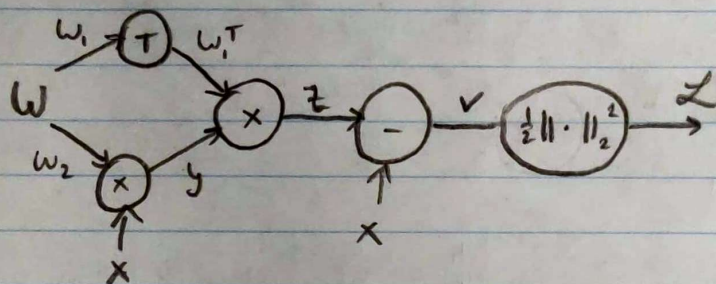


- 1) a. W^T is meant to be a kind of inverse operation to W . While W reduces x to a lower number of dimensions, W^T takes the m -dimensional "compressed" data and tries to restore it as well as possible. Therefore, we want the result of $W^T W x$ to recreate the original data x as well as possible, so we try to minimize the norm of the difference multiplied by $\frac{1}{2}$ as a regularization coefficient.

b.



- c. All of the derivatives for L which get backpropagated to W need to be summed, per the law of total derivatives, to get to the real $\frac{\partial L}{\partial w}$.

- d. let us use backpropagation.

Equations: $L = \frac{1}{2} \|v\|_2^2$

$$v = z - x$$

$$z = W_1^T y$$

$$y = W_2 x$$

Derivatives: $\frac{\partial L}{\partial v} = \frac{1}{2} \frac{\partial}{\partial v} (v^T v) = v$

$$\frac{\partial z}{\partial x} = I$$

$$\frac{\partial z}{\partial y} = W_1$$

$$\frac{\partial z}{\partial w_1} \approx y^T$$

per tensor trick from discussion

$$\frac{\partial z}{\partial w_2} \approx x^T$$

Backpropagation:

$$\frac{\partial L}{\partial z} = \frac{\partial v}{\partial z} \frac{\partial L}{\partial v} = I v = v$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial z}{\partial w_1} \frac{\partial L}{\partial z} = y v^T$$

$$\frac{\partial L}{\partial w_2} = \left(\frac{\partial z}{\partial w_2} \right)^T = y v^T$$

$$\frac{\partial L}{\partial y} = \frac{\partial z}{\partial y} \frac{\partial L}{\partial z} = W_1 v = W_1 v$$

$$\frac{\partial L}{\partial w_2} = \frac{\partial z}{\partial w_2} \frac{\partial L}{\partial z} = W_2 v x^T$$

per discussion trick

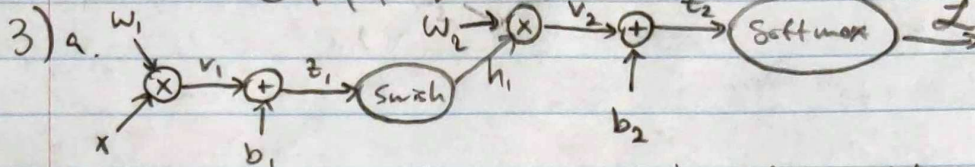
Then $\nabla_W L = \frac{\partial L}{\partial w} = \frac{\partial L}{\partial w_1} + \frac{\partial L}{\partial w_2} = y v^T + W_2 v x^T = W x (W^T W x - x)^T + W (W^T W x - x) x^T$

$$\begin{aligned}\nabla_W \mathcal{L} &= Wx((W^T W - I)x)^T + W(W^T W - I)xx^T \\ &= W(xx^T(W^T W - I) + (W^T W - I)xx^T)\end{aligned}$$

Since $(xx^T(W^T W - I))^T = (W^T W - I)xx^T$,

$$\nabla_W \mathcal{L} = 2W(W^T W - I)xx^T$$

2) I am a C147 student.



(c) $z_2 = v_2 + b_2$

(c) $\frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z_2}$

(c) $v_2 = w_2 h_1$

(c) $\frac{\partial \mathcal{L}}{\partial v_2} = \frac{\partial \mathcal{L}}{\partial z_2} \frac{\partial z_2}{\partial v_2} = I \frac{\partial \mathcal{L}}{\partial z_2} = \frac{\partial \mathcal{L}}{\partial z_2}$

(H) $h_1 = \text{Swish}(z_1)$

(H) $\frac{\partial \mathcal{L}}{\partial h_1} = \frac{\partial v_2}{\partial h_1} \frac{\partial \mathcal{L}}{\partial v_2} = w_2^T \frac{\partial \mathcal{L}}{\partial z_2}$

(H) $z_1 = v_1 + b_1$

(c) $\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial v_2}{\partial w_2} \frac{\partial \mathcal{L}}{\partial v_2} = \frac{\partial \mathcal{L}}{\partial z_2} h_1^T$ ← per the discussion trick

(H) $v_1 = w_1 x$

w_1
 $H \times D$

w_2
 $C \times H$

$$\begin{aligned}\frac{\partial h_i}{\partial z_i} &= \frac{\partial}{\partial z_i} \left(\frac{z_i}{1 + e^{-z_i}} \right) \delta_{ij} \\ &= \frac{1 + e^{-z_i} + z_i e^{-z_i}}{(1 + e^{-z_i})^2} \delta_{ij}\end{aligned}$$

$$\Rightarrow \frac{\partial h_i}{\partial z_i} = \begin{bmatrix} \frac{1 + e^{-z_{i1}} + z_{i1} e^{-z_{i1}}}{(1 + e^{-z_{i1}})^2} & 0 \\ \vdots & \vdots \\ 0 & \frac{1 + e^{-z_{iH}} + z_{iH} e^{-z_{iH}}}{(1 + e^{-z_{iH}})^2} \end{bmatrix} \quad (*)$$

(H) $\frac{\partial \mathcal{L}}{\partial z_1} = \frac{\partial h_1}{\partial z_1} \frac{\partial \mathcal{L}}{\partial h_1} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2}$

(H) $\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial z_1}{\partial b_1} \frac{\partial \mathcal{L}}{\partial z_1} = I \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2}$

(H) $\frac{\partial \mathcal{L}}{\partial v_1} = \frac{\partial z_1}{\partial v_1} \frac{\partial \mathcal{L}}{\partial z_1} = I \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2}$

(H x D) $\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial v_1}{\partial w_1} \frac{\partial \mathcal{L}}{\partial v_1} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2} x^T$ ← discussion trick

Then

a) $\nabla_{w_2} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial z_2} h_1^T$
 $\nabla_{b_2} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial z_2}$

c) $\nabla_{w_1} \mathcal{L} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2} x^T$

$\nabla_{b_1} \mathcal{L} = \frac{\partial h_1}{\partial z_1} w_2^T \frac{\partial \mathcal{L}}{\partial z_2}$

where $\frac{\partial h_i}{\partial z_i}$ is defined as above in (*).