

I/ Basic MLP expression:

- The matrix $X \in \mathbb{R}^{n \times d}$: + a mini batch has n examples
+ each example has d inputs (features)
- For one-hidden-layer MLP whose hidden layer has h hidden units,
denoted by $H \in \mathbb{R}^{n \times h}$ the outputs of the hidden layer (hidden representations)
- H : hidden-layer variable or hidden variable
- \Rightarrow Hidden and output layers are fully connected
- \Rightarrow Hidden-layer weights $W_1 \in \mathbb{R}^{d \times h}$
Biases $b_1 \in \mathbb{R}^{1 \times h}$
- Output-layer weights $W_2 \in \mathbb{R}^{h \times q}$
Biases $b_2 \in \mathbb{R}^{1 \times q}$
- \Rightarrow Outputs $O \in \mathbb{R}^{n \times q}$

- Thus, this can be represented as:

$$H = XW_1 + b_1$$

$$O = HW_2 + b_2 \rightarrow O = (XW_1 + b_1)W_2 + b_2$$

$$\Rightarrow O = XW_1W_2 + b_1W_2 + b_2$$

$$\Rightarrow O = XW + b$$

- In addition, we will denote non-linearity activation function σ .

The outputs of activation functions are called **activations**

$$H = \sigma(XW_1 + b_1)$$

$$\Rightarrow O = HW_2 + b_2$$

- We can approximate many functions by using deeper (wider) networks

II/ Activation functions:

- Activation functions decide whether a neuron should be activated or not by calculating the weighted sum and further adding biases to it
- They are differentiable operators; while most of them are non-linearity