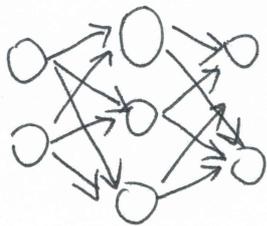


e.g. 2-layer Fully-connected NN (MLP):



input:  $x \in \mathbb{R}^D$

output:  $\hat{y} \in \mathbb{R}^C$

hidden layer:  $h \in \mathbb{R}^H$

activation function:  $g(\cdot), f(\cdot)$

parameters:  $\theta = \{ \underline{W}_I, \underline{b}_I, \underline{W}_O, \underline{b}_O \}$

feed-forward:

$$h = g(\underline{W}_I x + \underline{b}_I)$$

$$\hat{y} = f(\underline{W}_O h + \underline{b}_O)$$

note: I dropped transpose according to the advice I got from class to avoid confusion

parameter estimation:  $\theta = \min_{\theta} L(\hat{y}(x; \theta), y)$

gradient descent:  $\theta_{\text{new}} = \theta_{\text{old}} - \alpha \nabla_{\theta} L$

where  $\alpha$  is the learning rate

$\nabla_{\theta} L$  is the gradient of the loss w.r.t.  $\theta$

backpropagation:

Let  $\underline{\delta}_0 = \frac{\partial L}{\partial \underline{z}_0}$  for  $\underline{z}_0 = \underline{W}_0 \underline{h} + \underline{b}_0$  (pre-activation output)

$$\underline{\delta}_0 = \frac{\partial L}{\partial \hat{y}} \odot f'(\underline{z}_0) \quad \leftarrow \text{chain rule}$$

$$\therefore \nabla_{\underline{W}_0} L = \frac{\partial L}{\partial (\underline{W}_0)_{i,c}} = \frac{\partial L}{\partial (\underline{z}_0)_c} \cdot \frac{\partial (\underline{z}_0)_c}{\partial (\underline{W}_0)_{i,c}} = (\underline{\delta}_0)_c \cdot \underline{h}_i$$

$\Downarrow$

$$\nabla_{\underline{W}_0} L = \underline{h} \underline{\delta}_0^T$$

$$\nabla_{\underline{b}_0} L = \underline{\delta}_0$$

Let  $\underline{\delta}_I = \frac{\partial L}{\partial \underline{z}_I}$  where  $\underline{z}_I = \underline{W}_I \underline{x} + \underline{b}_I$

$$\underline{\delta}_I = \frac{\partial L}{\partial \underline{h}} \odot g'(\underline{z}_I) = \left( \frac{\partial \underline{z}_0}{\partial \underline{h}} \cdot \frac{\partial L}{\partial \underline{z}_0} \right) \odot g'(\underline{z}_I)$$

$$= \underline{W}_0 \underline{\delta}_0 \odot g'(\underline{z}_I)$$

$$\therefore \nabla_{\underline{W}_I} L = \underline{x} \underline{\delta}_I^T, \quad \nabla_{\underline{b}_I} L = \underline{\delta}_I$$