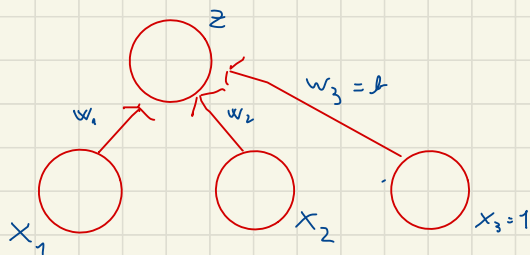


Backpropagation

GRADIENT DESCENT & LINEAR REGRESSION



Compute z in forward propagation:

$$z^{(n)} = x_1^{(n)} w_1 + x_2^{(n)} w_2 + b = \vec{w} \cdot \vec{x}^{(n)}$$

Compute $J(w) = \sum_{n \in \text{Train}} (t^{(n)} - z^{(n)})^2$

Let's compute the gradient for w_1

$$\frac{\partial J}{\partial w_1} = \frac{\partial}{\partial w_1} \frac{1}{2 N_{\text{Train}}} \sum_{n \in \text{Train}} (t^{(n)} - z^{(n)})^2 \quad \text{Sum rule}$$

$$= \frac{1}{2 N_{\text{Train}}} \sum_{n \in \text{Train}} \frac{\partial}{\partial w_1} (t^{(n)} - z^{(n)})^2 \quad \text{Power rule}$$

$$= \frac{1}{2 N_{\text{Train}}} \sum_{n \in \text{Train}} 2 (t^{(n)} - z^{(n)}) \frac{\partial}{\partial w_1} (t^{(n)} - z^{(n)}) \quad \frac{\partial}{\partial w_1} t^{(n)} = 0 \quad \frac{\partial}{\partial w_1} z^{(n)} = \frac{\partial}{\partial w_1} (x_1^{(n)} w_1 + x_2^{(n)} w_2 + b) = x_1^{(n)}$$

$$= \boxed{-} \frac{1}{N_{\text{Train}}} \sum_{n \in \text{Train}} (t^{(n)} - z^{(n)}) x_1^{(n)}$$

Similar derivation for w_2 and b .

$$\nabla_{\vec{w}} J = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial w_3} \right]$$

MATRIX NOTATION

$$\nabla_{\vec{w}} J = \frac{1}{N_{\text{Train}}} \sum_{n \in \text{Train}} (t^{(n)} - z^{(n)}) \vec{x}^{(n)T}$$

FOR GRADIENT ROW NOTATION: $\frac{\partial J}{\partial \vec{w}} = \left[\frac{\partial J}{\partial w_1}, \frac{\partial J}{\partial w_2}, \frac{\partial J}{\partial w_3} \right]$

Let's compute a numerical example.

$$T_{\text{Train}} = \left\{ \begin{bmatrix} 1 \\ 0 \end{bmatrix}, z \right\} \leftarrow \text{SIMPLE EXAMPLE. EXERCISE: TRY WITH 2!}$$

FORWARD: $z = 1 \cdot 1 + 0 \cdot 0 + 2 \cdot 1 = 3$

$$\frac{\partial J}{\partial w_1} = -(t - z) x_1 = (2 - 3) \cdot 1 = -1$$

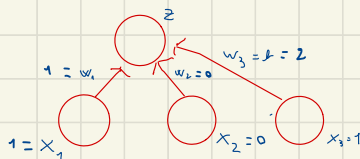
$$\frac{\partial J}{\partial w_2} = -(t - z) x_2 = (2 - 3) \cdot 0 = 0$$

$$\frac{\partial J}{\partial w_3} = (2 - 3) \cdot 1 = -1$$

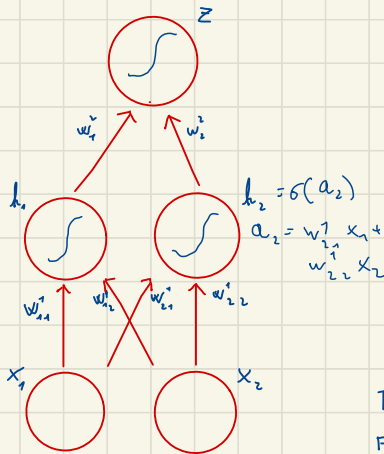
MATRIX NOTATION

$$\nabla_{\vec{w}} J = -(t - z) \vec{x} = (2 - 3) \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

We can update $\vec{w} \leftarrow \begin{bmatrix} 1 \\ 0 \end{bmatrix} - \epsilon \begin{bmatrix} -1 \\ 0 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 0 \end{bmatrix}$ NOTE: different learning rates would require more iterations!



BACKPROPAGATION



RECALL:

$$\sigma(\text{net}) = \frac{1}{1 + e^{-\text{net}}} \quad \sigma'(\text{net}) = \frac{d \sigma(\text{net})}{d \text{net}} = \sigma(\text{net})(1 - \sigma(\text{net}))$$

CHAIN RULE:

$$\frac{d f(g(x))}{d x} = \frac{d f(g(x))}{d g(x)} \cdot \frac{d g(x)}{d x}$$

$$\vec{a} = \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$T_n = \{(x^{(n)}, c^{(n)}) \dots (x^{(M_n)}, c^{(M_n)})\}$$

FORWARD PROP: COMPUTES \vec{z} and the hidden representation \vec{h}

Let's consider MSE loss: $J = \frac{1}{2} \frac{1}{N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} (c^{(p)} - z^{(p)})^2$

Let's compute the gradient w.r.t. the output weights \vec{w}_1, \vec{w}_2 is analogous

$$\frac{\partial J}{\partial w_1} = \frac{\partial}{\partial w_1} \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} (c^{(p)} - z^{(p)})^2 \quad \text{derivative is linear. Derivative of sum is the sum of the derivatives.}$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} \frac{\partial}{\partial w_1} (c^{(p)} - z^{(p)})^2 \quad \text{power rule: } \frac{d}{dx} f(x)^2 = 2[f(x)]^{1 \times 2-1} f'(x)$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} 2(c^{(p)} - z^{(p)}) \frac{\partial}{\partial w_1} (c^{(p)} - z^{(p)}) \quad z^{(p)} = \sigma(w_1^1 h_1^{(p)} + w_2^1 h_2^{(p)}) = \sigma(\vec{w}_1 \cdot \vec{h}^{(p)})$$

$$= \frac{1}{N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} (c^{(p)} - z^{(p)}) \frac{\partial}{\partial w_1} (c^{(p)} - \sigma(\vec{w}_1 \cdot \vec{h}^{(p)})) \quad \text{sum rule: } \frac{\partial}{\partial w_1} c^{(p)} = 0 \quad \text{chain rule: } \frac{\partial}{\partial w_1} \sigma(\vec{w}_1 \cdot \vec{h}^{(p)}) = \sigma'(\vec{w}_1 \cdot \vec{h}^{(p)}) \cdot \vec{h}_1^{(p)}$$

$$= \frac{1}{N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} \underbrace{(c^{(p)} - z^{(p)})}_{\nabla_{z^{(p)}} J} \underbrace{\sigma'(\vec{w}_1 \cdot \vec{h}^{(p)})}_{\text{BACKWARD PASS}} \underbrace{\vec{h}_1^{(p)}}_{\text{FORWARD PASS}} \quad \text{For sigmoid: } \sigma'(\vec{w}_1 \cdot \vec{h}^{(p)}) = \sigma(\vec{w}_1 \cdot \vec{h}^{(p)}) (1 - \sigma(\vec{w}_1 \cdot \vec{h}^{(p)}))$$

Let's now compute the gradient for the hidden units, w_{11}, w_{12}

$$\frac{\partial J}{\partial w_{11}} = \frac{\partial}{\partial w_{11}} \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} (c^{(p)} - z^{(p)})^2 \quad \text{Sum Rule}$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} \frac{\partial}{\partial w_{11}} (c^{(p)} - z^{(p)})^2 \quad \text{Power Rule and } \frac{\partial}{\partial w_{11}} c^{(p)} = 0 \quad \text{and Sum Rule}$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} 2(c^{(p)} - z^{(p)}) \frac{\partial}{\partial w_{11}} (-z^{(p)}) \quad z^{(p)} = \sigma(\vec{w}_1 \cdot \vec{h}^{(p)}) \quad \text{and Chain Rule}$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} 2(c^{(p)} - z^{(p)}) \sigma'(\vec{w}_1 \cdot \vec{h}^{(p)}) \frac{\partial}{\partial w_{11}} \vec{w}_1 \cdot \vec{h}^{(p)} \quad \vec{w}_1 \cdot \vec{h}^{(p)} = w_{11}^1 h_1^{(p)} + w_{12}^1 h_2^{(p)}$$

$$= \frac{1}{2 N_{\text{tr}}} \sum_{p \in T_{\text{tr}}} 2(c^{(p)} - z^{(p)}) \sigma'(\vec{w}_1 \cdot \vec{h}^{(p)}) \frac{\partial}{\partial w_{11}} w_{11}^1 h_1^{(p)} + \frac{\partial}{\partial w_{11}} w_{12}^1 h_2^{(p)} \quad h_1^{(p)} = \sigma(\vec{w}_2 \cdot \vec{x}^{(p)})$$

$$\frac{\partial J}{\partial z} \frac{dz}{da_1} \frac{da_1}{d \vec{w}_1} = \frac{1}{2} (c^{(p)} - z^{(p)}) \sigma'(\vec{w}_1 \cdot \vec{h}^{(p)}) \begin{bmatrix} w_{11}^1 & w_{12}^1 \\ w_{21}^1 & w_{22}^1 \end{bmatrix} \vec{h}^{(p)}$$

$$\nabla_{\vec{w}_1} a_1 = [h_1, h_2]^T$$

$$= \frac{1}{N_{\text{TN}}} \sum_{\text{pattern}} (c^{(n)} - z^{(n)}) \delta'(\vec{w}^2 \vec{h}^{(n)}) w_1^2 \frac{\partial}{\partial w_1^2} \delta(\sum_{k=1}^2 w_{1k}^1 x_k + w_2^1 x_2) \quad \text{CHAIN RULE}$$

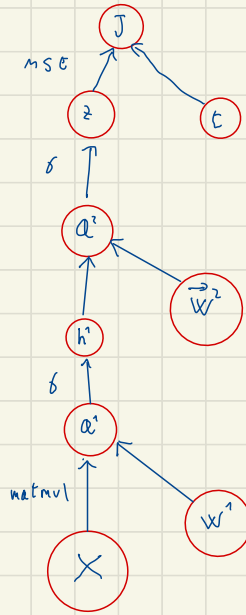
$$= \frac{1}{N_{\text{TN}}} \sum_{\text{pattern}} (c^{(n)} - z^{(n)}) \delta'(\vec{w}^2 \vec{h}^{(n)}) w_1^2 \delta'(\vec{w}_1^1 \vec{x}^{(n)}) \frac{\partial}{\partial w_1^1} (\underbrace{w_{11}^1 x_1}_{\substack{\uparrow \\ w_{11}^1}} + \underbrace{w_{12}^1 x_2}_{\substack{\uparrow \\ x_2}} + \underbrace{w_2^1}_{=0 \text{ odd bias}})$$

MATRIX NOTATION

on a single example for simplicity

$$\frac{dJ}{dw^2} = \frac{dJ}{dz} \frac{dz}{da^2} \frac{da^2}{dw^2} =$$

$$= \begin{bmatrix} -(t-z) \end{bmatrix} \begin{bmatrix} \delta(a^2) (1-\delta(a^2)) \end{bmatrix} \begin{bmatrix} w_1^2 & w_2^2 \\ h_1 & h_2 \end{bmatrix}$$



$$\frac{dJ}{dw^1} = \frac{dJ}{dz} \frac{dz}{da^2} \frac{da^2}{dh} \frac{dh}{da^1} \frac{da^1}{dw^1} \rightarrow \mathbb{R}^{2 \times (2 \times 2)}$$

$$= \begin{bmatrix} -(t-z) \end{bmatrix} \begin{bmatrix} \delta(a^2) (1-\delta(a^2)) \end{bmatrix} \begin{bmatrix} w_1^2 & w_2^2 \\ h_1 & h_2 \end{bmatrix} \begin{bmatrix} \frac{dh_1}{da^1} & \frac{dh_2}{da^1} \\ \frac{da^2}{da^1} \end{bmatrix} \begin{bmatrix} \frac{da^1}{dw_1^1} \\ \frac{da^1}{dw_2^1} \end{bmatrix}$$

$$\begin{bmatrix} \frac{dh_1}{da^1} & \frac{dh_2}{da^1} \\ \frac{da^2}{da^1} \end{bmatrix} \begin{bmatrix} \frac{da^1}{dw_1^1} \\ \frac{da^1}{dw_2^1} \end{bmatrix}$$

$\mathbb{R}^{2 \times (2 \times 2)}$

$$a_1 = w_{11}^1 x_1 + w_{12}^1 x_2$$

$$\frac{\partial a_1}{\partial w_{11}^1} = x_1$$

$$\frac{\partial a_1}{\partial w_{12}^1} = x_2 \in \mathbb{R}^{1 \times 1 \times 2}$$

$$\frac{\partial a_1}{\partial w_2^1} = 0 \in \mathbb{R}^{1 \times 1 \times 2}$$

$$\frac{\partial a_2}{\partial w_{22}^2} = h \in \mathbb{R}^{1 \times 1 \times 2}$$

JACOBIAN IS
DIAGONAL
SINCE ACT. FUNCTION
IS ELEMENT-WISE

$$\begin{bmatrix} \delta(a_1^2) (1-\delta(a_1^2)) & 0 \\ 0 & \delta(a_2^2) (1-\delta(a_2^2)) \end{bmatrix}$$

$$\begin{bmatrix} x_1 & x_2 \\ 0 & 0 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 0 \\ x_1 & x_2 \end{bmatrix}$$

ELEMENT-WISE
MULT. WITH
DIAGONAL

$$\odot \begin{bmatrix} \delta(a_1^2) w_1^2 & \delta(a_1^2) w_2^2 \end{bmatrix}$$

VERIFY CORRECTNESS OF MATRIX FORMULATION (DIFFICULT: PRODUCT WITH TENSORS)

UNF. OF TENSOR

$$\frac{da^1}{dw^1} \begin{bmatrix} x_1 & x_2 & 0 & 0 \\ 0 & 0 & x_1 & x_2 \end{bmatrix}$$

$$(t-z) (\sigma(a) (1-\sigma(a))) \left[w_1^2 \sigma(a_1) (1-\sigma(a_1)), w_2^2 \sigma(a_2) (1-\sigma(a_2)) \right] \cdot \frac{da^1}{dw^1}$$

$\uparrow \mathbb{R}^{1 \times 2 \times 2}$

$$= \begin{bmatrix} w_1^2 \sigma(\dots) x_1, & w_1^2 \sigma(\dots) x_2, & w_2^2 \sigma(\dots) x_1, & w_2^2 \sigma(\dots) x_2 \end{bmatrix}$$

$$= \begin{bmatrix} w_1^2 \sigma(\dots) x_1 & w_1^2 \sigma(\dots) x_2 \\ w_2^2 \sigma(\dots) x_1 & w_2^2 \sigma(\dots) x_2 \end{bmatrix}$$

OR EQUIVALENTLY

$$\left((t-z) (\sigma(a) (1-\sigma(a))) \begin{bmatrix} w_1^2 & w_2^2 \end{bmatrix} \odot \begin{bmatrix} \sigma(a_1) (1-\sigma(a_1)) & \sigma(a_2) (1-\sigma(a_2)) \end{bmatrix} \right)^T x^T$$

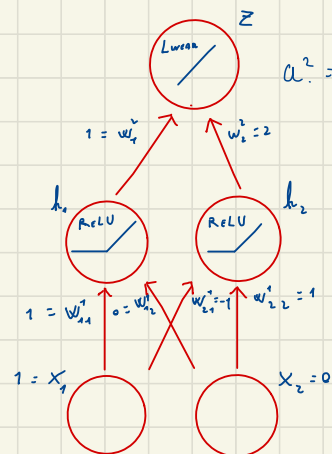
$\uparrow w_i^T$

EXTENSION TO MULTIPLE EXAMPLES IN MATRIX NOTATION

$$J = \frac{1}{n} \begin{bmatrix} \underbrace{J^{(1)}}_{\substack{\text{what} \\ \text{we computed before} \\ \text{for a single example}}} \dots J^{(n)} \end{bmatrix} \begin{bmatrix} 1 \\ i \\ 1 \end{bmatrix}$$

EXERCISE

Let's consider an instantiation of our NN, and let's compute the gradients numerically ^{on a single example} and update the weights.



$$a^2 = w^2 \cdot h$$

$$z = i(a^1) = a^2$$

DERIVATIVE OF LINEAR

$$l(x) = 1$$

non. of ReLU

$$\text{ReLU}'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

Let's assume $c = [2]$

$$N_{\text{trn}} = 1$$

Let's start with the FORWARD propagation.

$$h_1 = \max(0, 1 \cdot 1 + 0 \cdot 0) = \max(0, 1) = 1$$

$$z = 1 \cdot 1 + 2 \cdot 0 = 1$$

$$h_2 = \max(0, 1 \cdot 0 + 1 \cdot 0) = \max(0, 0) = 0$$

$$J = (c - z)^2 = 1$$

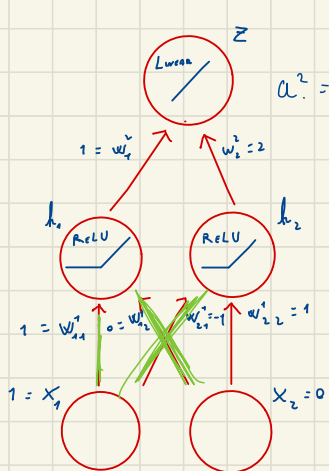
$$\frac{\partial J}{\partial w_{11}^2} = \frac{1}{N_{\text{trn}}} \sum_{i \in \text{trn}} - (c^{(i)} - z^{(i)}) \underbrace{\delta'(\vec{w}^1 \cdot \vec{x}^{(i)})}_{\text{linear}} h_1^{(i)} = -(c - z) \cdot 1 \cdot h_1 = (2 - 1) \cdot 1 \cdot 1 = 1$$

$$\frac{\partial J}{\partial w_{12}^2} = (2 - 1) \cdot 1 \cdot 0 = 0$$

EXERCISE

compute the partial derivatives w.r.t. the weights of the first layer

... compute them also in matrix notation



$$a^2 = w^2 \cdot h \quad \frac{d(1-z)^2}{dz} = -(1-z)$$

$$J = \text{MSE} = \frac{1}{2} (1-2)^2$$

$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

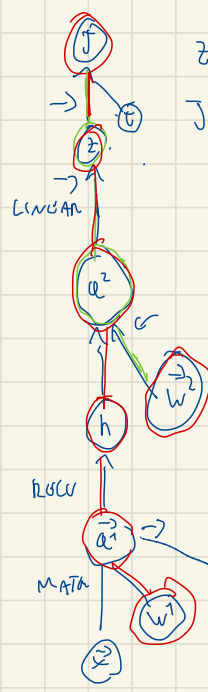
$$h = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

start with the FORWARD

$$\frac{\partial J}{\partial \vec{w}^2} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial a^2} \cdot \frac{\partial a^2}{\partial \vec{w}^2}$$

-1 1 $h = \begin{bmatrix} 1 \\ 0 \end{bmatrix}^T$

$$= \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$



$$z = 1 \quad J = \frac{1}{2}$$

$$a^2 = \vec{w}^2 \cdot \vec{h} \quad \frac{\partial a^2}{\partial \vec{w}^2} = \vec{h}$$

$$\vec{w}^2 = \vec{w}^2 - \eta \begin{bmatrix} -1 \\ 0 \end{bmatrix}$$

$$\rightarrow \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 2 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} \frac{dh_1}{da_1} & 0 \\ 0 & \frac{dh_2}{da_2} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$$

$$\vec{w}_1^{[1,2]} = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 0 \end{bmatrix} = 2 \cdot 0 = 2$$

$$\frac{\partial J}{\partial \vec{w}^1} = \frac{\partial J}{\partial z} \cdot \frac{\partial z}{\partial a^2} \cdot \frac{\partial a^2}{\partial \vec{h}} \cdot \frac{\partial \vec{h}}{\partial a^1} \cdot \frac{\partial a^1}{\partial \vec{w}^1}$$

-1 1 $\begin{bmatrix} 1 & 2 \end{bmatrix}$ 0 $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$ $\begin{bmatrix} 1 & 0 \end{bmatrix}$ $\begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}$

$$= \begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix}$$

Trick $\rightarrow XT$

$$\begin{bmatrix} -1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 0 \end{bmatrix}$$