

# Neural networks :

## Linear models :

$$y = f \left( \sum_{i=1}^d w_i \phi_i(x) \right) \begin{cases} \underline{f: \text{linear}}, \phi: \text{identity} \rightarrow \text{Linear regression} \\ \underline{f: \text{linear}}, \phi: \text{nonlinear} \rightarrow \text{with basis f.} \\ \underline{f: \text{logistic}}, \phi: \text{identity} \rightarrow \text{logistic regression} \end{cases}$$

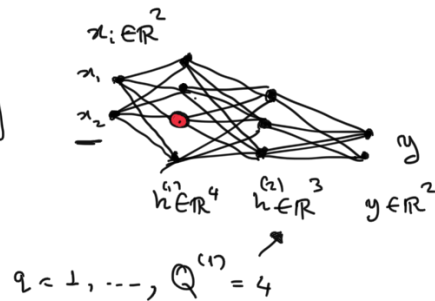
Goal : Make  $\phi_i(x)$  depend on trainable parameters.

## Multi-layer perceptron (MLP) :

$$y = w^T \phi(x; \theta)$$

$$\underline{h_q^{(1)}} = \underbrace{f \left( \sum_{i=1}^d w_{iq}^{(1)} + b_q^{(1)} \right)}_{\text{linear}}$$

activations



vectorize

$$\Rightarrow \mathcal{D} := \{ X, y \} \quad , \quad X = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{md} \end{bmatrix}, \quad y = \begin{bmatrix} y_1 & \dots & y_s \\ \vdots & & \vdots \\ y_{m1} & \dots & y_{ms} \end{bmatrix}$$

$n \times d \quad n \times s$

$$H^{(1)} = \begin{bmatrix} h_1^{(1)} \\ \vdots \\ h_q^{(1)} \end{bmatrix}_{n \times Q^{(1)}} = f^{(1)} \left( \underbrace{X W^{(1)}}_{n \times d \times Q^{(1)}} + \underbrace{b^{(1)}}_{1 \times Q^{(1)}} \right) \quad \begin{matrix} X: \text{input} \\ \downarrow \\ H^{(1)} \end{matrix}$$

$$H^{(2)} = f^{(2)} \left( \underbrace{H^{(1)} W^{(2)}}_{n \times Q^{(1)} \times Q^{(2)}} + \underbrace{b^{(2)}}_{1 \times Q^{(2)}} \right)$$

$\vdots$

$\vdots$

$$y = H^{(L)} W^{(L)} + b^{(L)}$$

$\downarrow$   
 $H^{(2)}$

$\downarrow$   
 $\vdots$   
 $H^{(L)}$

$\downarrow$   
 $u \dots$

output

Trainable parameters:

$$\theta := \left\{ \underbrace{W^{(1)}}_{2 \times 4}, \underbrace{b^{(1)}}_{1 \times 4}, \underbrace{W^{(2)}}_{4 \times 3}, \underbrace{b^{(2)}}_{1 \times 3}, \dots, \underbrace{W^{(L)}}_{3 \times 2}, \underbrace{b^{(L)}}_{1 \times 2} \right\}$$

Network hyper-parameters:

1.) Number of hidden layers:  $L$

2.) Number of neurons/dimensionality:  $\underline{Q^{(1)}}, \dots, \underline{Q^{(L)}}$

3.) Choice of activation functions: sigmoid, tanh, ReLU, etc.

Training:

$$y = f_{\theta}(x) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\Rightarrow p(y|x, \theta, \sigma^2) = \mathcal{N}(y | f_{\theta}(x), \sigma^2)$$

$$L(\theta, \sigma^2) := -\log p(y|x, \theta, \sigma^2) = \frac{1}{2\sigma^2} \sum_{i=1}^n [f_{\theta}(x_i) - y_i]^2 + \frac{n}{2} \log_{(2\pi\sigma^2)}$$

$\downarrow$

$$L(\theta) := \frac{1}{n} \sum_{i=1}^n [f_{\theta}(x_i) - y_i]^2 \quad (\text{mean squared loss})$$

(regression)

$$L(\theta) := -\sum_{i=1}^n y_i \log \sigma(f_{\theta}(x_i)) + (1 - y_i) \log (1 - \sigma(f_{\theta}(x_i)))$$

(classification)

$\theta^*$  ... S.O.D. ...  $n+1$

$$\theta = \arg \min_{\theta \in \Theta} \mathcal{L}(\theta) \implies \theta^{n+1} = \theta^n - \eta \nabla_{\theta} \mathcal{L}(\theta^n)$$

## Automatic differentiation

$$F: \mathbb{R}^d \rightarrow \mathbb{R} \quad (\text{e.g. } \mathcal{L}(\theta))$$

$$F: \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \rightarrow \square \quad \begin{matrix} y \in \mathbb{R} \\ x \in \mathbb{R}^d \end{matrix}$$

$$F = D \circ C \circ B \circ A$$

$$y = F(x) = D(C(B(A(x))))$$

$$\hookrightarrow y = D(c), \quad c = C(b), \quad b = B(a), \quad a = A(x)$$

Jacobian:  $\nabla_x F = \frac{dF}{dx} = \frac{dy}{dx} = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_d} \right]^T$

Chain rule:  $\nabla_x F = \begin{bmatrix} \frac{\partial y}{\partial c} & \frac{\partial c}{\partial b} \end{bmatrix} \cdot \begin{bmatrix} \frac{\partial b}{\partial a} & \frac{\partial a}{\partial x} \end{bmatrix} \xrightarrow{\text{reverse mode}}$

$\frac{\partial y}{\partial c} = D'(c)$   
 $\underline{1 \times d_c}$

$\frac{\partial c}{\partial b} = C'(b)$   

$\underline{d_c \times d_b}$

$\frac{\partial a}{\partial x} = A'(x)$   
 $\underline{d_a \times d}$

### • Forward accumulation:

$$\nabla_x F = \frac{\partial y}{\partial c} \left( \left( \frac{\partial c}{\partial b} \left( \frac{\partial b}{\partial a} \left( \frac{\partial a}{\partial x} \right) \right) \right) \right) \quad \text{use when}$$

$$\dim \{x\} \ll \dim \{y\}$$

- Reverse mode :

$$\nabla_x F = \left( \left( \frac{\partial y}{\partial c} \cdot \frac{\partial c}{\partial b} \right) \frac{\partial b}{\partial a} \right) \frac{\partial a}{\partial x}$$

use  
 $\dim \{x\} \gg \dim \{y\}$

- FW-AD  $\Rightarrow$  JVPs (Jacobian-vector product)

Constructs the Jacobian one column at a time

$$\nabla_x F: \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_d} \right] \underset{[0, \dots, 0]}{\vec{v}} = \frac{\partial y}{\partial c} \left( \frac{\partial c}{\partial b} \left( \left( \frac{\partial b}{\partial a} \left( \frac{\partial a}{\partial x} \vec{v} \right) \right) \right) \right)$$

- RV-AD  $\Rightarrow$  VJP, (Vector-Jacobian product)

Construct the Jacobian one row at a time:

$$\nabla_x F: \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_d} \right] = \left( \left( \left( \vec{v}^T \frac{\partial y}{\partial c} \right) \frac{\partial c}{\partial b} \right) \frac{\partial b}{\partial a} \right) \frac{\partial a}{\partial x}$$

$\mathbb{R}^d \rightarrow \mathbb{R}$