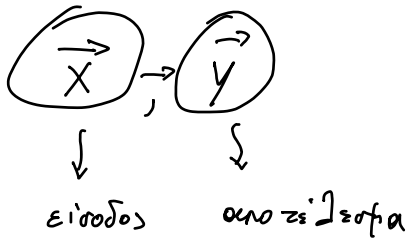


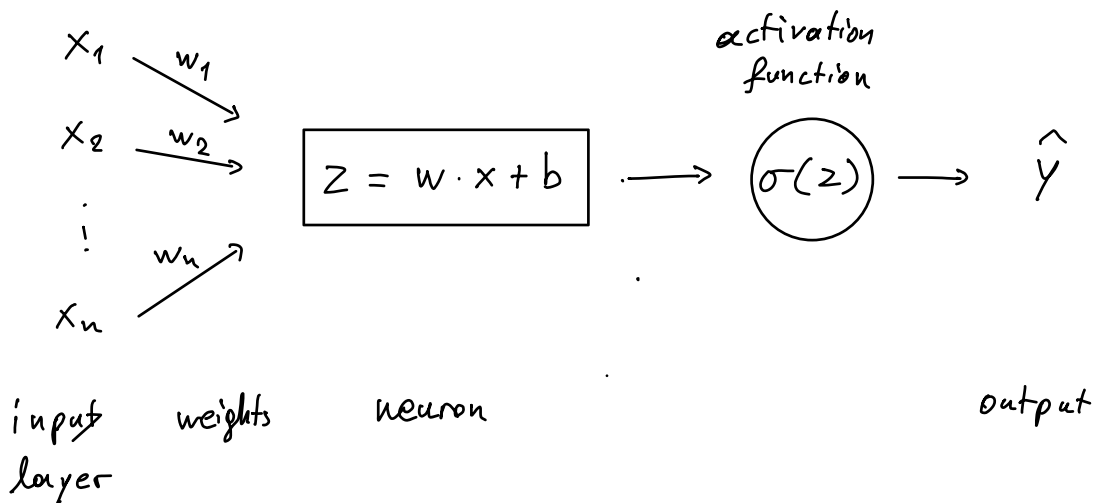
## Εισαγωγή στα Νευρωνικά Δίκτυα



$$\vec{X} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

$$\text{ή} \quad \vec{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ x_{21} & & \ddots & \vdots \\ \vdots & & & \vdots \\ x_{n1} & \dots & & x_{nm} \end{bmatrix}$$

$$\vec{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$



Βήμα 1.

Για κάθε input πολλαπλασιάζεται το  $x_i$  με το αντίστοιχο βάρος  $w_i$  και αθροίζεται ό)ες τις τιμές σε έναν αθροισμό.

Τα βάρη (weight)  $\rightarrow$  strength of connection

δηλ. όσο μεγαλύτερο είναι το  $w$  τόσο σημαντικότερο είναι το  $x_i$  για το  $y$

$w_1 > w_2 \Rightarrow x_1$  πιο σημαντικό από το  $x_2$  για την πρόβλεψη του  $y$

$$\Sigma = (x_1 \cdot w_1) + (x_2 \cdot w_2) + \dots + (x_n \cdot w_n)$$

$$\vec{x}, \vec{w} \quad i = 1 \dots n$$

$$\vec{x} \cdot \vec{w} = (x_1 \cdot w_1) + (x_2 \cdot w_2) + \dots + (x_n \cdot w_n)$$

$$\Sigma = \vec{x} \cdot \vec{w}$$

Βίας  $b$  : χρησιμοποιείται για να αναγλυτοποιήσει συνολικό shift των τιμών.

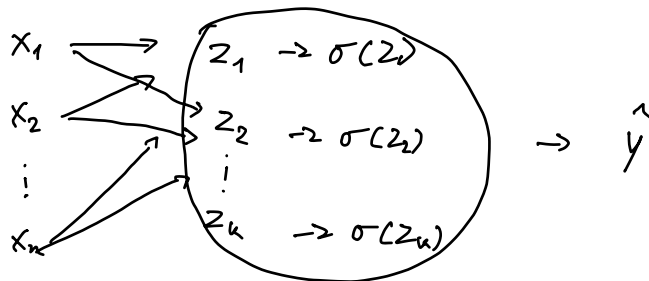
Βήμα 2

$$Z = \vec{X} \cdot \vec{W} + b \quad \Leftrightarrow y$$

Παρατήρηση: με τον παραπάνω ορισμό του  $z$  η αντιστροφή μεταξύ των inputs και των outputs είναι ανολύτως γραμμική

Εφαρμόζουμε λοιπόν ένα non-linear activation function.

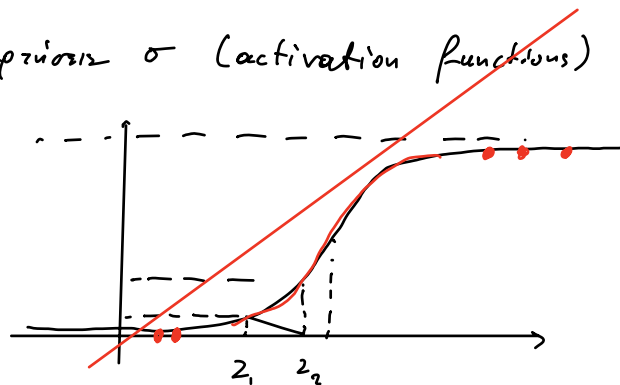
$$\hat{y} = \sigma(z) \quad \text{όπου } \sigma \text{ η γραμμική συνάρτηση}$$



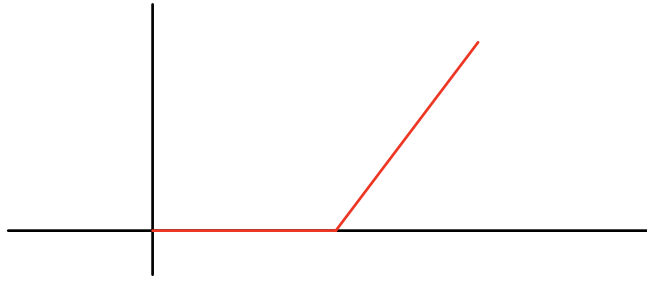
$\sigma$ : η γραμμική συνάρτηση

Συνήθεις συναρτήσεις  $\sigma$  (activation functions)

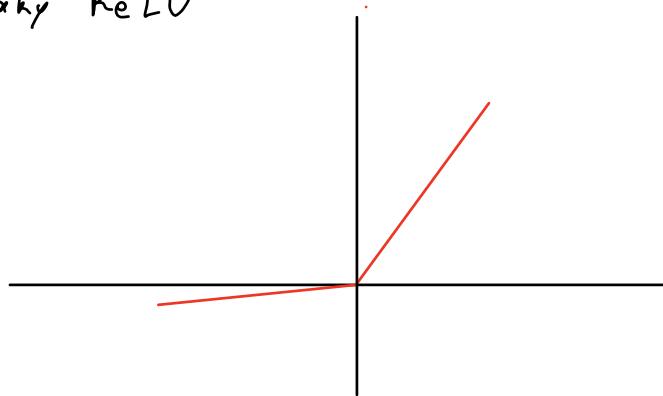
- Sigmoid



- ReLU: rectified linear unit activation function



• Leaky ReLU



activation function  $\sigma$

ας θυμίσω ότι ενδιόχουσε η σιγμοειδ

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Αλγόριθμος εναείδευσης

1ο βήμα για την εναείδευση είναι η ενιδογύ κάρηοιου  
loss function

Τι είναι το loss function;

Είναι συνάρτηση "δυναί > loss ενιδογύ's" ηε βάρη των

οποία κρινόμαστε αν το αποζητίωμα που έχουμε προβλέψει είναι κοντά ή όχι στην πραγματική τιμή.

Στις περισσότερες περιπτώσεις το loss function που χρησιμοποιούμε είναι η μέση τετραγωνική απόκλιση (mean squared error - MSE)

$$MSE_i = (y_i - \hat{y}_i)^2$$

Cost function : μέση τιμή του loss function  
(C)

$$C = MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} MSE_i$$

Έχουμε ξεκινήσει με random  $w$ . Βρίσκουμε με αυτά τα τυχαία βάρη πόσο ακριβής η πρόβλεψή μας από την πραγματική τιμή.

Πρέπει να μεταβάλλουμε τα  $w$  ώστε να πετύχουμε ένα συνολικά μικρότερο  $C$ .

$$\frac{\partial C}{\partial w_i} = \frac{\partial C}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z} \cdot \frac{\partial z}{\partial w}$$

$$\frac{\partial C}{\partial \hat{y}} = ? \quad \frac{\partial \hat{y}}{\partial z} = ? \quad \frac{\partial z}{\partial w} = ?$$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 2 \cdot \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$

expose  $\vec{y} = [y_1, y_2 \dots y_n]$ ,  $\vec{\hat{y}} = [\hat{y}_1, \hat{y}_2 \dots \hat{y}_n]$

$$\frac{\partial \mathcal{L}}{\partial \hat{y}} = \frac{2}{n} \sum (y - \hat{y})$$

$$\frac{\partial \hat{y}}{\partial z} = \frac{\partial}{\partial z} \hat{y} = \frac{\partial}{\partial z} \sigma(z) = \frac{\partial}{\partial z} \left( \frac{1}{1+e^{-z}} \right)$$

$$= - \frac{1}{(1+e^{-z})^2} (e^{-z})' = \frac{e^{-z}}{(1+e^{-z})^2} =$$

$$= \frac{1}{1+e^{-z}} \cdot \frac{e^{-z}}{1+e^{-z}} =$$

$$= \frac{1}{1+e^{-z}} \cdot \left( 1 - \frac{1}{1+e^{-z}} \right) =$$

$$= \sigma(z) \cdot (1 - \sigma(z))$$

$$\frac{\partial z}{\partial w_i} = \frac{\partial}{\partial w_i} \sum_{i=1}^n (x_i \cdot w_i + b) = x_i$$

Συνολικά έχουμε:

$$\frac{\partial \mathcal{L}}{\partial w_i} = \frac{2}{n} \cdot \sum (y - \hat{y}) \cdot \sigma(z) \cdot (1 - \sigma(z)) x_i$$

Βρίσκουμε το  $\frac{\partial \mathcal{L}}{\partial w_i}$  το οποίο θα είναι η βάση μας

για το gradient descent που θα χρησιμοποιούμε  
για το optimization μέχρι ένα όριο σύγκλισης

$$w'_i = w_i - \alpha \frac{\partial \mathcal{L}}{\partial w_i}$$

$\alpha$ : learning rate

όσο μεγαλύτερο είναι το  $\alpha$  τόσο πιο γρήγορα ο  
αλγόριθμος προσπαθεί να μείδσει

αν θέλουμε να κάνουμε optimize και στο bias  
χρησιμοποιούμε αντίστοιχα το gradient descent για το  
bias

$$b' = b - \alpha \frac{\partial \mathcal{L}}{\partial b}$$