

NEURAL NETWORKS

REGRESSION

Samples: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

$$x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

Find the best fit:

$$f(x_i; \underline{\theta}) \approx y_i$$

parameters.

Loss function: example

$$L(\underline{\theta}) = \sum_{i=1}^n (y_i - f(x_i; \underline{\theta}))^2$$

Goal: Find $\underline{\theta}$ that minimise $L(\underline{\theta})$

↳ Solve: $\nabla L(\underline{\theta}) = 0 \rightarrow$ not always solvable!

GRADIENT DESCENT

⊗ "Searching" for the minimum

- Choose $\underline{\theta}^{(0)}$
- $\underline{\theta}^{(h+1)} = \underline{\theta}^{(h)} - \beta \nabla L(\underline{\theta})$

NEURAL NETWORKS!

- "Complicated" function f
- $\underline{\theta} \in \mathbb{R}^D$ D - very large. (many params.)

Feedforward networks:

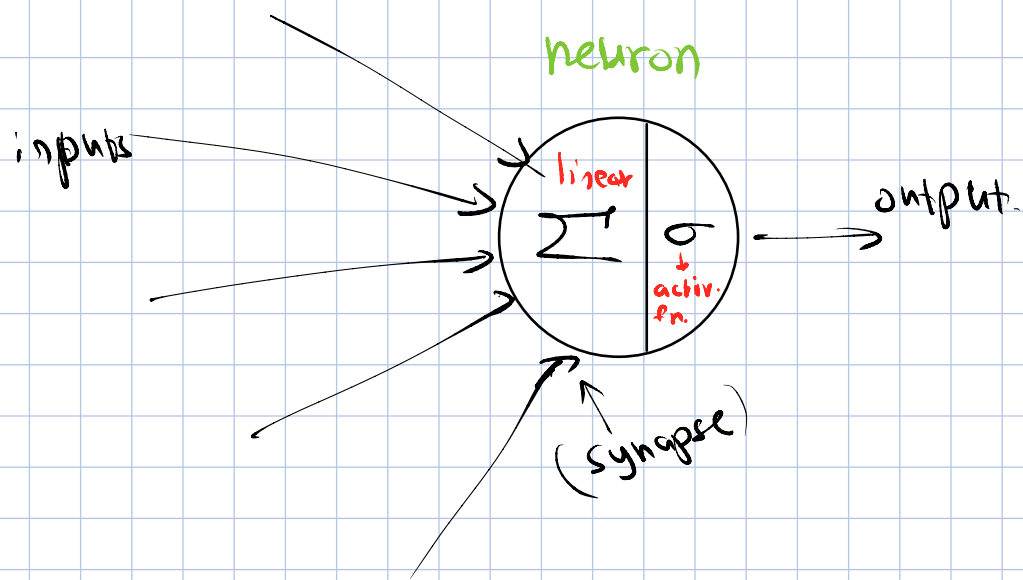
$$\underline{x}^{(0)} \in \mathbb{R}^d \rightarrow \text{activation function}$$

$$\underline{x}^{(1)} = \sigma^{(1)}(W^{(1)} \underline{x}^{(0)} + b^{(1)})$$

$$\underline{x}^{(2)} = \sigma^{(2)}(W^{(2)} \underline{x}^{(1)} + b^{(2)})$$

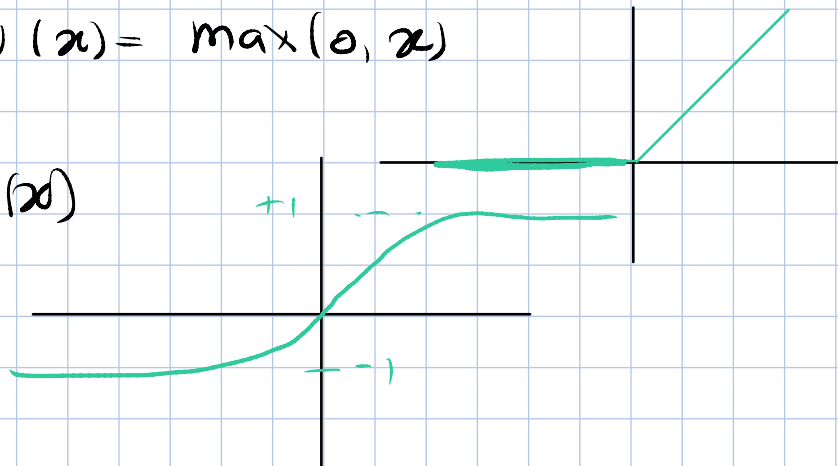
$$\underline{x}^{(L)} = \sigma^{(L)}(W^{(L)} \underline{x}^{(L-1)} + b^{(L)}) \rightarrow f(\underline{x}^{(0)}; \underline{\theta})$$

$\underline{\theta} = (W^{(1)}, \dots, W^{(L)}, b^{(1)}, \dots, b^{(L)})$

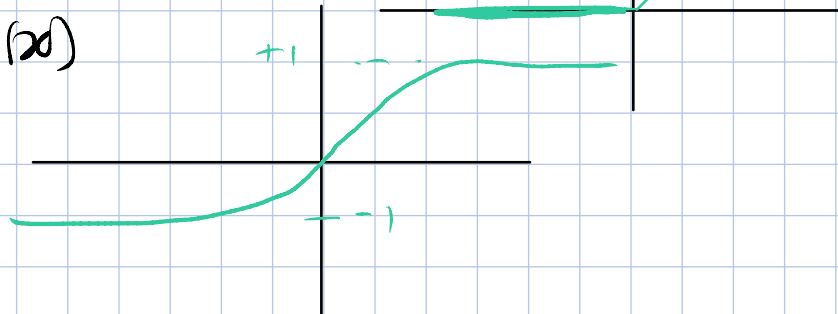


ACTIVATION FUNCTIONS:

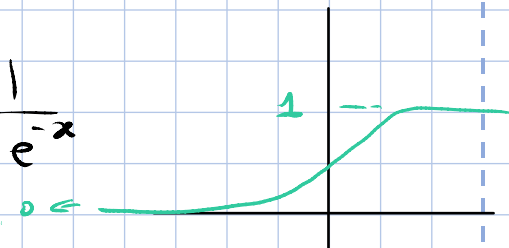
- $\text{ReLU}(x) = \max(0, x)$



- $\tanh(x)$



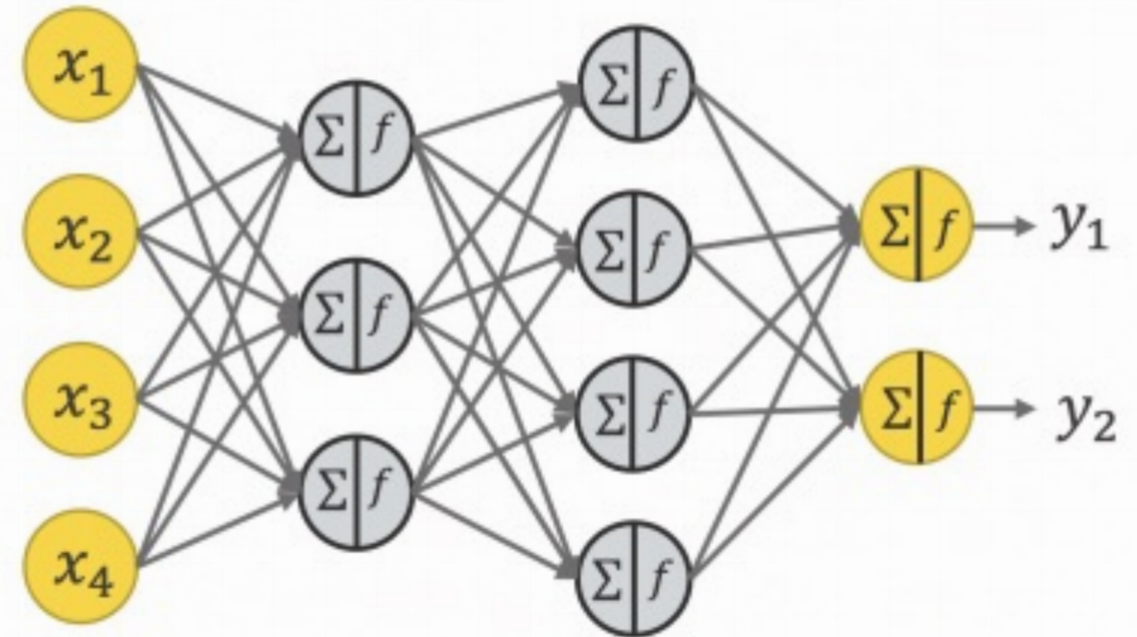
- logistic function - $\frac{1}{1 + e^{-x}}$



Input layer

Hidden layers

Output layer



TRAINING A NN:

Training set $(x_1, y_1), \dots, (x_n, y_n)$

Initialise $\theta^{(0)}$ (possibly random)

Iterate:

(1) Feed forward:

$$x_1, \dots, x_n \rightarrow f(x_1, \theta^{(k)}), \dots, f(x_n, \theta^{(k)})$$

(2) Compute loss:

$$L(\theta) = \sum_{i=1}^n (y_i - f(x_i; \theta^{(k)}))^2$$

(3) Back-propagation: $\nabla L(\theta^{(k)})$

(4) Update parameters:

$$\theta^{(k+1)} = \theta^{(k)} - \beta \nabla L(\theta^{(k)})$$

step-size

efficient way
to compute
gradients for NN.

STOCHASTIC GRADIENT DESCENT

$$L(\theta) = L(\theta; X) = \sum_{i=1}^n (y_i - f(x_i; \theta))^2 \text{ (MSE)}$$

$$\nabla L(\theta) = \sum_{i=1}^n \nabla L(\theta; x_i)$$

n -large \Rightarrow slow

Gradient descent:

$$\theta^{(k+1)} = \theta^{(k)} - \beta \nabla L(\theta; x)$$

Stochastic gradient descent:

• Choose $i, 1 \leq i \leq n$ at random

• Update: $\theta^{(k+1)} = \theta^{(k)} - \beta \nabla L(\theta^{(k)}; x_i)$

Mini-batch Gradient Descent:

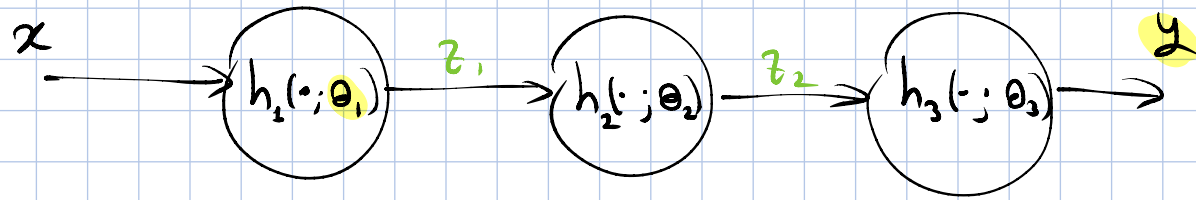
- Set m (commonly 32-256)
- Take a random subset $X^{(m)} = \{x_{i_1}, \dots, x_{i_m}\}$
- Update!

$$\underline{\theta}^{(k+1)} = \underline{\theta}^{(k)} - \beta \nabla \mathcal{L}(\underline{\theta}; X^{(m)})$$

Common practices:

- Shuffle X .
 - Split X into - $X_1^{(m)}, X_2^{(m)}, \dots, X_{\frac{N}{m}}^{(m)}$
 - Update $\underline{\theta}$ using $X_1^{(m)}$
" " " $X_2^{(m)}$
" " " \vdots
" " " $X_{\frac{N}{m}}^{(m)}$
 - repeat
- epoch

BACK-PROPAGATION



$$y = y(x; \theta_1, \theta_2, \theta_3) \quad (\text{ex. } y = \text{loss})$$

Q: What is $\frac{\partial y}{\partial \theta_1}$?

$$\begin{aligned} y &= h_3(z_2; \theta_3) = h_3(h_2(z_1; \theta_2); \theta_3) \\ &= h_3(h_2(h_1(x; \theta_1); \theta_2); \theta_3) \end{aligned}$$

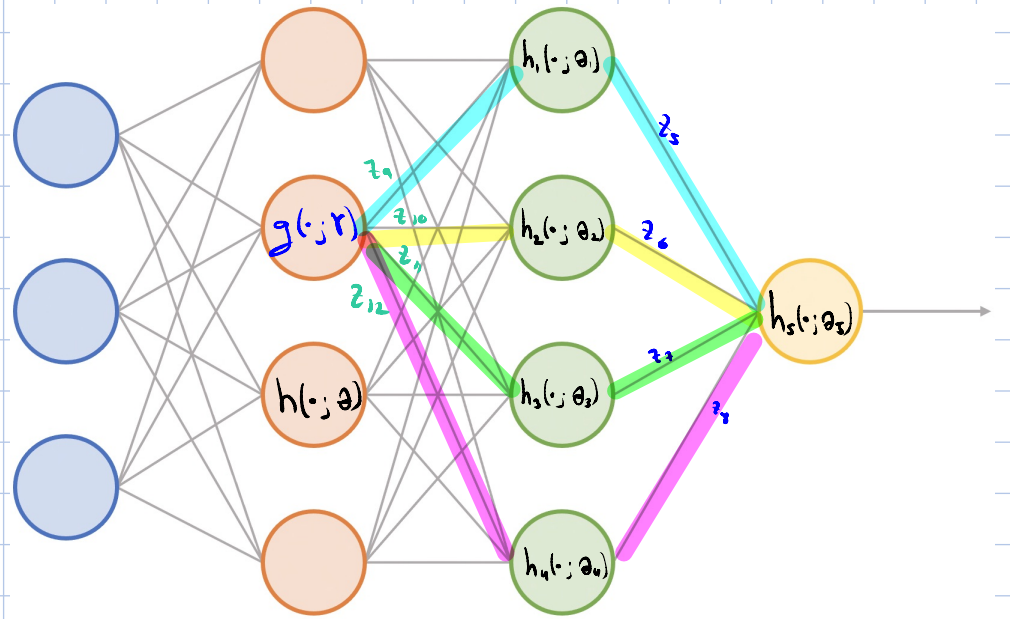
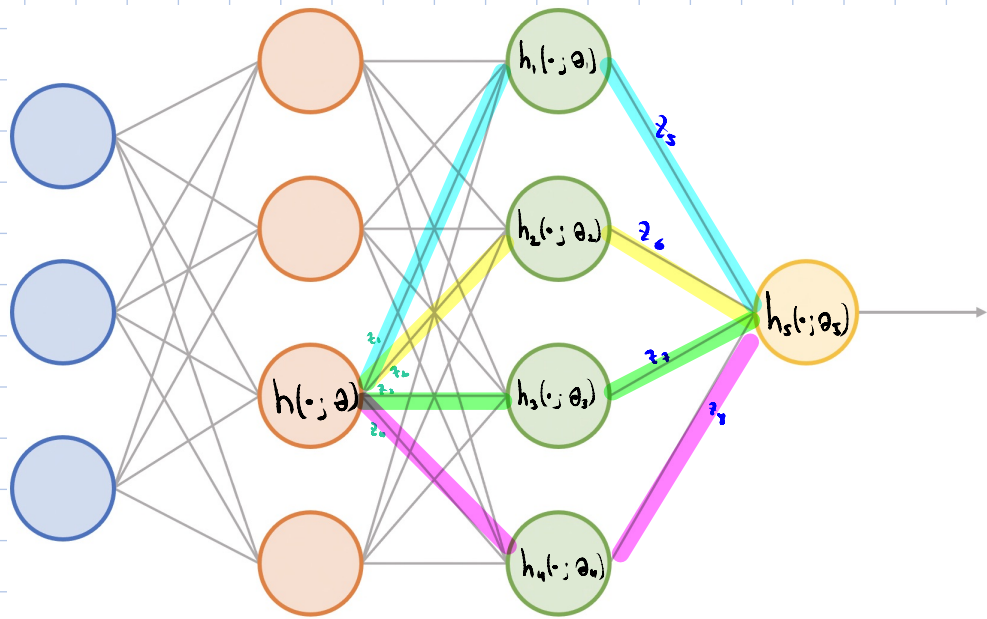
$$\begin{aligned} \frac{\partial y}{\partial \theta_1} &= \frac{\partial y}{\partial z_2} \cdot \frac{\partial z_2}{\partial \theta_1} = \frac{\partial h_3}{\partial z_2} \cdot \frac{\partial z_2}{\partial \theta_1} \\ &= \frac{\partial h_3}{\partial z_2} \cdot \left[\frac{\partial h_2}{\partial z_1} \cdot \frac{\partial z_1}{\partial \theta_1} \right] \end{aligned}$$

$$\frac{\partial y}{\partial \theta_1} = \frac{\partial h_3}{\partial z_2} \cdot \frac{\partial h_2}{\partial z_1} \cdot \frac{\partial h_1}{\partial \theta_1}$$

Feed-forward:

Take $x \rightarrow$ compute z_1, z_2, y .

$$\frac{\partial h_3}{\partial z_2}(z_2; \theta_3) \cdot \frac{\partial h_2}{\partial z_1}(z_1; \theta_2) \cdot \frac{\partial h_1}{\partial \theta_1}(x; \theta_1)$$



$$\frac{\partial y}{\partial \theta} = \frac{\partial h_5}{\partial z_5} \cdot \frac{\partial h_4}{\partial z_1} \cdot \frac{\partial h}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_6} \cdot \frac{\partial h_2}{\partial z_2} \cdot \frac{\partial h}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_7} \cdot \frac{\partial h_3}{\partial z_3} \cdot \frac{\partial h}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_8} \cdot \frac{\partial h_4}{\partial z_4} \cdot \frac{\partial h}{\partial \theta}$$

same values
computed once

$$\frac{\partial y}{\partial \theta} = \frac{\partial h_5}{\partial z_5} \cdot \frac{\partial h_4}{\partial z_1} \cdot \frac{\partial g}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_6} \cdot \frac{\partial h_2}{\partial z_{10}} \cdot \frac{\partial g}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_7} \cdot \frac{\partial h_3}{\partial z_{11}} \cdot \frac{\partial g}{\partial \theta}$$

$$+ \frac{\partial h_5}{\partial z_8} \cdot \frac{\partial h_4}{\partial z_{12}} \cdot \frac{\partial g}{\partial \theta}$$