# 1. Backpropagation for autoencoders

With an autoencoder, try to reconstruct the original data dimensions after some operation that reduces the data's dimensionality

E.g. Consider $x \in \mathbb{R}^n$ and $W \in \mathbb{R}^{m \times n}$ where $m < n$. Then $Wx$ is of lower dimensionality than $x$.

One way to design $W$ s.t. $Wx$ still contains key features of $x$ is to minimize $\mathcal{L}$ w.r.t. $W$

$$\mathcal{L} = \tfrac{1}{2}\|W^T W x - x\|^2 \qquad \mathcal{L} = \tfrac{1}{2}\|f(W^T f(Wx)) - x\|^2$$

↑ Linear Example       ↑ Nonlinear Example
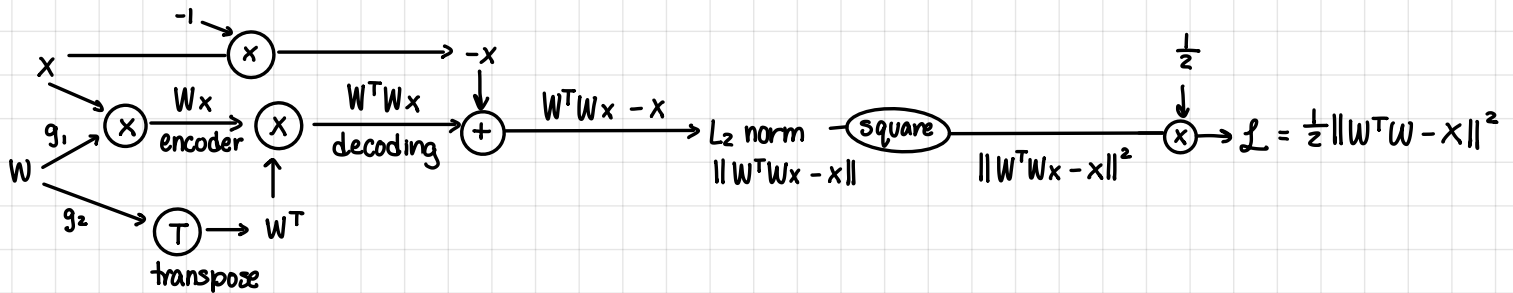
$$\overset{(n \times m)(m \times n)}{\underset{(n \times 1)}{\hookrightarrow W^T W\, x}} \in \mathbb{R}^n$$
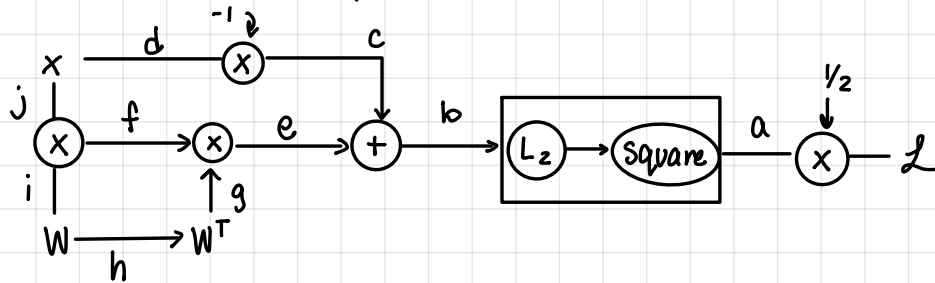
Use the linear example for the following:

## (a) Why does the minimization find a W that ought to preserve info about x

In this minimization of $\mathcal{L} = \tfrac{1}{2}\|W^T W x - x\|^2$, we ensure that we find a matrix $W$ that will preserve the features of $X$ because $Wx$ will reduce the dimensions of $x$ to $m$, but $W^T$ will attempt to reconstruct $x$ from the compressed representation. In other words $Wx$ will result in an $n \times 1$ vector from a $(m \times n)(n \times 1)$ multiplication, whereas $W^T W x$ will result in an $n \times 1$ vector from a $(n \times m)(m \times n)(n \times 1)$ multiplication. If $W$ were to be poorly chosen, important information would be lost and have a high reconstruction error. Minimizing $\mathcal{L}$ forces $W$ to learn an optimal low-dimensional representation where it preserves key features (similar to PCA analysis)

## (b) Draw the computational Graph for $\mathcal{L}$



Setup so that I can solve for part (d)



## (c) In the computational graph, there should be 2 paths to W. How do we account for these two paths when calculating $\nabla_W \mathcal{L}$? Should include mathematical argument.

In the computational graph, the matrix $W$ appears when $W$ maps $x$ to a lower dimension $(Wx)$ and when we reconstruct $(W \xrightarrow{T} W^T \rightarrow W^T W x)$. Both will ultimately converge at $W^T W x$.

Mathematically, we defined $g_1$ to be the path that $W$ takes to become $Wx$ and $g_2$ to be the path that $W$ takes to get to $W \xrightarrow{T} W^T \rightarrow W^T W x$

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial g_1}{\partial W} \cdot \frac{\partial \mathcal{L}}{\partial g_1} + \frac{\partial g_2}{\partial W} \cdot \frac{\partial \mathcal{L}}{\partial g_2}$$

# (d) Calculate the gradient $\nabla_W \mathcal{L}$

$\mathcal{L} = \frac{1}{2}\|W^TWx - x\|^2$

🟩 = Forward Propagation
🟥 = Back Propagation

$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial b}{\partial c} \cdot \frac{\partial \mathcal{L}}{\partial b}$

**First diagram (forward propagation):**

x —d— $\times$ (−1) ——−x—— c

x —j— $\times$ —Wx (f)— $\times$ —$W^TWx$ (e)— $+$ —b— $L_2$ Norm —$\|W^TWx - x\|$— squared —$\|W^TWx-x\|^2$ (a)— $\times$ ($\frac{1}{2}$) — $\frac{1}{2}\|W^TWx-x\|^2$ → $\mathcal{L}$
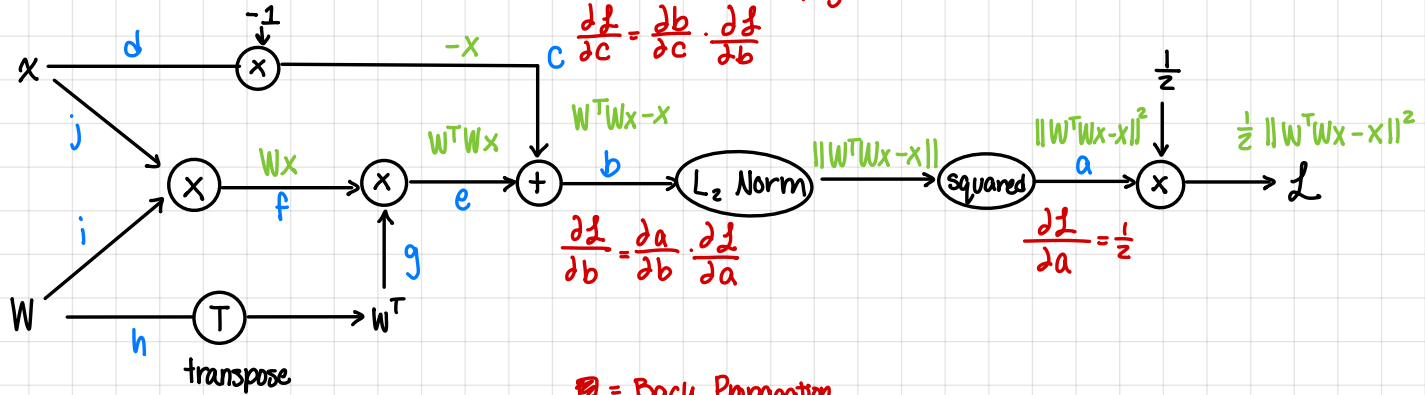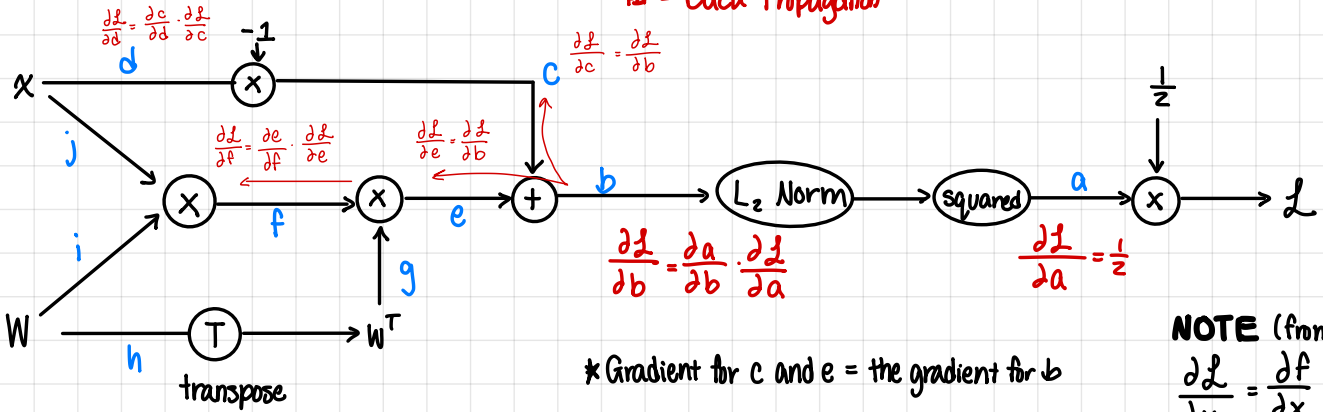
W —h— T —$W^T$ (g)—

transpose

$W^TWx - x$ (b)

$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial a}{\partial b} \cdot \frac{\partial \mathcal{L}}{\partial a}$

$\frac{\partial \mathcal{L}}{\partial a} = \frac{1}{2}$

**Second diagram (Back Propagation):**

🟥 = Back Propagation

$\frac{\partial \mathcal{L}}{\partial d} = \frac{\partial c}{\partial d} \cdot \frac{\partial \mathcal{L}}{\partial c}$

x —d— $\times$ (−1) —— c

$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial \mathcal{L}}{\partial b}$

$\frac{\partial \mathcal{L}}{\partial f} = \frac{\partial e}{\partial f} \cdot \frac{\partial \mathcal{L}}{\partial e}$

$\frac{\partial \mathcal{L}}{\partial e} = \frac{\partial \mathcal{L}}{\partial b}$

x —j— $\times$ —f— $\times$ —e— $+$ —b— $L_2$ Norm — squared —a ($\frac{1}{2}$)— $\times$ → $\mathcal{L}$

W —h— T —$W^T$ (g)—

transpose

$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial a}{\partial b} \cdot \frac{\partial \mathcal{L}}{\partial a}$

$\frac{\partial \mathcal{L}}{\partial a} = \frac{1}{2}$

\* Gradient for c and e = the gradient for b

**NOTE (from Lecture 6)**

$\frac{\partial \mathcal{L}}{\partial x} = \underset{\text{local}}{\frac{\partial f}{\partial x}} \underset{\text{upstream}}{\frac{\partial \mathcal{L}}{\partial f}}$

$\frac{\partial \mathcal{L}}{\partial a} = \frac{1}{2}$

$\frac{\partial \mathcal{L}}{\partial b} = \underset{\text{local gradient}}{\frac{\partial a}{\partial b}} \cdot \underset{\text{upstream gradient}}{\frac{\partial \mathcal{L}}{\partial a}} = 2b \cdot \frac{1}{2} = b = \underset{\text{value comes from forward propagation}}{W^TWx - x}$

$\left(a = b^T b \quad \therefore \frac{\partial a}{\partial b} = 2b\right)$

$\frac{\partial \mathcal{L}}{\partial c} = \frac{\partial \mathcal{L}}{\partial b} = b = W^TWx - x$  and  $\frac{\partial \mathcal{L}}{\partial e} = \frac{\partial \mathcal{L}}{\partial b} = b = W^TWx - x$

NOTE: $\frac{\partial \mathcal{L}}{\partial c}$ and $\frac{\partial \mathcal{L}}{\partial e}$ come from $\oplus$ $\therefore$ distributed gradient

$\frac{\partial \mathcal{L}}{\partial d} = \frac{\partial c}{\partial d} \cdot \frac{\partial \mathcal{L}}{\partial c} = -1 \cdot b = -b = -(W^TWx - x)$

$\frac{\partial \mathcal{L}}{\partial f} = \frac{\partial e}{\partial f} \cdot \frac{\partial \mathcal{L}}{\partial e} = \frac{\partial e}{\partial f} \cdot b = \frac{\partial W^TWx}{\partial Wx} \cdot b = (W^T)^T b = Wb = W(W^TWx - x)$

$\frac{\partial \mathcal{L}}{\partial g} = \frac{\partial e}{\partial g} \times \frac{\partial \mathcal{L}}{\partial e} = b \cdot (Wx)^T = bx^TW^T = (W^TWx - x)x^TW^T$  Uses the trick in class

$\frac{\partial \mathcal{L}}{\partial h} = \left(\frac{\partial \mathcal{L}}{\partial g}\right)^T = (bx^TW^T)^T = Wxb^T = Wx(W^TWx - x)^T$

$\frac{\partial \mathcal{L}}{\partial i} = \frac{\partial f}{\partial i} \cdot \frac{\partial \mathcal{L}}{\partial f} = Wbx^T$

**NOTE (Lecture 7)**

W
$\searrow$
$\times$ → y
$\nearrow$
x

$\frac{\partial \mathcal{L}}{\partial x} = W^T \frac{\partial \mathcal{L}}{\partial y}$

$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial y} x^T$

$b = W^TWx - x$

$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial h} + \frac{\partial \mathcal{L}}{\partial i} = Wxb^T + Wbx^T = Wx(W^TWx - x) + W(W^TWx - x)x^T$

# Problem #2: I am a C147 Student

# Problem #3: NNDL

$D$ = # of neurons in input layer, $H$ = # of neurons in the hidden layer, $C$ = # of neurons in the oupt $(C = 7)$

Swish activation function $\quad \text{swish}(k) = \dfrac{k}{1+e^{-k}} = k\sigma(k) \quad$ where $\sigma(k)$ is sigmoid activation function
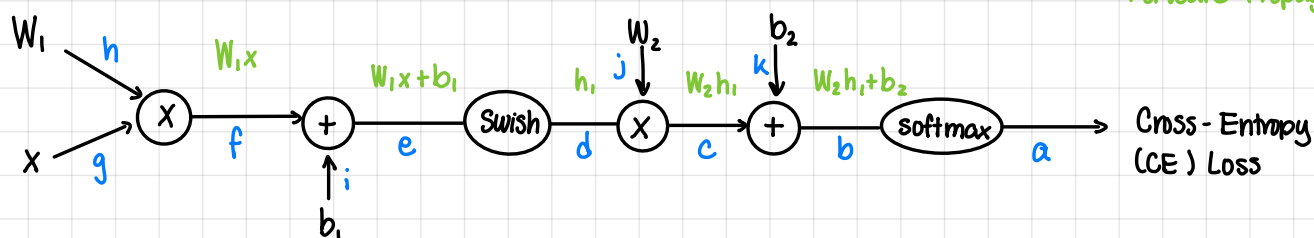


Hidden Layer: $z_1 = W_1 x + b_1 \qquad h_1 = \text{Swish}(z_1) \qquad z_2 = W_2 h_1 + b_2$
$$z_1 \in \mathbb{R}^H \qquad\qquad\qquad z_2 \in \mathbb{R}^C$$

## (a) Draw the computational graph for the 2-layer FC Net

■ = Forward Propagation



## (b) Compute $\nabla_W \mathcal{L} \;\; \nabla_{b_2}\mathcal{L}$ (For the gradient computations you can keep it as $\frac{\partial \mathcal{L}}{\partial z_2}$)

Cross-entropy Loss $= \mathcal{L} = -\sum y_i \log \hat{y}_i$
where $\hat{y}_i$ is the predicted value

$$\hat{y}_i = \frac{e^{z_{2,i}}}{\sum_j e^{z_{2,j}}}$$

$$\therefore \mathcal{L} = -\sum y_i \log\left(\frac{e^{z_{2,i}}}{\sum_j e^{z_{2,j}}}\right) = -\sum y_i \left(z_{2,i} - z_{2,j}\right)$$

$$\frac{d\mathcal{L}}{dz_{2,i}} = \frac{\partial}{\partial z_{2,i}}\left(-\sum y_i(z_{2,i} - z_{2,j})\right) = \hat{y}_i - y_i$$

$$\boxed{\therefore \frac{\partial \mathcal{L}}{\partial z_2} = \hat{y} - y}$$

$$\frac{d\mathcal{L}}{dW_2} = \frac{d\mathcal{L}}{dz_2}\cdot\frac{dz_2}{dW_2} \quad \text{we know that } \frac{d\mathcal{L}}{dz_2} = \hat{y}-y \text{ and } \frac{dz_2}{dW_2} = \frac{\partial W_2 h_1 + b_2}{\partial W_2} = h_1^T$$

$$\hookrightarrow \nabla_{W_2}\mathcal{L} = (\hat{y}-y)\,h_1^T$$

$$\nabla_{b_2}\mathcal{L} = \frac{\partial \mathcal{L}}{\partial b_2} = \frac{\partial \mathcal{L}}{\partial z_2} = (\hat{y}-y)$$

$$\boxed{\begin{array}{l}\textbf{Results}\\[4pt] \nabla_{W_2}\mathcal{L} = (\hat{y}-y)\,h_1^T \\[8pt] \nabla_{b_2}\mathcal{L} = (\hat{y}-y)\end{array}}$$

$y$ comes from the Jacobian

$$J = \frac{\partial y}{\partial x} = \begin{bmatrix} \dfrac{\partial y_1}{\partial x_1} & \dfrac{\partial y_1}{\partial x_2} & \cdots & \dfrac{\partial y_1}{\partial x_m} \\[10pt] \dfrac{\partial y_2}{\partial x_1} & \dfrac{\partial y_2}{\partial x_2} & \cdots & \dfrac{\partial y_2}{\partial x_m} \\[6pt] \vdots & & & \\[6pt] \dfrac{\partial y_n}{\partial x_1} & \dfrac{\partial y_n}{\partial x_2} & \cdots & \dfrac{\partial y_n}{\partial x_m} \end{bmatrix}$$

# (c) Compute $\nabla_{W_1}\mathcal{L}, \nabla_{b_1}\mathcal{L}$

Calculating $\nabla_{W_1}\mathcal{L}$ requires $\frac{\partial}{\partial z_1}$

$h_1 = \text{Swish}(z_1) \rightarrow \text{Swish}(x) = x\sigma(x)$ and $\sigma(x) = \dfrac{1}{1+e^{-z}}$

$\dfrac{\partial \text{Swish}(x)}{\partial x} = \sigma(x) + x\sigma(x)(1-\sigma(x))$

Replacing $x$ with $z_1$ we then have $\sigma(z_1) + z_1\sigma(z_1)(1-\sigma(z_1))$

$\dfrac{\partial \mathcal{L}}{\partial z_1} = \dfrac{\partial \mathcal{L}}{\partial h_1} \cdot \dfrac{\partial h_1}{\partial z_1} = \left(W_2^T(\hat{y}-y)\right) \circ \dfrac{\partial h_1}{\partial z_1} = \left((W_2^T(\hat{y}-y)) \circ \left[\sigma(z_1) + z_1\sigma(z_1)(1-\sigma(z_1))\right]\right)$

we use $\dfrac{\partial \mathcal{L}}{\partial z_1}$ to find $\dfrac{\partial \mathcal{L}}{\partial W_1}$ and $\dfrac{\partial \mathcal{L}}{\partial b_1}$

$\nabla_{W_1}\mathcal{L} = \dfrac{\partial \mathcal{L}}{\partial z_1} \cdot \dfrac{\partial z_1}{\partial W_1} = (\nabla_{z_1}\mathcal{L})\left(x^T\right) = \left[(W_2^T(\hat{y}-y)) \circ \left(\sigma(z_1) + z_1\sigma(z_1)(1-\sigma(z_1))\right)\right]$

Comes from the trick in class $\dfrac{\partial z_1}{\partial W_1} = \dfrac{\partial W_1 x + b_1}{\partial W_1} = x^T$

$\nabla_{b_1}\mathcal{L} = \nabla_{z_1}\mathcal{L} = (W_2^T(\hat{y}-y)) \circ ($