

# Neural Networks :

The linear models we covered so far are based on linear combinations of some fixed basis functions/features :

$$y = \sigma \left( \sum_{i=1}^d w_i \overbrace{\phi_i(x)}^{\text{fixed}} \right) + \varepsilon$$

$\downarrow$  model params       $\downarrow$  features/basis functions

$$\left\{ \begin{array}{l} \bullet \sigma: \text{linear}, \phi: \text{identity} \rightarrow \text{Linear reg. with 1 fun.} \\ \bullet \sigma: \text{linear}, \phi: \text{nonlinear} \rightarrow \text{Linear reg. with } d \text{ fun.} \\ \bullet \sigma: \text{logistic sigmoid}, \phi: \text{identity} \rightarrow \text{Logit reg.} \end{array} \right.$$

Goal : Allow the features/basis functions  $\phi_i(x)$  to depend on parameters, which can be "trained"/adjusted (along with the  $w_i$ ) from the observed data.

This idea leads to the basic Multi-Layer Perceptron (MLP) model which can be described as a series of functional transformations :

$$y = \mathbf{w}^T \underset{=}{\phi}(\mathbf{x}; \underset{=}{\theta}) + \varepsilon, \text{ with model params: } \Theta := \{\mathbf{w}, \theta\}$$

$\hookrightarrow$  Non-linear regression model

Forward pass for an MLP model :

$$h_q^{(1)} = \underset{\substack{\downarrow \\ \text{activation} \\ \text{in the 1st} \\ \text{layer}}}{f^{(1)}} \left( \underbrace{\sum_{i=1}^d w_{iq}^{(1)} x_i + b_q^{(1)}}_{\text{linear transform.}} \right), \quad x_i \in \mathbb{R}^d, \quad q=1, \dots, Q^{(1)}$$

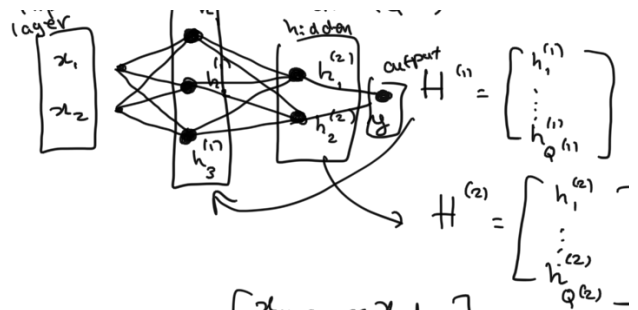
$h \in \mathbb{R}^{Q^{(1)}}$

$$\vdots$$

$$h_q^{(2)} = \underset{\substack{\text{hidden} \\ \text{in } \{1\}}}{f^{(2)}} \left( \sum_{i=1}^{Q^{(1)}} w_{iq}^{(2)} h_i^{(1)} + b_q^{(2)} \right)$$

inout

for e.g. assume:  $x \in \mathbb{R}^2$   
 $x = (x_1, x_2)$   
 hidden  $\dim(Q^{(1)}) = 3$   
 $\dim(Q^{(2)}) = 2$



In matrix-vector notation:  
(1st hidden layer):

$$H^{(1)}_{n \times Q^{(1)}} = f^{(1)}(X W^{(1)} + b^{(1)})$$

$$X = \begin{bmatrix} x_{11} & \dots & x_{1d} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nd} \end{bmatrix}_{n \times d}$$

$$W^{(1)}_{d \times Q^{(1)}} = \begin{bmatrix} w_{11} & \dots & w_{1Q^{(1)}} \\ \vdots & & \vdots \\ w_{d1} & \dots & w_{dQ^{(1)}} \end{bmatrix}, \quad b^{(1)}_{1 \times Q^{(1)}} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{Q^{(1)}} \end{bmatrix}$$

$$H^{(1)}_{n \times Q^{(1)}} = \begin{bmatrix} h_{11}^{(1)} & \dots & h_{1Q^{(1)}}^{(1)} \\ \vdots & & \vdots \\ h_{n1}^{(1)} & \dots & h_{nQ^{(1)}}^{(1)} \end{bmatrix}$$

2nd hidden layer:  $H^{(2)}_{n \times Q^{(2)}} = f^{(2)}(H^{(1)} W^{(2)} + b^{(2)})$

$W^{(2)}_{Q^{(1)} \times Q^{(2)}}$  weights,  $b^{(2)}_{1 \times Q^{(2)}}$  biases

Output layer:

$$Y_{n \times S} = f^{(0)}(H^{(L)} W^{(0)} + b^{(0)})$$

In practice the bias is added via broadcasting

$$b^{(0)}_{n \times S} = \begin{bmatrix} b_1^{(0)} & \dots & b_S^{(0)} \\ b_1^{(0)} & \dots & b_S^{(0)} \\ \vdots & & \vdots \\ b_1^{(0)} & \dots & b_S^{(0)} \end{bmatrix}$$

Setup: Given data:

$$D := \{(x_1, y_1), \dots, (x_n, y_n)\},$$

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

We have to make the following choices:

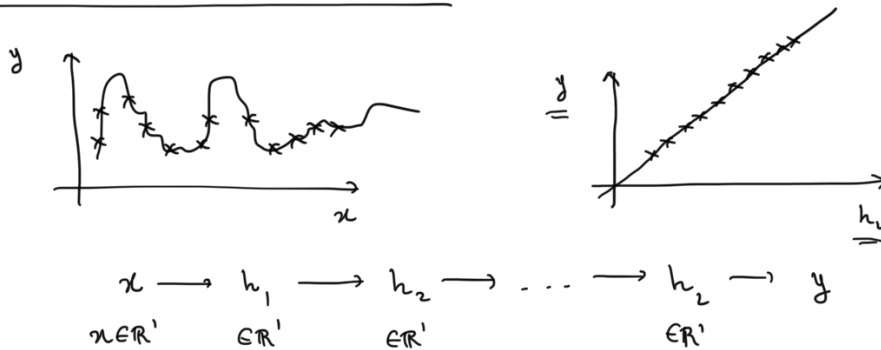
- 1.) # of hidden layers ( $L$ ) (i.e. how "deep" is our model)
- 2.) Dimensionality / # of "neurons":  $Q^{(1)}, \dots, Q^{(L)}$  (i.e. how wide is every hidden layer)
- 3.) Activation functions:  $f^{(1)}, f^{(2)}, \dots, f^{(L)}, f^{(0)}$

choices

typically the same depends on the task

#1,2,3 : define the MLP architecture.

Intuition on regression:



• Most common activation functions for the output layer:

- Linear :  $\underline{Y}_{n \times s} = f(H^{(L)} W^{(o)} + b^{(o)})$ , where  $f(x) = x$   
 $\underbrace{n \times Q^{(L)} \quad Q^{(L)}_{xs} \quad 1 \times s}_{\text{typical choice for regression}}$

- Logistic sigmoid :  $Y_{n \times 2} = f(H^{(L)} W^{(o)} + b^{(o)})$ ,  $f(x) = \frac{1}{1 + e^{-x}}$   
 typical choice for binary classification.

- Soft-max  $Y_{n \times q} = f(H^{(L)} W^{(o)} + b^{(o)})$   
 typical choice for multi-class classification,  $f(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^q \exp(x_j)}$