

I) forward pass

- Input \vec{x}
- Compute $\vec{\text{net}}_z = \begin{pmatrix} \text{net}_1 \\ \text{net}_2 \\ \vdots \\ \text{net}_m \end{pmatrix}$ $\vec{\text{net}}_z = \underline{W}_h^T \vec{x} + \underline{\vec{b}}_h$
- $\vec{z} = f^h(\vec{\text{net}}_z) = \begin{pmatrix} f^h(\text{net}_1) \\ f^h(\text{net}_2) \\ \vdots \\ f^h(\text{net}_m) \end{pmatrix}$ $\vec{z} = f^h(\vec{\text{net}}_z) = \underline{f}^h(\underline{W}_h^T \vec{x} + \underline{\vec{b}}_h)$
- \vec{z} is the "input" to the \vec{o} layer
 $\vec{\text{net}}_o = \begin{pmatrix} \text{net}_{o1} \\ \text{net}_{o2} \\ \vdots \\ \text{net}_{ok} \end{pmatrix}$ $\vec{\text{net}}_o = \underline{W}_o^T \vec{z} + \underline{\vec{b}}_o$
- $\vec{o} = f^o(\vec{\text{net}}_o) = \underline{f}^o(\underline{W}_o^T \vec{z} + \underline{\vec{b}}_o)$

II) Back propagation of error

$$\vec{o} = \underline{f}^o(\underline{W}_o^T (\underline{f}^h(\underline{W}_h^T \vec{x} + \underline{\vec{b}}_h)) + \underline{\vec{b}}_o)$$

$\vec{\text{net}}_z$
 $\vec{\text{net}}_o$

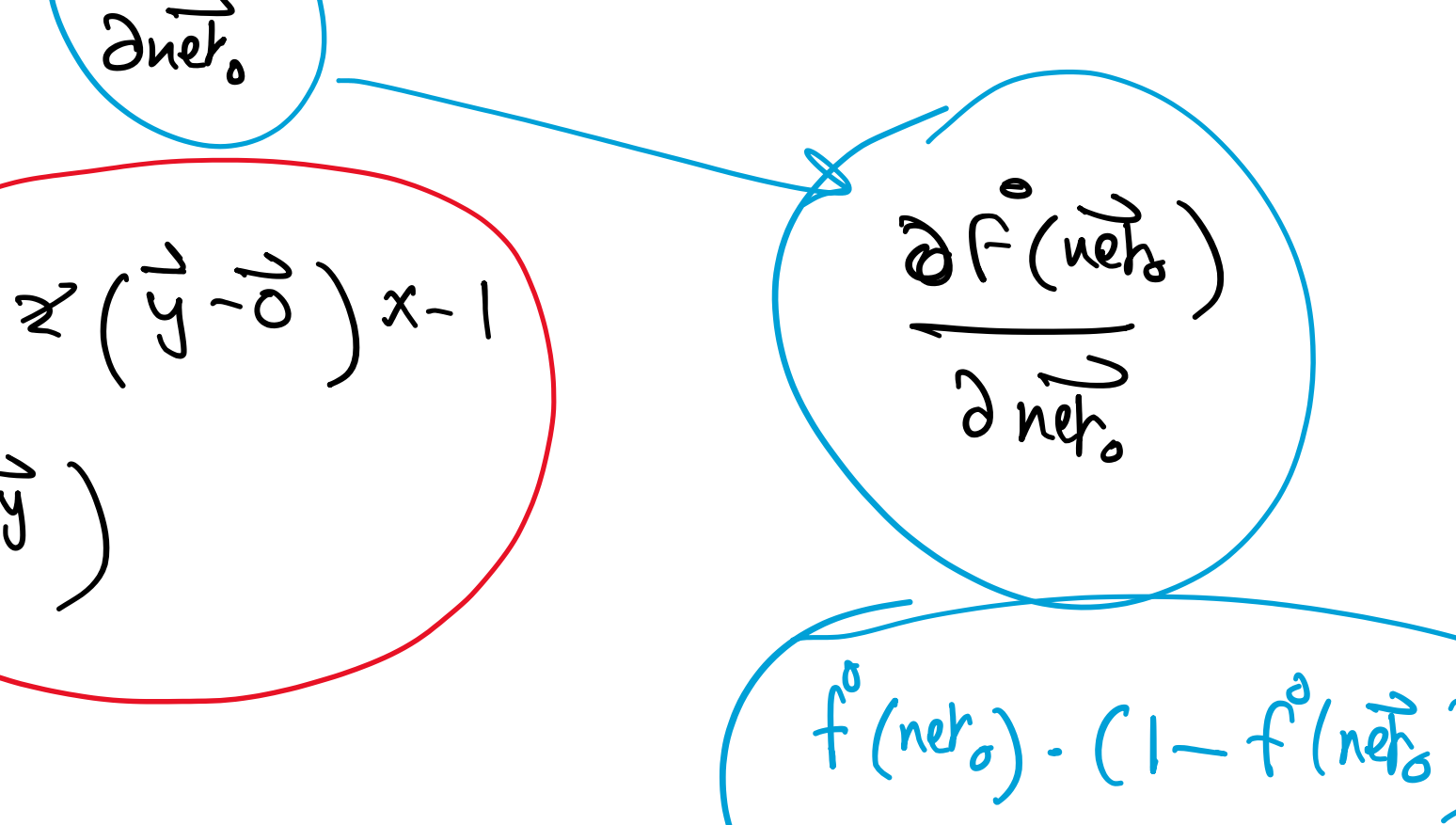
$\nabla_{W_h} = \frac{\partial \varepsilon}{\partial W_h}$ matrix $d \times m$	$\nabla_{b_h} = \frac{\partial \varepsilon}{\partial b_h}$ vector m	$\nabla_{W_o} = \frac{\partial \varepsilon}{\partial W_o}$ matrix $m \times k$	$\nabla_{b_o} = \frac{\partial \varepsilon}{\partial b_o}$ vector k
gradient descent $W_h = W_h - \eta \cdot \nabla_{W_h}$ $b_h = b_h - \eta \cdot \nabla_{b_h}$		$W_o = W_o - \eta \cdot \nabla_{W_o}$ $b_o = b_o - \eta \cdot \nabla_{b_o}$	

$$\vec{o} = \underline{f}^o(\underline{W}_o^T \vec{z} + \underline{\vec{b}}_o) = \underline{f}^o(\vec{\text{net}}_o)$$

$$\varepsilon = \frac{1}{2} \|\vec{y} - \vec{o}\|^2$$

$$\vec{\delta}_o \equiv \text{net gradient vector for output} = \frac{\partial \varepsilon}{\partial \vec{\text{net}}_o}$$

$$\vec{\delta}_h \equiv \text{net gradient for hidden} = \frac{\partial \varepsilon}{\partial \vec{\text{net}}_z}$$



$$\varepsilon = \frac{1}{2} \|\vec{y} - \vec{o}\|^2 = \frac{1}{2} \|\vec{y} - \underline{f}^o(\vec{\text{net}}_o)\|^2$$

$$\vec{o} = \underline{f}^o(\vec{\text{net}}_o)$$

$$\frac{\partial \varepsilon}{\partial \vec{\text{net}}_o} = \frac{\partial \varepsilon}{\partial \vec{o}} \times \frac{\partial \vec{o}}{\partial \vec{\text{net}}_o}$$

e.g. f^o to sigmoid

$$= \frac{1}{2} \cdot 2(\vec{y} - \vec{o}) \cdot 1 = (\vec{o} - \vec{y})$$

$$\frac{\partial \vec{f}^o(\vec{\text{net}}_o)}{\partial \vec{\text{net}}_o}$$

$$\vec{o} \odot (1 - \vec{o})$$

$$\vec{\delta}_o = \frac{\partial \varepsilon}{\partial \vec{\text{net}}_o} = (\vec{o} - \vec{y}) \odot \vec{o} \odot (1 - \vec{o})$$

element-wise multiplication

e.g. f^o to tanh
 $\frac{\partial f^o(\vec{\text{net}}_o)}{\partial \vec{\text{net}}_o} = 1 - \vec{\text{net}}_o^2$
 element-wise

$$\frac{\partial \varepsilon}{\partial \vec{\text{net}}_o} = \begin{cases} 1 & \text{if } \vec{\text{net}}_o > 0 \\ 0 & \text{otherwise} \end{cases}$$

ReLU

$$\frac{\partial \varepsilon}{\partial \vec{\text{net}}_z} = \frac{\partial \varepsilon}{\partial \vec{o}} \cdot \frac{\partial \vec{o}}{\partial \vec{\text{net}}_o} \cdot \frac{\partial \vec{\text{net}}_o}{\partial \vec{z}} \cdot \frac{\partial \vec{z}}{\partial \vec{\text{net}}_z}$$

$$\frac{\partial \varepsilon}{\partial \vec{\text{net}}_z} = \delta \cdot \frac{\partial \vec{\text{net}}_o}{\partial \vec{z}} \cdot \frac{\partial \vec{z}}{\partial \vec{\text{net}}_z}$$

Sigmoid

$$\vec{\delta}_h = \underline{W}_o^T \cdot \vec{\delta}_o \odot \vec{z} \odot (1 - \vec{z})$$

$$\vec{\delta}_h = \underline{W}_o^T \cdot \vec{\delta}_o \odot \vec{z} \odot (1 - \vec{z})$$

$$\nabla_{W_o} = \frac{\partial \varepsilon}{\partial W_o} = \frac{\partial \varepsilon}{\partial \vec{\text{net}}_o} \cdot \frac{\partial \vec{\text{net}}_o}{\partial W_o} = \vec{z} \odot \vec{\delta}_o^T$$

Outer product

$$\nabla_{b_o} = \vec{\delta}_o$$

$$\nabla_{W_h} = \vec{x} \cdot \vec{\delta}_h^T$$

$$\nabla_{b_h} = \vec{\delta}_h$$

$$\vec{\delta}_h = (\underline{W}_o^T \cdot \vec{\delta}_o) \odot \vec{z} \odot (1 - \vec{z})$$

$$\vec{\delta}_o = \frac{\partial \varepsilon}{\partial \vec{o}} \odot \frac{\partial \vec{o}}{\partial \vec{\text{net}}_o}$$

$$\nabla_{W_o} = \vec{z} \cdot \vec{\delta}_o^T$$

$$\nabla_{W_h} = \vec{x} \cdot \vec{\delta}_h^T$$

$$\nabla_{b_h} = \vec{\delta}_h$$

deep MLP

h : # of hidden layers (info bottleneck)



W_i : weight matrix is between $l=i$ and $l=i+1 \in \mathbb{R}^{m_i \times m_{i+1}}$
 \vec{b}_i : bias vector for $l=i+1$

I) forward pass

$$\vec{z}^0 = \vec{x}$$

for $i = 0, \dots, h$
 $\vec{\text{net}}_{i+1} = \underline{W}_i^T \vec{z}^i + \underline{\vec{b}}_i$ $(\vec{\text{net}}_o = \underline{W}_o^T \vec{z} + \underline{\vec{b}}_o)$
 $\vec{z}^{i+1} = f^{i+1}(\vec{\text{net}}_{i+1})$

II) error

$$\vec{o} = \vec{z}^{h+1} \in \mathbb{R}^k$$

$$\varepsilon = \frac{1}{2} \|\vec{y} - \vec{o}\|^2 \quad \text{or} \quad \varepsilon = C\varepsilon = -\sum_{i=1}^k y_i \log o_i$$

III) backprop

$$\vec{\delta}^{h+1} = \frac{\partial \varepsilon}{\partial \vec{o}} \odot \frac{\partial \vec{o}}{\partial \vec{\text{net}}_o}$$

for $l = h, h-1, \dots, 1$ (reverse)
 $\vec{\delta}^l = \frac{\partial \varepsilon}{\partial \vec{\text{net}}_l} \odot \frac{\partial \vec{\text{net}}_l}{\partial \vec{z}^l}$ $\frac{\partial \varepsilon}{\partial \vec{\text{net}}_l} = f'^l \odot \frac{\partial \varepsilon}{\partial \vec{\text{net}}_{l+1}}$

IV gradient for parameters

for $l = 0, 1, \dots, h$

$$\nabla_{W_l} = \vec{z}^l \cdot (\vec{\delta}^{l+1})^T$$

outer product of two vectors

$$\nabla_{b_l} = \vec{\delta}^{l+1}$$

V gradient descent

for $l = 0, 1, \dots, h$

$$W_l = W_l - \eta \cdot \nabla_{W_l}$$

$$b_l = b_l - \eta \cdot \nabla_{b_l}$$

repeat I - V $\forall x \in \text{Domain}$

of repeats \equiv epochs

for $e = 1, \dots$ epochs

- I
- II
- III
- IV
- V

Stopping criteria

- fixed # of epochs
- monitor the loss

derivatives

Sigmoid $f(z) = \frac{1}{1 + e^{-z}}$ $f' = \vec{z} \odot (1 - \vec{z})$

tanh $f(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ $f' = 1 - \vec{z}^2$

ReLU $f(z) = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$ $f' = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{else} \end{cases}$

Softmax $(\text{net}_1, \text{net}_2, \dots, \text{net}_k) = \frac{e^{\text{net}_i}}{\sum_{j=1}^k e^{\text{net}_j}}$

$$\frac{\partial \varepsilon}{\partial o} \cdot \frac{\partial f^o}{\partial \text{net}_i}$$

for softmax

$$\begin{pmatrix} \frac{\partial o_1}{\partial \text{net}_1} & \frac{\partial o_1}{\partial \text{net}_2} & \dots & \frac{\partial o_1}{\partial \text{net}_k} \\ \frac{\partial o_2}{\partial \text{net}_1} & \frac{\partial o_2}{\partial \text{net}_2} & \dots & \frac{\partial o_2}{\partial \text{net}_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial o_k}{\partial \text{net}_1} & \frac{\partial o_k}{\partial \text{net}_2} & \dots & \frac{\partial o_k}{\partial \text{net}_k} \end{pmatrix}$$