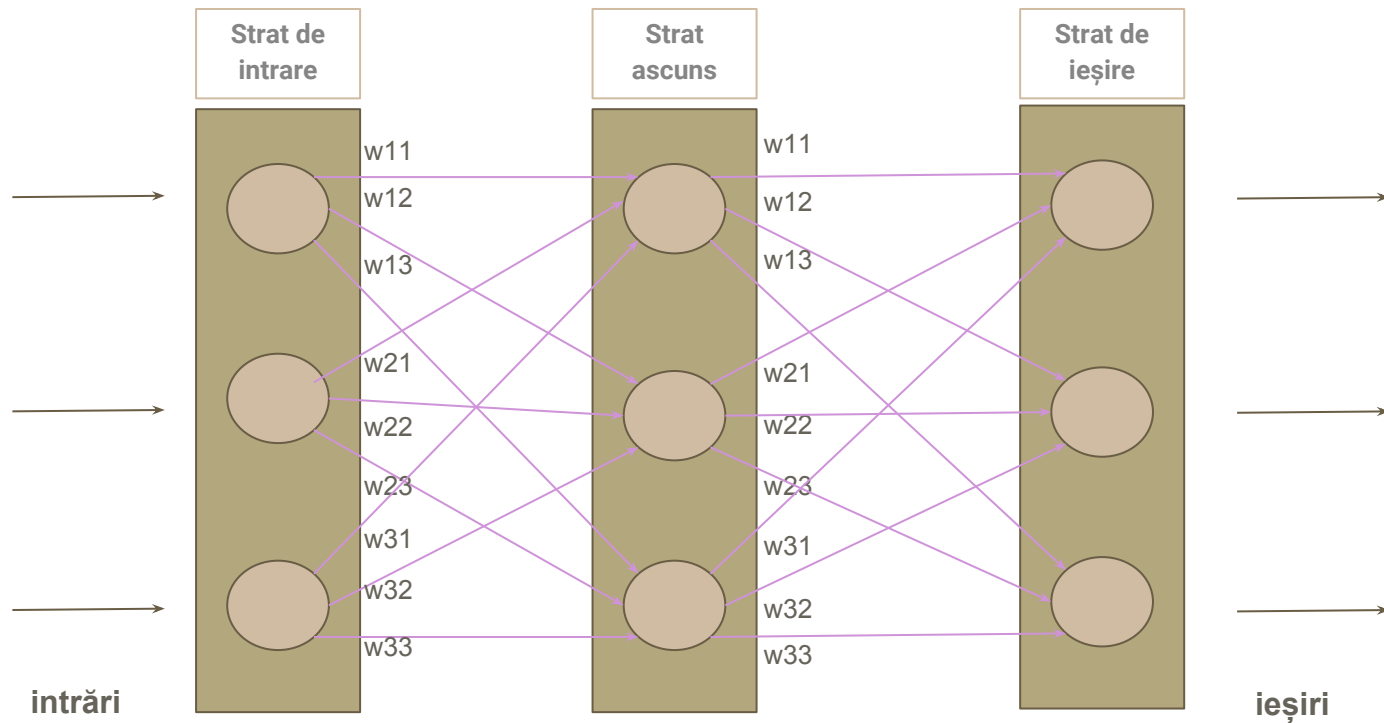
- Folosind calcul matriceal:

$$W * I = X$$

$$\begin{pmatrix} w11 & w21 \\ w12 & w22 \end{pmatrix} \begin{pmatrix} \text{input_1} \\ \text{input_2} \end{pmatrix} =$$
$$= \begin{pmatrix} \text{input_1} * w11 + \text{input_2} * w21 \\ \text{input_1} * w12 + \text{input_2} * w22 \end{pmatrix}$$

Semnalul moderat combinat care intră în layer 2

Propagarea semnalelor prin rețea (3)



Propagarea semnalelor prin rețea (4)

$$I = \begin{pmatrix} 0.9 \\ 0.1 \\ 0.8 \end{pmatrix}$$

Input-uri din
setul de
training

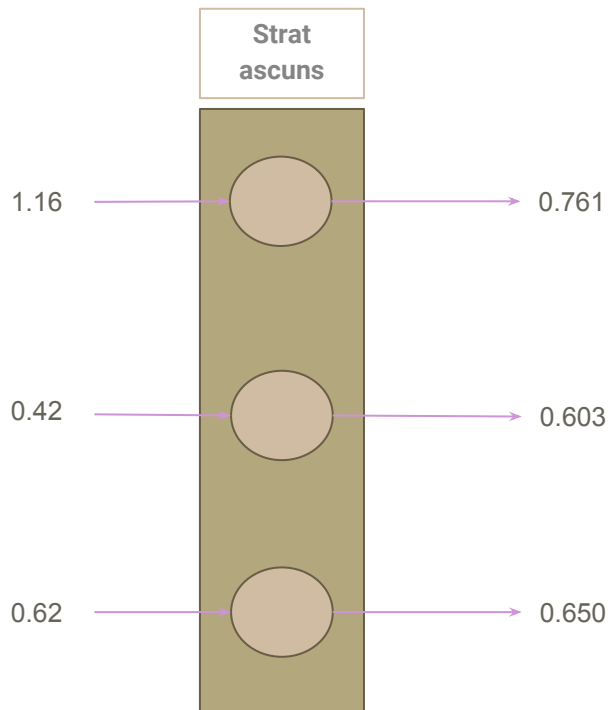
$$W_{\text{input_hidden}} = \begin{pmatrix} 0.9 & 0.3 & 0.4 \\ 0.2 & 0.8 & 0.2 \\ 0.1 & 0.5 & 0.6 \end{pmatrix}$$

Alegem aleatoriu weight-urile
conexiunilor dintre neuronii stratului
de **intrare** și cei de pe stratul **ascuns**.

$$W_{\text{hidden_output}} = \begin{pmatrix} 0.3 & 0.7 & 0.5 \\ 0.6 & 0.5 & 0.2 \\ 0.8 & 0.1 & 0.9 \end{pmatrix}$$

Alegem aleatoriu weight-urile
conexiunilor dintre neuronii
stratului **ascuns** și cel de **ieșire**.

Propagarea semnalelor prin rețea (5)



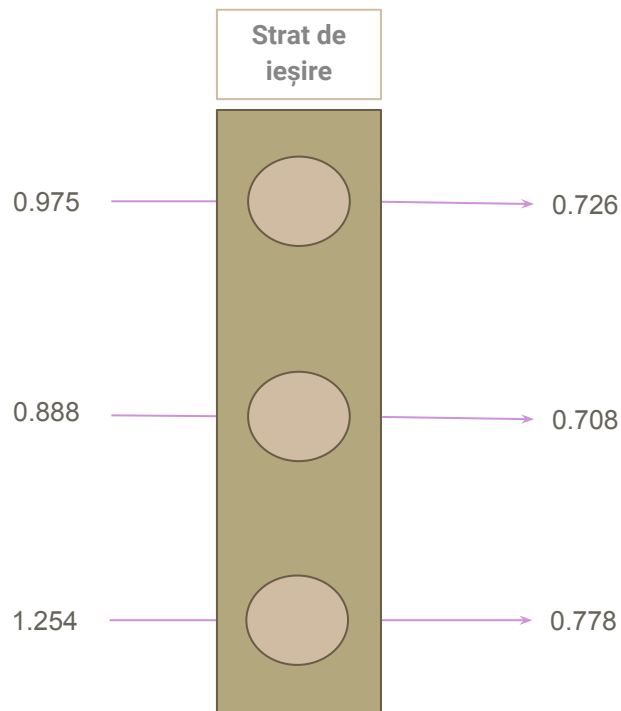
$$X_{\text{hidden}} = W_{\text{input_hidden}} * I$$

$$O_{\text{hidden}} = \text{sigmoid}(X_{\text{hidden}})$$

$$X_{\text{hidden}} = \begin{pmatrix} 0.9 & 0.3 & 0.4 \\ 0.2 & 0.8 & 0.2 \\ 0.1 & 0.5 & 0.6 \end{pmatrix} \begin{pmatrix} 0.9 \\ 0.1 \\ 0.8 \end{pmatrix} = \begin{pmatrix} 0.16 \\ 0.42 \\ 0.62 \end{pmatrix}$$

$$O_{\text{hidden}} = \text{sigmoid} \begin{pmatrix} 0.16 \\ 0.42 \\ 0.62 \end{pmatrix} = \begin{pmatrix} 0.761 \\ 0.603 \\ 0.650 \end{pmatrix}$$

Propagarea semnalelor prin rețea (6)



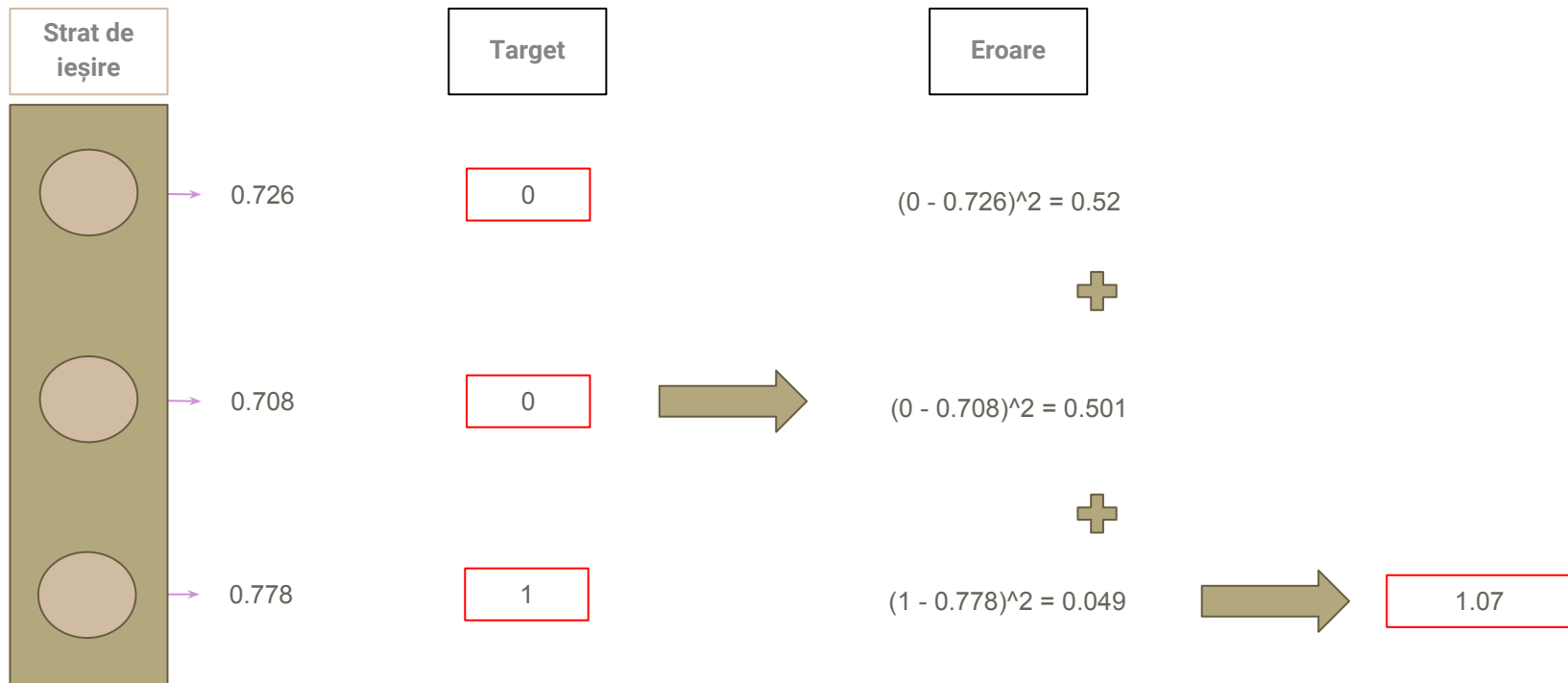
$$X_{\text{output}} = W_{\text{hidden_output}} * O_{\text{hidden}}$$

$$O_{\text{output}} = \text{sigmoid}(X_{\text{output}})$$

$$X_{\text{output}} = \begin{pmatrix} 0.3 & 0.7 & 0.5 \\ 0.6 & 0.5 & 0.2 \\ 0.8 & 0.1 & 0.9 \end{pmatrix} \begin{pmatrix} 0.761 \\ 0.603 \\ 0.650 \end{pmatrix} = \begin{pmatrix} 0.975 \\ 0.888 \\ 1.254 \end{pmatrix}$$

$$O_{\text{output}} = \text{sigmoid} \begin{pmatrix} 0.975 \\ 0.888 \\ 1.254 \end{pmatrix} = \begin{pmatrix} 0.726 \\ 0.708 \\ 0.778 \end{pmatrix}$$

Eroare



Problema

Sa se scrie ecuatile

Input - vector cu 784 de neuroni

Strat ascuns 1 - cu 100 de neuroni

Strat ascuns 2 - cu 10 neuroni

Eroare

Solutie

$$y1 = I * W1$$

$$a1 = \text{sigmoid}(y1)$$

$$y2 = a1 * W2$$

$$a2 = \text{sigmoid}(y2)$$

$$E = \text{SSE}(a2, T)$$

1.3. Rețele Neurale (Backpropagation)

Ruxandra Burtică

Computer Scientist

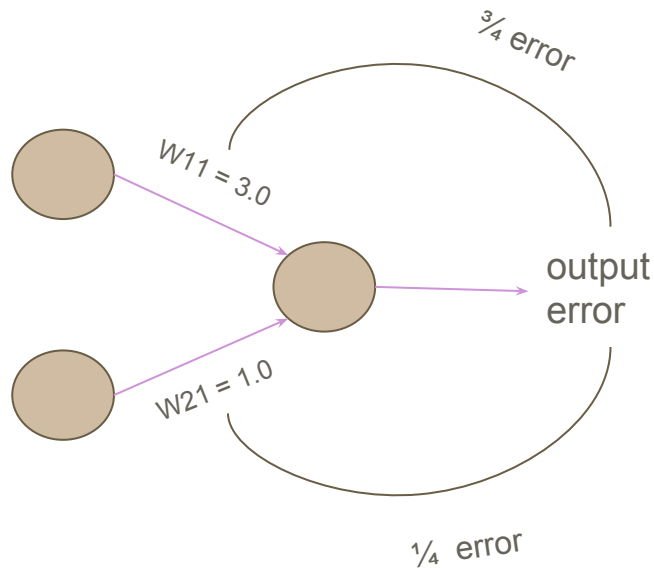
@Adobe



Backpropagation

- În general, pentru a calcula gradientul fiecărui strat vom avea doua operații
- **Forward pass:**
 - Calculăm toate rezultatele intermediare, ieșirile și eroarea
- **Backward pass:**
 - Calculăm pe rând gradientii fiecărui strat, de la ultimul la primul

Backpropagation (2)



Propagarea înapoi a erorilor de la un nod:

- Update proporțional
- Propagare cantitate eroare finală în funcție de cât de puternic este link-ul.

Notă:

- Abia pe stratul de ieșire se poate vorbi despre o eroare!

Calcularea erorii

Ce funcție alegem pentru eroare?

1. $(\text{target} - \text{actual})$
2. $|\text{target} - \text{actual}|$
3. $(\text{target} - \text{actual})^2$

Calcularea erorii

Ce funcție alegem pentru eroare?

1. **(target - actual)** - NU! Anulare erori cu + vs erori cu -.
2. **|target - actual|** - NU! Nu este continuă în punctul de minim.
3. **(target - actual)²** - DA. De ce?
 - a. Ecuații simple.
 - b. Funcție netedă și continuă
 - i. Fără sărituri abrupte
 - c. Gradientul se micșorează în apropierea minimului
 - i. Fără overshoot dacă micșorăm rata de învățare aici (pasul)

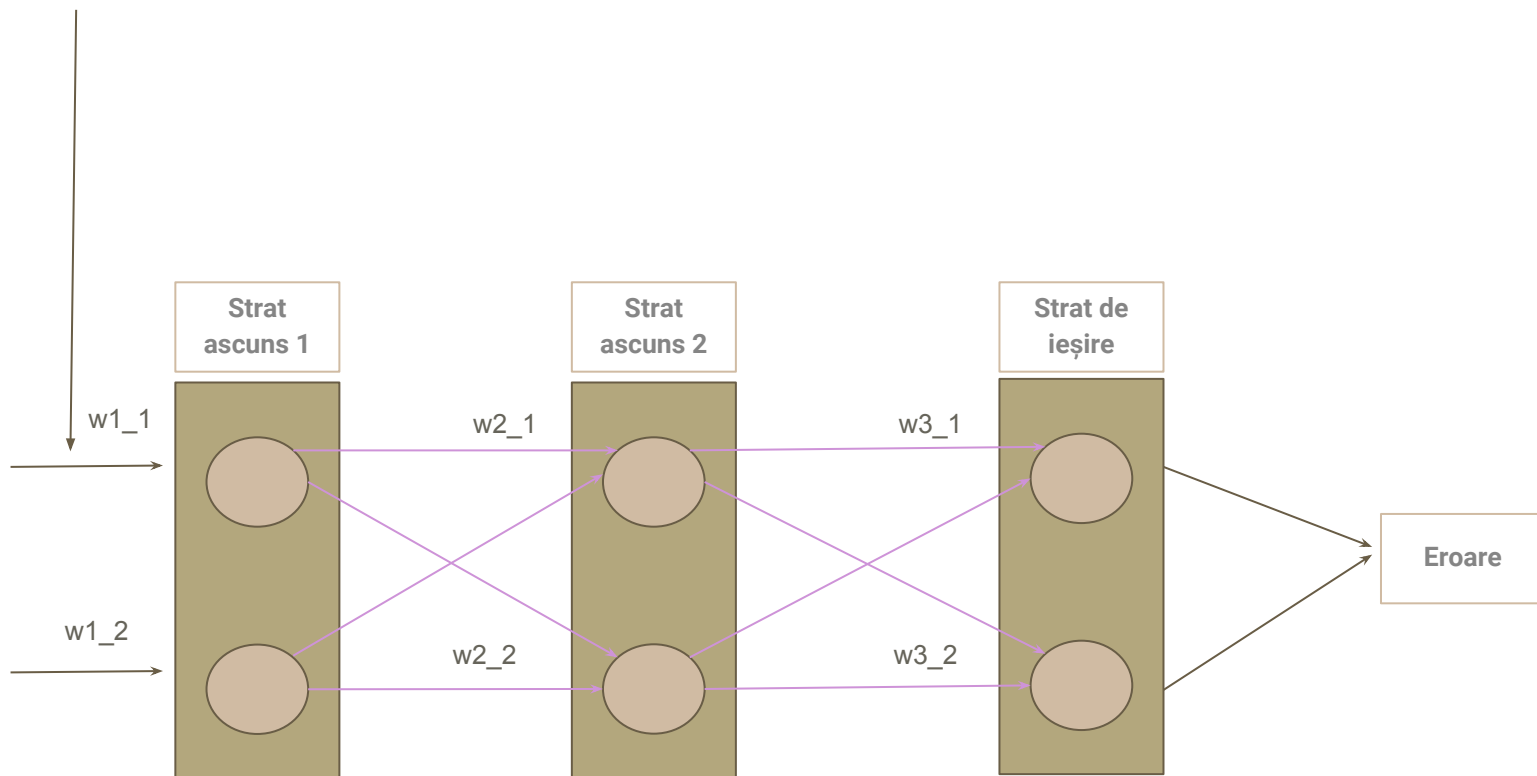
Cum updatăm parametrii rețelei?

- Pentru orice rețea neurală, trebuie să știm ce efect au perturbații mici ale parametrilor în funcția de cost.
- Variem fiecare parametru în parte

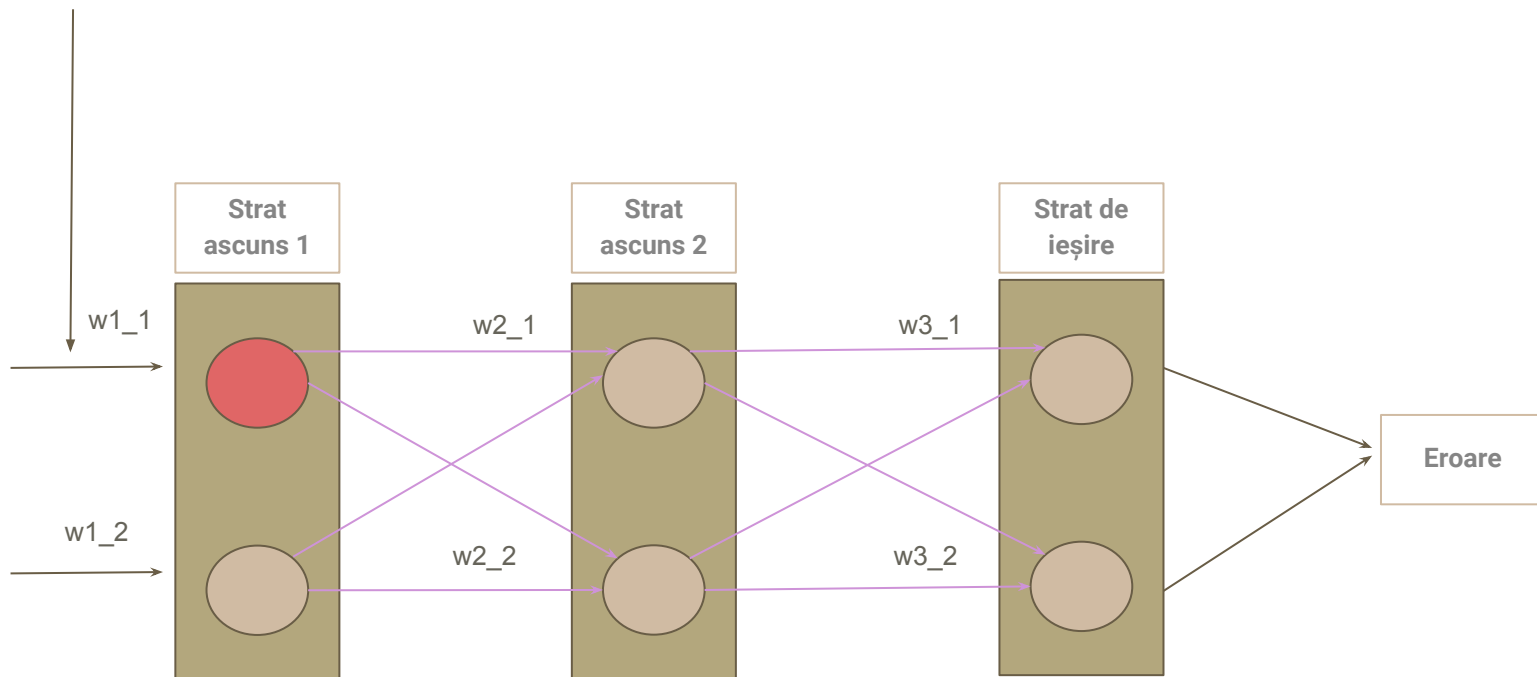
$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} \rightarrow E(W)$$

$$\begin{bmatrix} w_1 + \Delta w_1 \\ w_2 \\ \vdots \\ w_N \end{bmatrix} \rightarrow E(W + \Delta w_1)$$

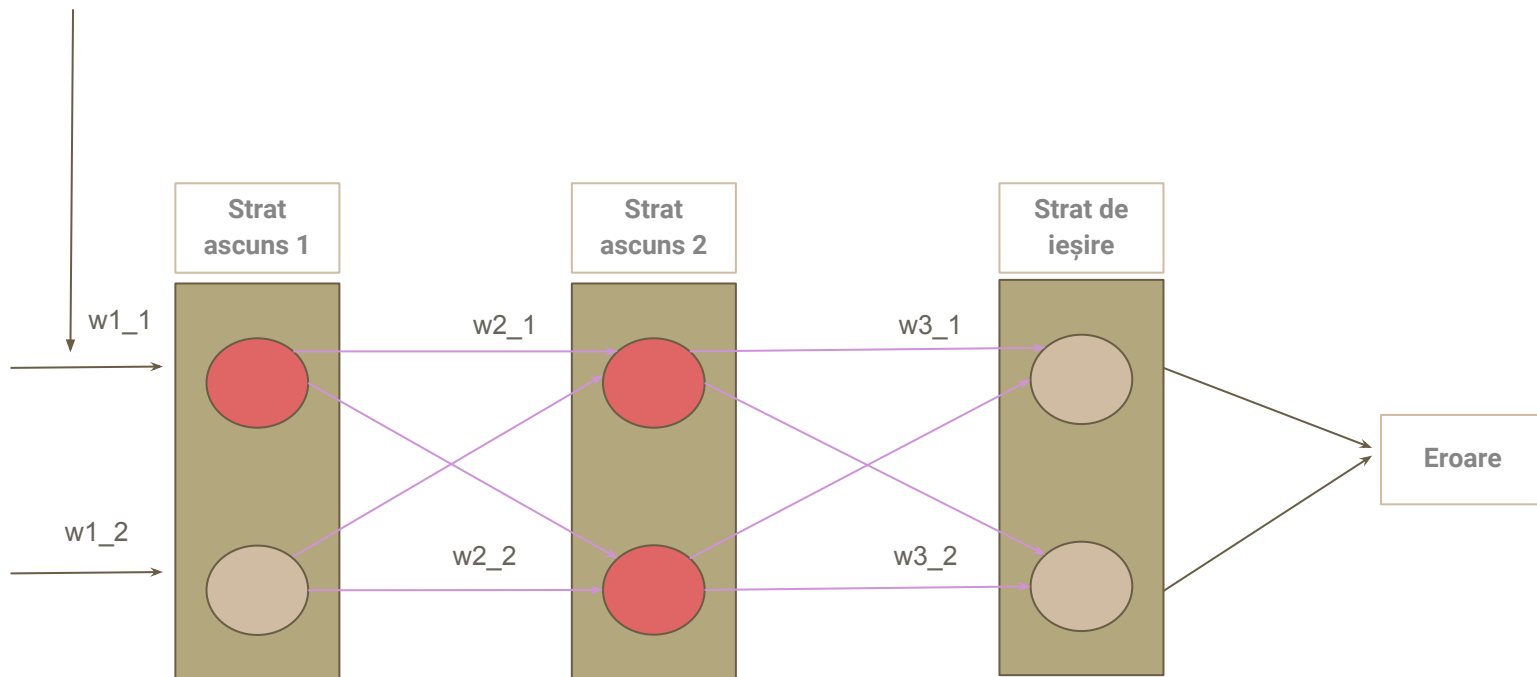
- Perturbăm un singur parametru



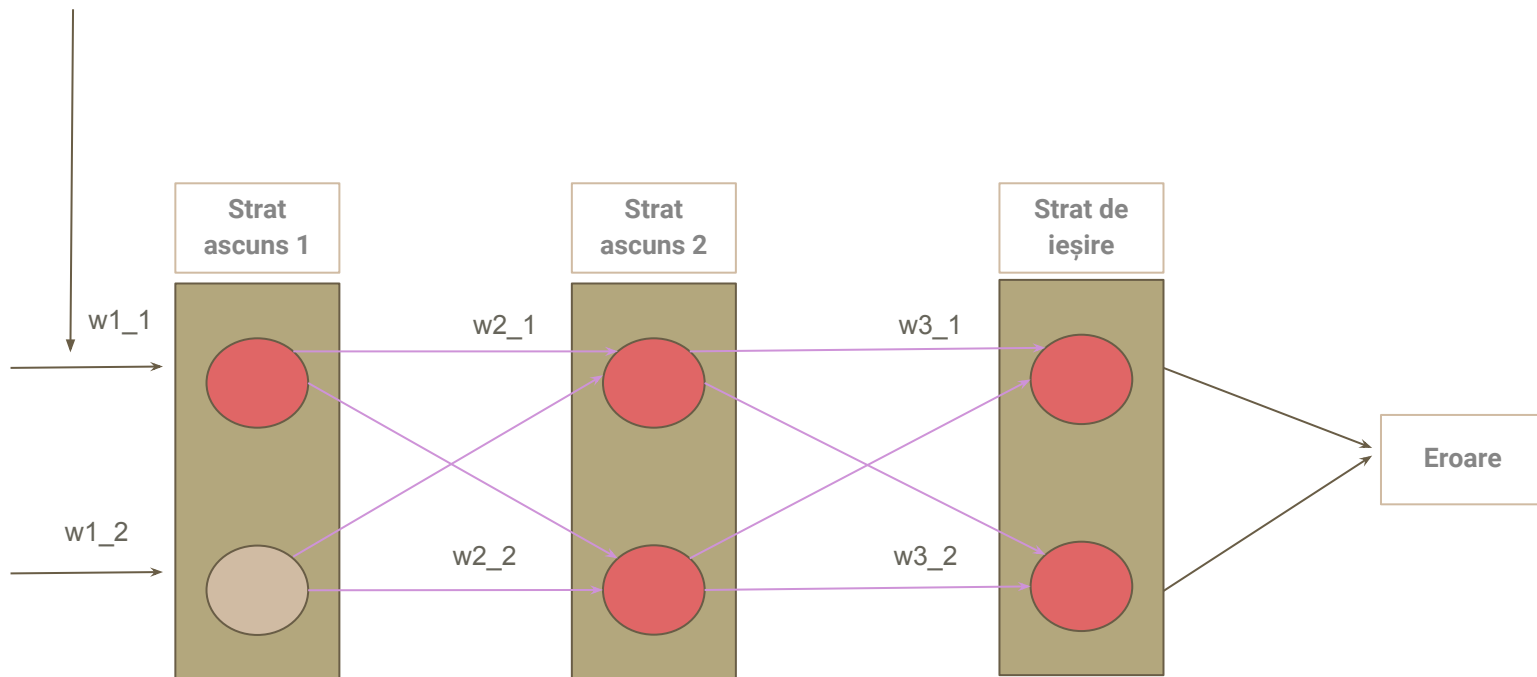
- Perturbăm un singur parametru
- Perturbațiile se propagă prin toții neuronii influențați de primul parametru



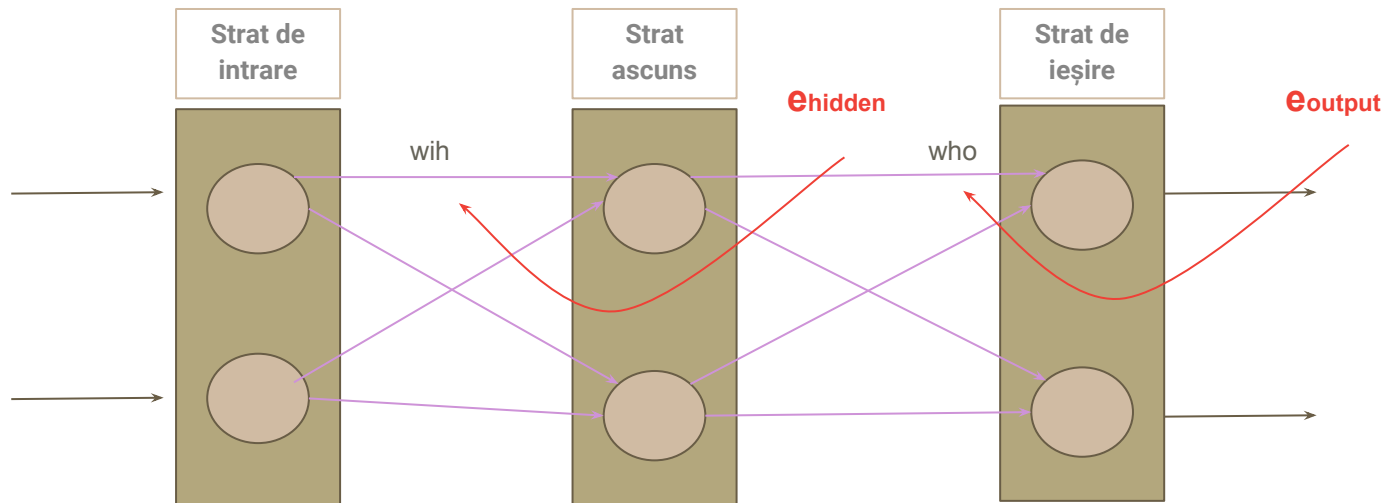
- Perturbăm un singur parametru
- Perturbațiile se propagă prin toții neuronii influențați de primul parametru



- Perturbăm un singur parametru
- Perturbațiile se propagă prin toții neuronii influențați de primul parametru

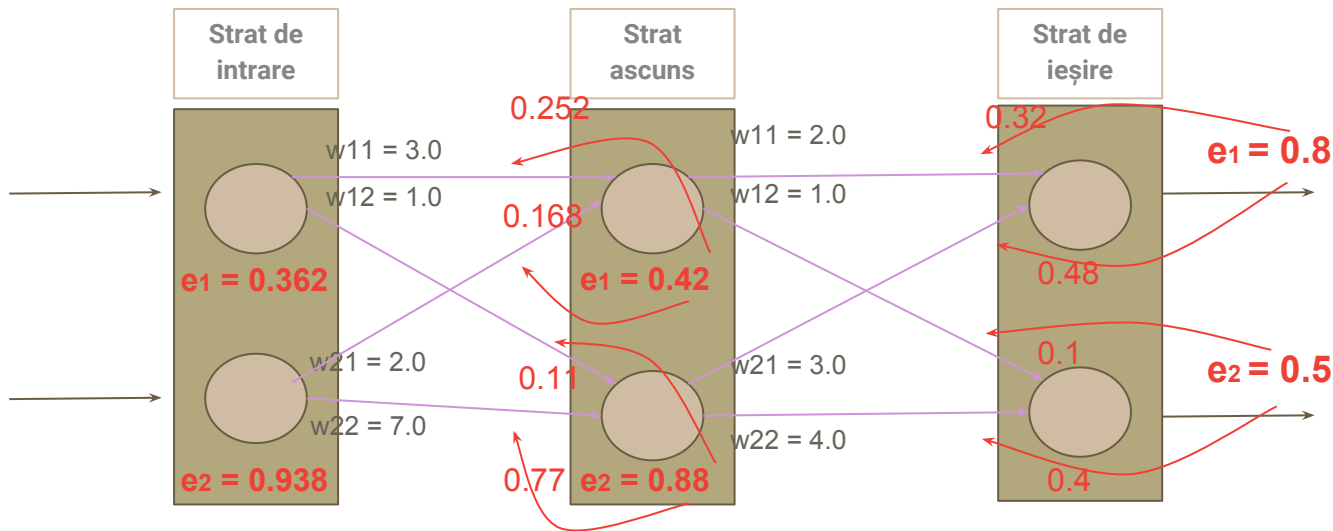


Propagarea erorii



- wih - weight-urile link-urilor dintre neuroni input și hidden
- who - weight-urile link-urilor dintre neuroni hidden și output

Propagarea erorii (2)



$e_{\text{hidden}, 1} = \text{sum}(\text{split errors on links } w_{11} \text{ and } w_{12})$

$$= e_{\text{output}, 1} * w_{11} / (w_{11} + w_{21}) + e_{\text{output}, 2} * w_{12} / (w_{12} + w_{22})$$

Gradient

- Ar fi ineficient să calculăm numeric fiecare din aceste perturbații
- Precum și la modelele liniare, putem observa cum se modifică eroarea în funcție de perturbații prin derivatele (parțiale)

- Avem câte o derivată parțială pentru fiecare parametru $\frac{\partial E}{\partial w_i} \in \mathbb{R}$

- Vectorul tuturor derivatelor parțiale: gradientul $\nabla E = \begin{bmatrix} \frac{\partial E}{\partial w_1} \\ \vdots \\ \frac{\partial E}{\partial w_N} \end{bmatrix}$

Gradient

- Calculăm gradienti pentru fiecare strat din rețea

$$\begin{bmatrix} \frac{\partial E}{\partial w_1^1} \\ \vdots \\ \frac{\partial E}{\partial w_N^1} \end{bmatrix}$$

- Vom vedea că gradientii din stratul 3 depind de gradientii din stratul 2, cei din stratul 2 depind de cei din stratul 1
- Spunem că **gradientii se propagă** din ultimele straturi spre primele
- După calcularea gradientilor, putem face update la weight-uri:

$$w_{jk} = w_{jk} - \alpha \frac{\partial E}{\partial w_{jk}}$$

Chain rule

- Dacă o funcție $\mathbf{f}(\mathbf{x}, \mathbf{y})$ depinde de de mai multe variabile, putem calcula derivate parțiale pentru fiecare
- Dacă o variabilă a funcției \mathbf{y} depinde de altă variabila \mathbf{x} , atunci putem calcula derivatele parțiale ale lui \mathbf{f} fata de \mathbf{x} folosind regula lanțului (chain rule)

$$f = f(y), y = y(x)$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial y} \frac{\partial y}{\partial x}$$

Backpropagation

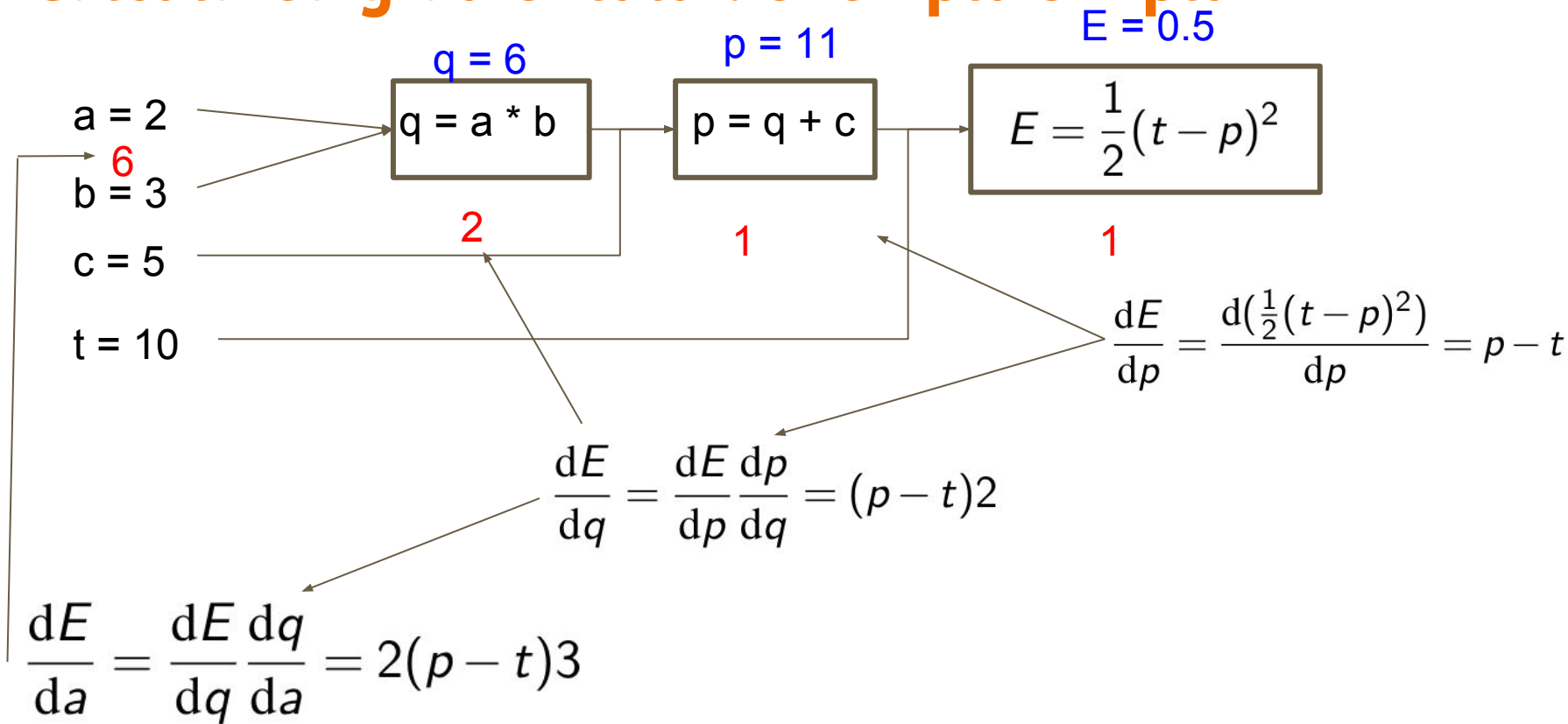
- Dacă f depinde de x_1 care depinde de x_2 care depinde de x_3 care depinde

...

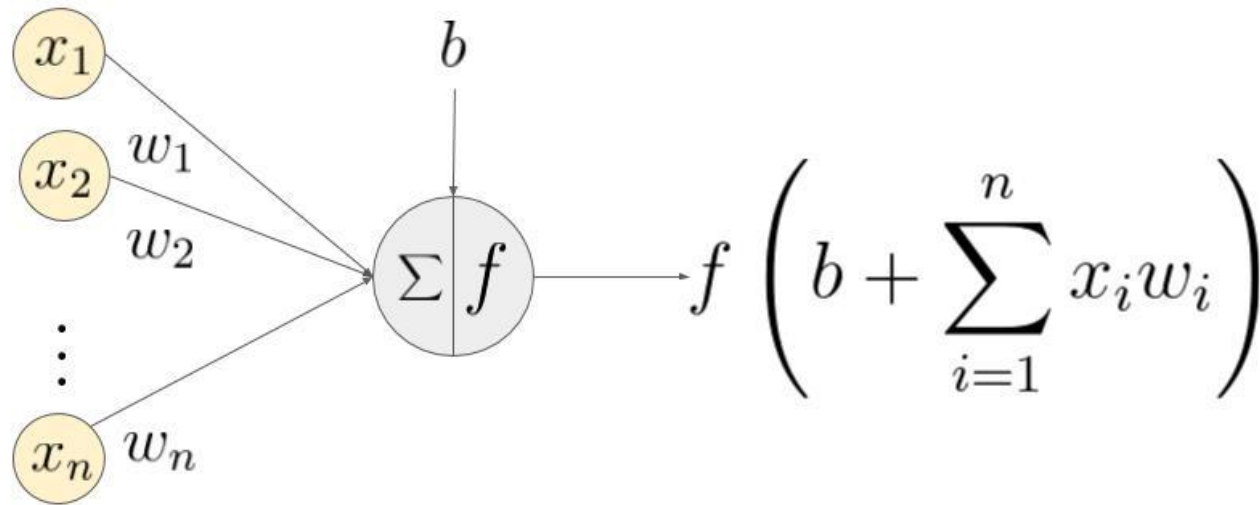
$$\frac{\partial f}{\partial x_N} = \frac{\partial f}{\partial x_1} \frac{\partial x_1}{\partial x_2} \frac{\partial x_2}{\partial x_3} \cdots \frac{\partial x_{N-1}}{\partial x_N}$$

- Cazul rețelelor neurale, unde avem intrările unor straturi depind de ieșirile straturilor precedente, iar ieșirile straturilor depind de parametrii și intrări
- Vrem să calculăm gradientul fiecărui parametru, trebuie să avem un lanț de la funcția de eroare, prin toate operațiile intermediare

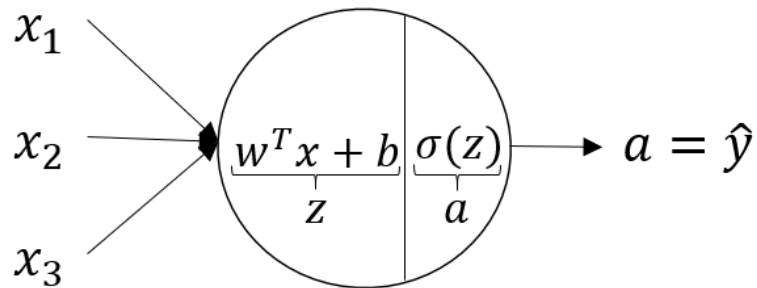
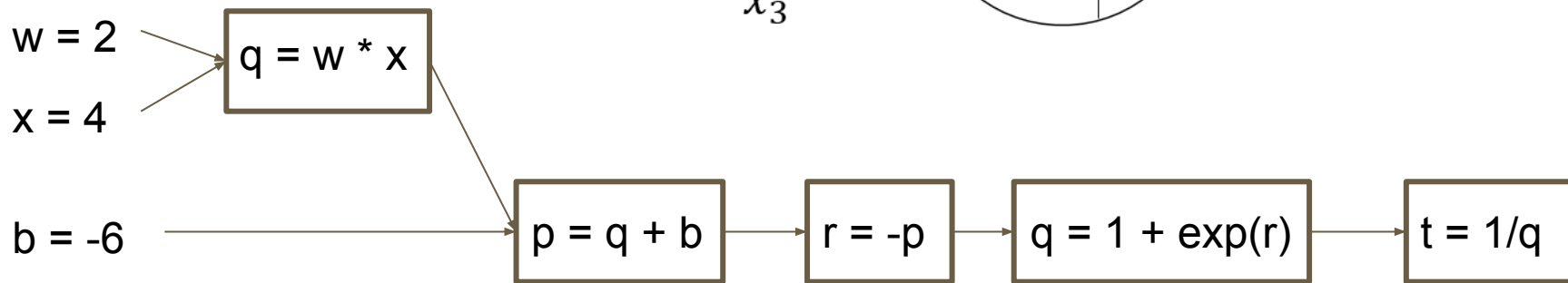
Calcularea gradientului: exemplu simplu



Calcularea gradientului: exemplu simplu (2)

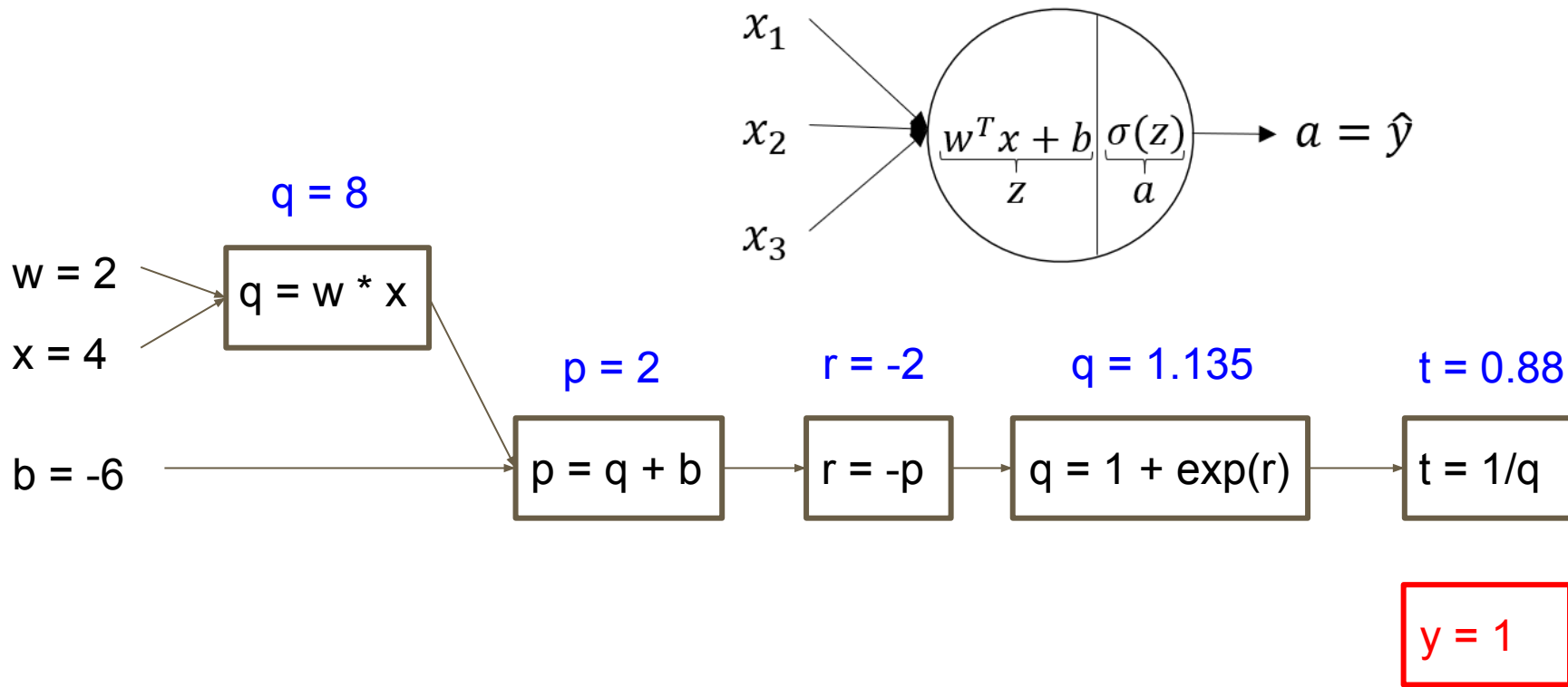


Calcularea gradientului (exercitiu)



$$y = 1$$

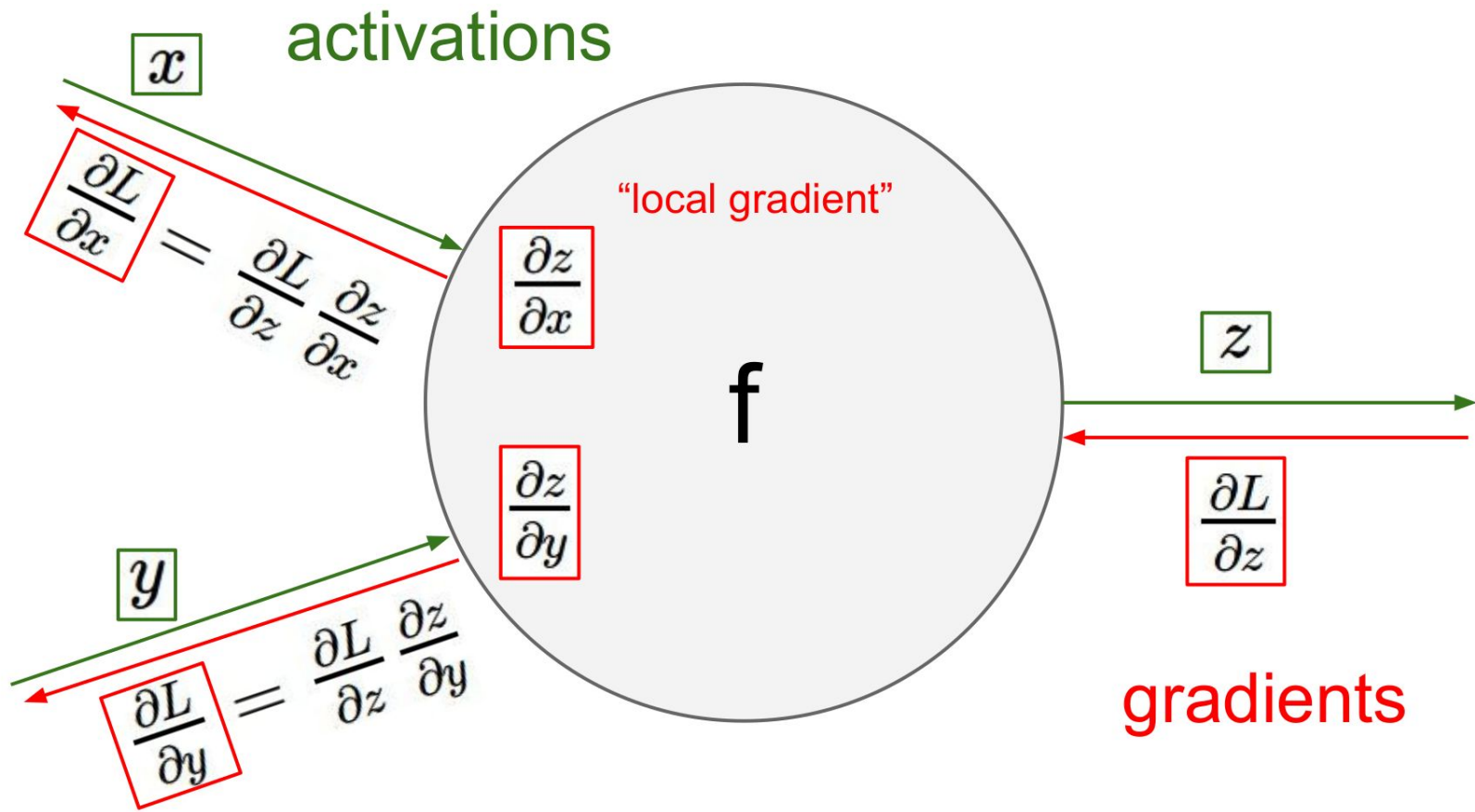
Calcularea gradientului (exercitiu)



Calcularea gradientului (hints)

$$\left(\frac{1}{x}\right)' = -\frac{1}{x^2}$$

$$\frac{\partial t}{\partial w} = \frac{\partial t}{\partial q} \cdot \frac{\partial q}{\partial r} \cdot \frac{\partial r}{\partial p} \cdot \frac{\partial p}{\partial w}$$



Actualizarea weight-urilor

(forward propagation)

$$h_1 = xW_1 + b_1$$

$$z_1 = \sigma(h_1)$$

$$z_2 = z_1W_2 + b_2$$

$$Loss = (z_2 - y)^2$$



(backward propagation)

$$\frac{\partial h_1}{\partial x} = W_1^T$$

$$\frac{\partial z_1}{\partial h_1} = \sigma'(h_1) = z_1 \circ (1 - z_1)$$

$$\frac{\partial z_2}{\partial z_1} = W_2^\top$$

$$\frac{\partial Loss}{\partial z_2} = 2(z_2 - y)$$



Actualizarea weight-urilor (2)

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} * \sum_n (t_n - o_n)^2 = \frac{\partial (t_k - o_k)^2}{\partial w_{jk}} \text{---}$$

Din suma anterioară doar $(t_k - o_k)$ are treabă w_{jk} . Restul se anulează la derivare.

↓
E - eroarea tuturor nodurilor din output

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial o_k} \frac{\partial o_k}{\partial w_{jk}} = -2(t_k - o_k) \frac{\partial o_k}{\partial w_{jk}}$$

↓
Chain rule

$$o_k = \text{sigmoid}(\sum_j w_{jk} o_j)$$

↓
Output-ul nodului k

Actualizarea weight-urilor (3)

$$\frac{\partial \text{sigmoid}(x)}{\partial x} = \text{sigmoid}(x)(1 - \text{sigmoid}(x))$$



$$\frac{\partial E}{\partial w_{jk}} = -2(t_k - o_k) \text{sigmoid}(\sum_j w_{jk} o_j) (1 - \text{sigmoid}(\sum_j w_{jk} o_j)) \frac{\partial}{\partial w_{jk}} (\sum_j w_{jk} o_j)$$



$$\frac{\partial E}{\partial w_{jk}} = -2(t_k - o_k) \text{sigmoid}(\sum_j w_{jk} o_j) (1 - \text{sigmoid}(\sum_j w_{jk} o_j)) o_j$$

Putem scăpa de doi-ul din față. Pe noi oricum ne interesează doar direcția pantei.

Actualizarea weight-urilor (4)

$$\frac{\partial E}{\partial w_{jk}} = -(\underbrace{t_k - o_k}_{\substack{\mathbf{e}_j = \text{target} - \text{actual} = \\ \text{eroarea}}}) \underbrace{\text{sigmoid}(\sum_j w_{jk} o_j)}_{\substack{\text{Semnalul ce intră în ultimul nod} \\ \mathbf{ij}}} (1 - \text{sigmoid}(\sum_j w_{jk} o_j)) \underbrace{o_j}_{\substack{\mathbf{o}_j - \text{output-ul} \\ \text{layer-ului ascuns} \\ \text{anterior } j}}$$

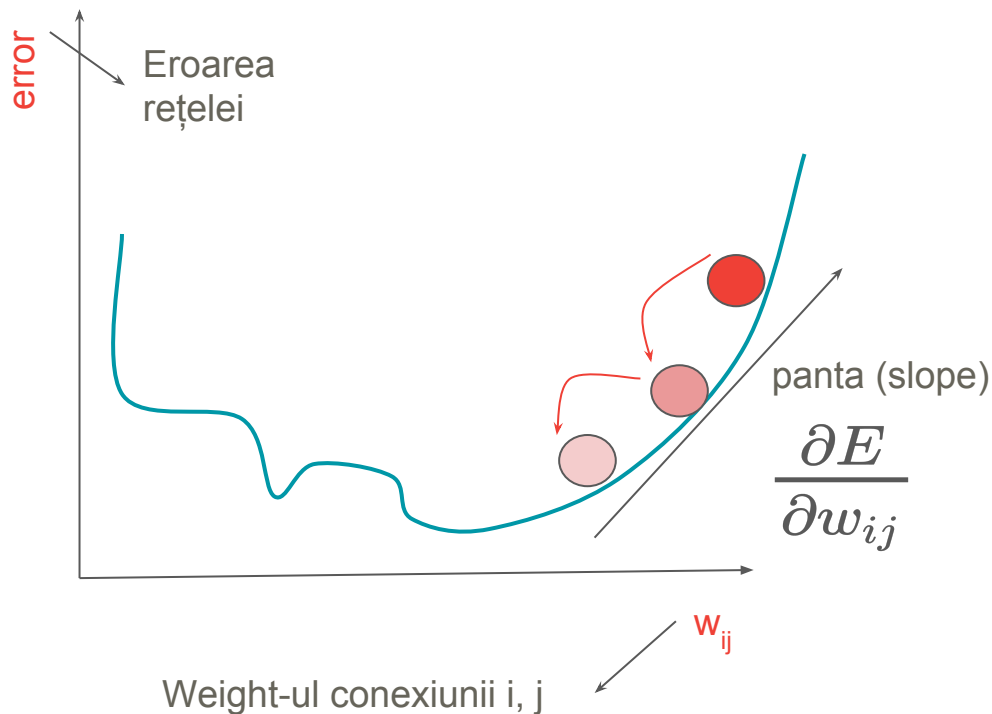
Panta dintre output și hidden layer:

Vom folosi formula asta la actualizarea weight-urilor.

$$\frac{\partial E}{\partial w_{ij}} = -e_j \text{sigmoid}(\sum_i w_{ij} o_i) (1 - \text{sigmoid}(\sum_i w_{ij} o_i)) o_i$$

Panta dintre hidden și input layer

Actualizarea weight-urilor (5)



Gradient Descent

- Metodă bună de găsit minimul unei funcții.
- Minimizarea lui E , prin rafinarea w_{ij} .
- Derivata lui E în raport cu w_{ij} ne spune cum se schimbă E odată cu modificarea weight-urilor.

Actualizarea weight-urilor (6)

$$w_{jk} = w_{jk} - \alpha \frac{\partial E}{\partial w_{jk}}$$

Learning rate - moderează dimensiunea salturilor către minim.

Ajustăm vechiul weight cu minus panta erorii.

Matricea schimbărilor de weights

$$\begin{pmatrix} \text{delta_}w_{11} & \text{delta_}w_{21} \\ \dots & \dots \\ \text{delta_}w_{12} & \text{delta_}w_{22} \\ \dots & \dots \\ \text{delta_}w_{13} & \text{delta_}w_{jk} \dots \end{pmatrix}$$

Leagă nodul j de pe un layer cu nodul k de pe următorul.

$$= \begin{pmatrix} E_1 S_1 (1 - S_1) \\ E_2 S_2 (1 - S_2) \\ E_k S_k (1 - S_k) \\ \dots \end{pmatrix} \begin{pmatrix} o_1 & o_2 & o_j & \dots \end{pmatrix}$$

Valori de pe stratul următor Valori de pe stratul anterior

$$\Delta w_{jk} = \alpha E_k o_k (1 - o_k) o_j^T$$

Sigmoidele înlocuite cu output-urile nodurilor o_k

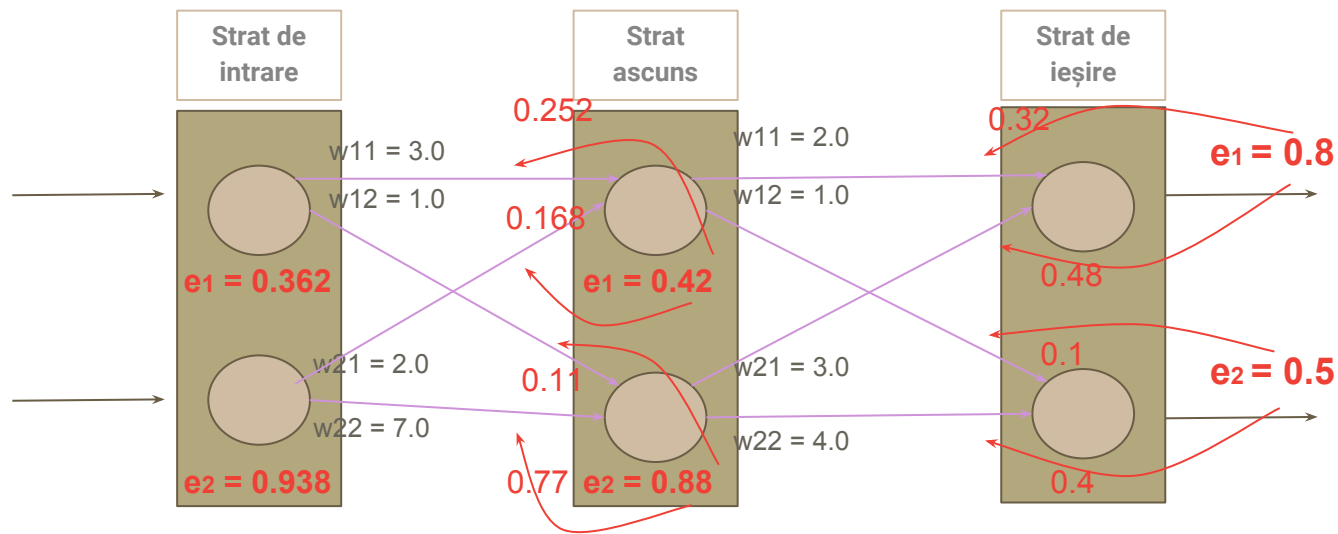
Actualizarea weight-urilor (7)

NOTE

- Rețea neuronală \Leftrightarrow funcție liniară de weight-uri
- Îmbunătățirea rețelei \Leftrightarrow minimizarea erorii \Leftrightarrow rafinare weights
- Îmbunătățim weight-urile iterativ prin descreșterea treptată a erorii folosind Gradient Descent.
- Fiecare pas spre minimul erorii este făcut în direcția în care panta este cea mai abruptă în jos!

Actualizarea weight-urilor (8)

Exemplu



Actualizarea weight-urilor (9)

Exemplu

Vrem să actualizăm $w_{11} = 2$ (dintre hidden și output layer).

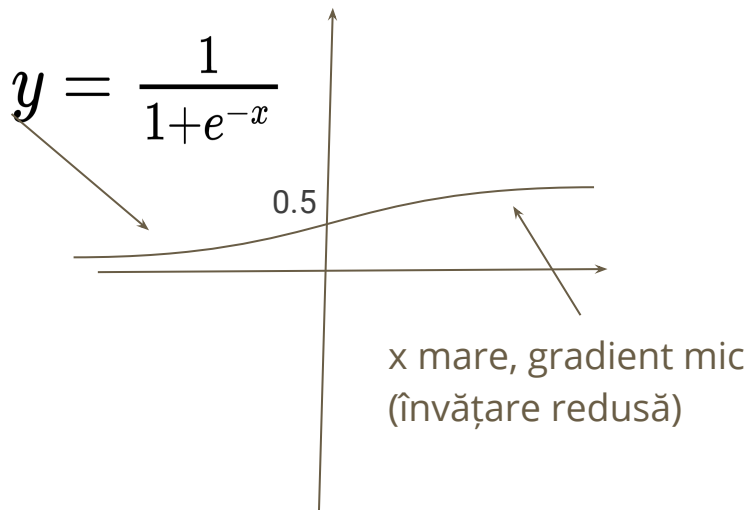
$$\frac{\partial E}{\partial w_{jk}} = -(t_k - o_k) \underbrace{\text{sigmoid}(\sum_j w_{jk} o_j)}_{= 2.0 * 0.4 + 3.0 * 0.4 = 2} (1 - \underbrace{\text{sigmoid}(\sum_j w_{jk} o_j)}_{= \text{sigmoid}(2) = 0.88}) \underbrace{o_j}_{= o_1 = 0.4}$$

$= e_1 = 0.8$

$= -0.0675$

$w_{11} = 2.0 + 0.1 * 0.0675 = 2.00675$ → Observăm o schimbare destul de mică, dar pe măsură ce iterăm, ea va crește.

Pregătirea datelor



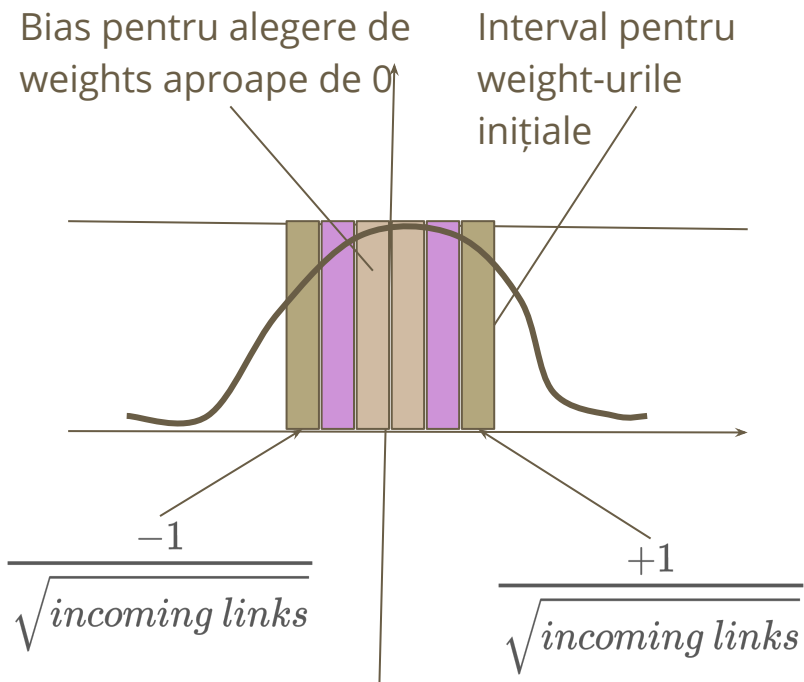
INPUT

- Valori foarte mari
=> aplatizare grafic
=> greu de anticipat încotro ne îndreptăm cu update-urile
=> **saturare**
- Țineți input-urile mici!
- Recomandare: input în **(0.0, 1.0)**

OUTPUT

- Grijă la valorile target!
- Funcția de activare trebuie să poată genera outputs în același interval cu target-urile!

Pregătirea datelor (2)



Weight-uri inițiale aleatorii

- Nu setăm toate weight-urile la aceeași valoare constantă. Cu atât mai puțin zero!
 - 0 să anuleze input-urile.
- Ideal e să facem sample dintr-o distribuție normală cu media aproape de zero și deviația standard $1/\sqrt{\text{incoming links}}$.

Part II

2NClub #1

Hands on:

Recunoașterea Cifrelor Scrise de Mână

