

Introductory Machine Learning

Some Notes on Backpropagation

$$\frac{\partial \mathcal{L}}{\partial B}(A) = \frac{\partial \mathcal{L}}{\partial B}(BC) =$$

$$A = BC$$

$$\frac{\partial \mathcal{L}}{\partial B} =$$

If \mathcal{L} is a scalar function and $A = BC$, then

$$\frac{\partial \mathcal{L}}{\partial B} = \frac{\partial \mathcal{L}}{\partial A} C^T \quad (1)$$

$$\frac{\partial \mathcal{L}}{\partial C} = B^T \frac{\partial \mathcal{L}}{\partial A}. \quad (2)$$

This can be shown directly using the “usual” chain rule. You should check that the dimensions match up!

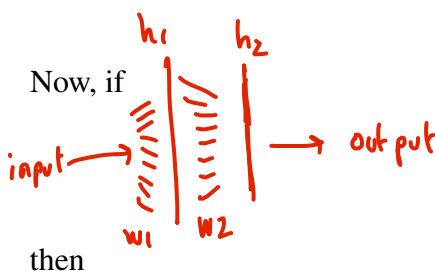
The function \mathcal{L} is our loss function. The use of this is called “backpropagation” because we start with the derivatives in the last layer of the network, and recursively send these “back” to compute the derivatives in the earlier part of the network.

So, if $f = Wx + b$ then

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial f} x^T. \quad (3)$$

If the loss \mathcal{L} is squared error $\mathcal{L} = \frac{1}{2}(y - f)^2$ then

$$\frac{\partial \mathcal{L}}{\partial f} = (f - y). \quad (4)$$



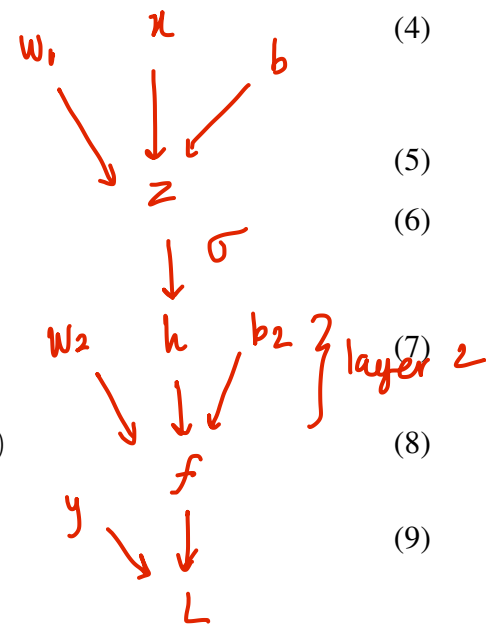
$$f = W_2 h + b_2 \quad (5)$$

$$h = W_1 x + b_1 \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial f} h^T = (f - y) h^T$$

$$\frac{\partial \mathcal{L}}{\partial h} = W_2^T \frac{\partial \mathcal{L}}{\partial f} = W_2^T (f - y) \quad (8)$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial f} = (f - y). \quad (9)$$



$$\frac{\partial \mathcal{L}}{\partial b_2} = 1 \cdot \frac{\partial \mathcal{L}}{\partial f} \cdot \frac{\partial f}{\partial b_2} \quad \text{because } f = W_2 h + b_2 = 1$$

Then, we have that

$$\frac{\partial \mathcal{L}}{\partial W_1} = \frac{\partial \mathcal{L}}{\partial h} x^T = W_2^T \frac{\partial \mathcal{L}}{\partial f} x^T = W_2^T (f - y) x^T \quad (10)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial \mathcal{L}}{\partial h} = W_2^T (f - y). \quad (11)$$

This is just another linear model, so it will train to be equivalent to least squares regression. But if we add a ReLU nonlinearity, we get the two-layer neural network

$$h = \text{ReLU}(W_1 x + b_1) \quad (12)$$

$$f = W_2 h + b_2. \quad (13)$$

Let the loss for an example (x, y) be $\mathcal{L} = \frac{1}{2}(y - f(x))^2$. From the above calculations we then have for the second layer

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial f} h^T \quad \checkmark \quad \text{I get it} \quad (14)$$

$$\frac{\partial \mathcal{L}}{\partial h} = W_2^T \frac{\partial \mathcal{L}}{\partial f} \quad \checkmark \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial f}. \quad \checkmark \quad (16)$$

For the hidden layer we have

$$\frac{\partial \mathcal{L}}{\partial W_1} = \mathbb{1}(h > 0) \frac{\partial \mathcal{L}}{\partial h} x^T \quad (17)$$

$$= \mathbb{1}(h > 0) W_2^T \frac{\partial \mathcal{L}}{\partial f} x^T \quad (18)$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \mathbb{1}(h > 0) \frac{\partial \mathcal{L}}{\partial h} \quad (19)$$

$$= \mathbb{1}(h > 0) W_2^T \frac{\partial \mathcal{L}}{\partial f}. \quad (20)$$

For a two-layer network *for classification* we have

$$h = \text{ReLU}(W_1 x + b_1) \quad (21)$$

$$f = W_2 h_1 + b_2 \quad (22)$$

$$p = \text{Softmax}(f) \quad (23)$$

where again W_j and b_j are matrices or vectors of the appropriate dimensions. Let the loss for an example (x, y) be the log-loss $\mathcal{L} = -\log p(y | x)$. Then

$$\frac{\partial \mathcal{L}}{\partial f} = \begin{pmatrix} p_1 - \mathbb{1}(y = 1) \\ p_2 - \mathbb{1}(y = 2) \\ p_3 - \mathbb{1}(y = 3) \end{pmatrix} \in \mathbb{R}^3 \quad (24)$$

Backpropagating to the second layer we have

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{\partial \mathcal{L}}{\partial f} h_1^T \quad (25)$$

$$\frac{\partial \mathcal{L}}{\partial h} = W_2^T \frac{\partial \mathcal{L}}{\partial f} \quad (26)$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial f} \quad (27)$$

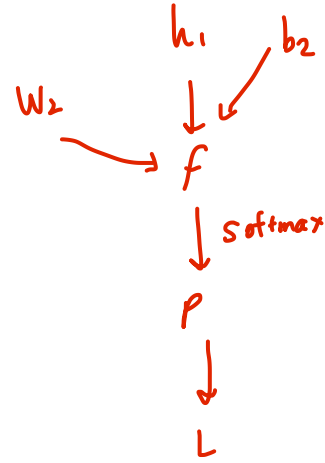
Then backpropagating to the first layer

$$\frac{\partial \mathcal{L}}{\partial W_1} = \mathbb{1}(h > 0) \frac{\partial \mathcal{L}}{\partial h} x^T \quad (28)$$

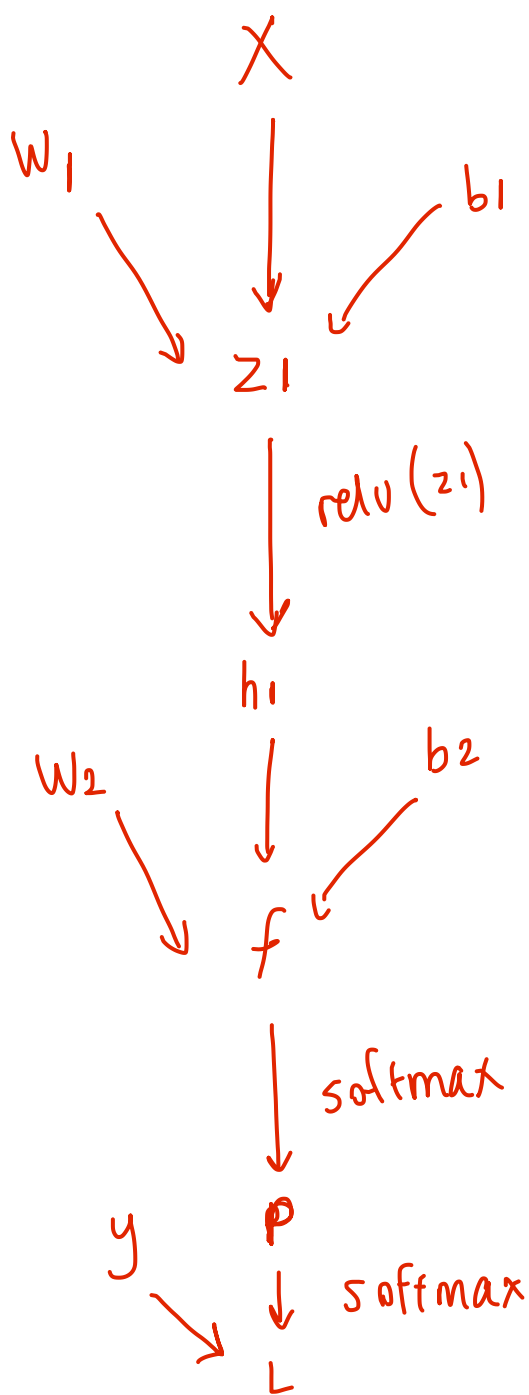
$$= \mathbb{1}(h > 0) W_2^T \frac{\partial \mathcal{L}}{\partial f} x^T \quad (29)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \mathbb{1}(h > 0) \frac{\partial \mathcal{L}}{\partial h} \quad (30)$$

$$= \mathbb{1}(h > 0) W_2^T \frac{\partial \mathcal{L}}{\partial f}. \quad (31)$$



We'll leave it to you to extend the calculations (and implementation!) to a three-layer network, if you are so-inclined.



relu

softmax

Loss

$$f = W_2 h_1 + b_2$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial b_2}$$

Suppose

$$h_1 = \text{relu}(w_1 x + b_1)$$

$$f = w_2 h_1 + b_2$$

$$L = \frac{1}{2} (y - f(x))^2$$

$$\begin{aligned} \frac{\partial L}{\partial f} &= 2 \times \frac{1}{2} (y - f(x)) (0 - 1) \\ &= (y - f) (-1) \\ &= (f - y) \end{aligned}$$

$$\frac{\partial f}{\partial b_2} = 0 + 1 = 1$$

$$\frac{\partial f}{\partial w_2} = h_1$$

$$\frac{\partial f}{\partial h_1} = w_2$$

$$\frac{\partial h_1}{\partial w_1} = x$$

$$\frac{\partial h_1}{\partial x_1} = \text{not needed}$$

$$\frac{\partial h_1}{\partial b_1} = 0 + 1 = 1$$

LAYER 1

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1} \cdot \frac{\partial h_1}{\partial b_1} \\ &= (f - y) \cdot w_2^T \cdot 1 \end{aligned}$$

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial b_2}$$

$$= \underbrace{(f - y)}_{\text{"d scores"}} (1)$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_2}$$

$$= (f - y) \cdot h_1^T$$

idle why transpose but probs for dimension reasons

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1}$$

$$= (f - y) \cdot w_2^T$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_1}$$

$$= (f - y) \cdot w_2^T \cdot x^T$$

in code.

Layer 2

- $dw_2 = \text{np.dot}(h^T, \text{dscores})$
makes sense because $\frac{dL}{dw_2} = \frac{dL}{df} \cdot \frac{df}{dw_2} = (f-y) \cdot h_1^T$

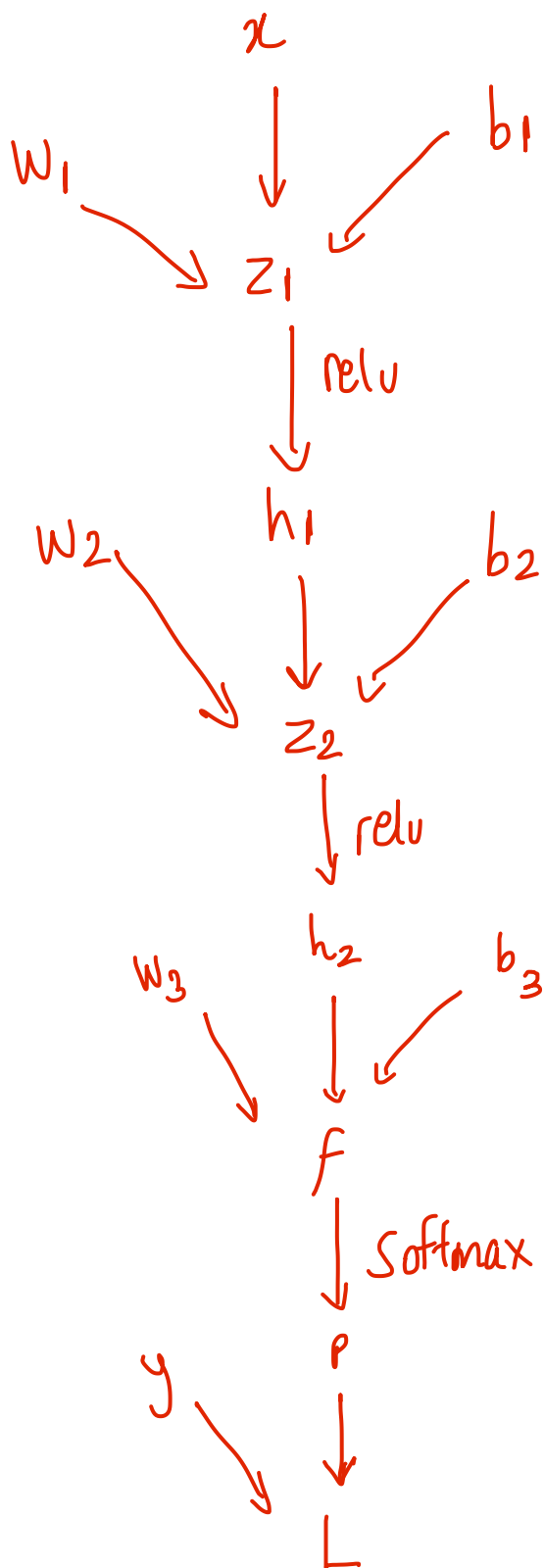
- $db_2 = \text{np.sum}(\text{dscores}, \text{axis}=0, \text{keepdims}=\text{True})$
makes sense because $\frac{dL}{db_2} = \frac{dL}{df} \cdot \frac{df}{db_2} = f-y$

- $dh_1 = \text{np.dot}(\text{dscores}, w_2^T)$
makes sense because $\frac{dL}{dh_1} = \frac{dL}{df} \cdot \frac{df}{dh_1} = (f-y) \cdot w_2^T$

- $dw_1 = \text{np.dot}(x^T, dh_1)$
makes sense because $\frac{dL}{dw_1} = \frac{dL}{df} \cdot \frac{df}{dh_1} \cdot \frac{dh_1}{dw_1}$
 $= (f-y) \cdot w_2^T \cdot x^T$

- $db_1 = \text{np.sum}(dh_1, \text{axis}=0, \text{keepdims}=\text{True})$

makes sense because $\frac{dL}{db_1} = \frac{dL}{df} \cdot \frac{df}{dh_1} \cdot \frac{dh_1}{db_1}$
 $= (f-y) \cdot w_2^T$



$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_2} = \frac{\partial L}{\partial f} \cdot w_3^T$$

$\hat{=}$

$$\frac{\partial L}{\partial h_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} = \frac{\partial L}{\partial f} \cdot w_3^T$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_3} = \frac{\partial L}{\partial f} \cdot h_2^T$$

$$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial b_3} = \frac{\partial L}{\partial f}$$

$$h_1 = \text{relu}(w_1 x + b_1)$$

$$h_2 = \text{relu}(w_2 h_1 + b_2)$$

$$f = w_3 h_2 + b_3$$

$$L = \frac{1}{2} (y - f(x))^2$$

$\frac{\partial f}{\partial w_3} = h_2^T$	$\frac{\partial f}{\partial h_2} = w_3^T$
---	---

$$\frac{\partial h_2}{\partial b_2} = 0 + 1 = 1$$

$\frac{\partial h_2}{\partial w_2} = h_1$	$\frac{\partial h_2}{\partial h_1} = w_2$
---	---

$$\frac{\partial h_1}{\partial b_1} = 0 + 1 = 1$$

$$\frac{\partial h_1}{\partial w_1} = x$$

LAYER 3

$$\frac{\partial L}{\partial b_3} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial b_3} = (f - y) \cdot 1$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_3} = (f - y) \cdot h_2^T$$

$$\frac{\partial L}{\partial h_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} = (f - y) \cdot w_3^T$$

LAYER 2

$$\frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial b_2} = (f - y) \cdot w_3^T \cdot 1$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_2} = (f - y) \cdot w_3^T \cdot h_1^T$$

$$\frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} = (f - y) \cdot w_3^T \cdot w_2^T$$

LAYER 1

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial b_1} \\ &= (f - y) \cdot w_3^T \cdot w_2^T \cdot 1 \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_1} \\ &= (f - y) \cdot w_3^T \cdot w_2^T \cdot x^T \end{aligned}$$

$$\frac{\partial L}{\partial x} = \text{not required}$$

