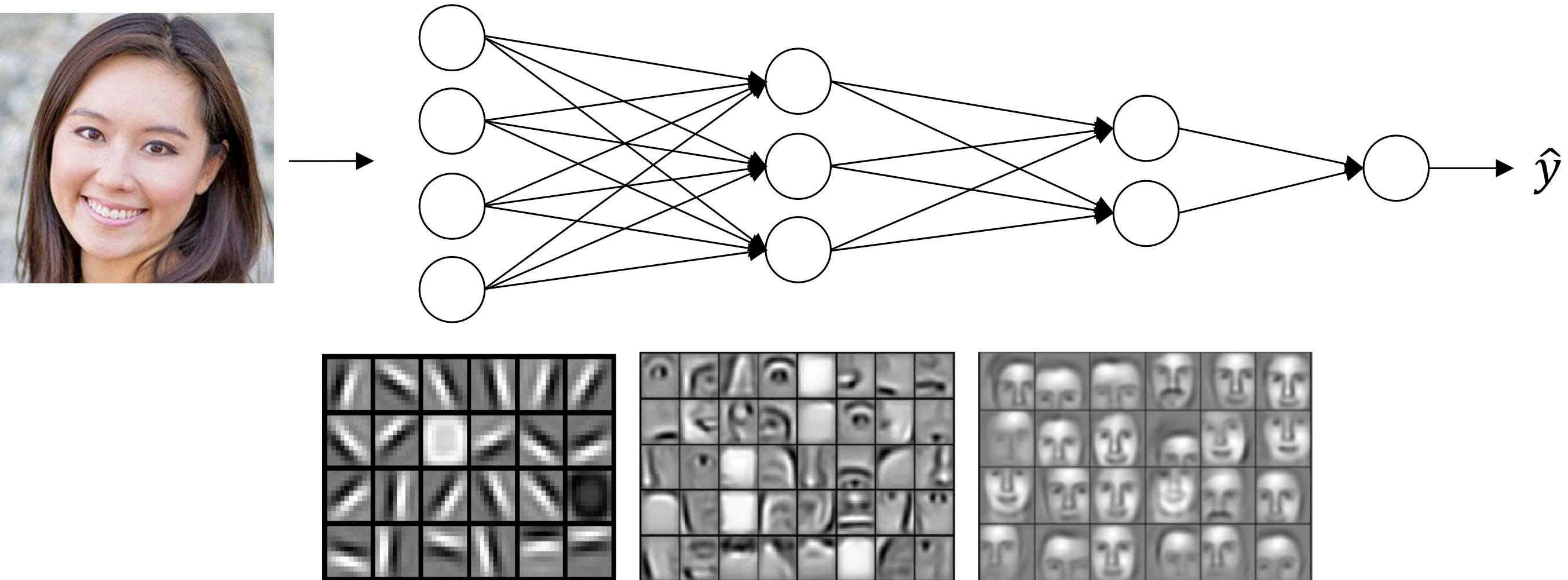


# Backpropagation

Seyoung Yun

- [http://cs231n.stanford.edu/slides/2017/cs231n\\_2017\\_lecture4.pdf](http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture4.pdf)
- <http://cs231n.github.io/optimization-2/>
- <http://yann.lecun.com/exdb/publis/pdf/lecun-98b.pdf>

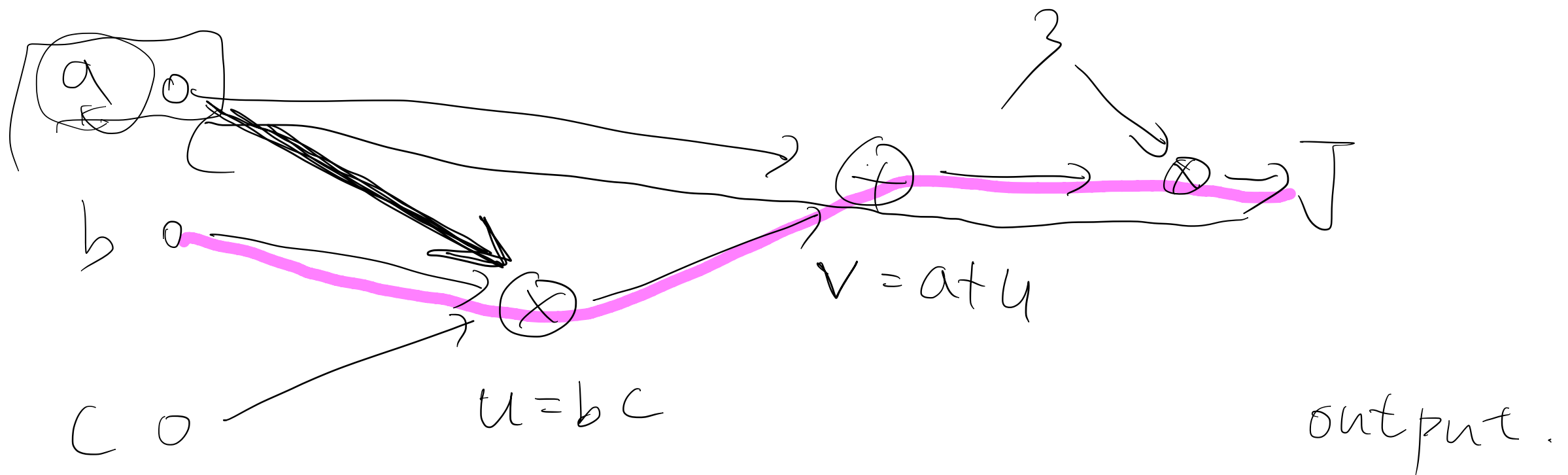
# Intuition about deep representation



Informally: There are functions you can compute with a “small” L-layer deep neural network that shallower networks require exponentially more hidden units to compute.

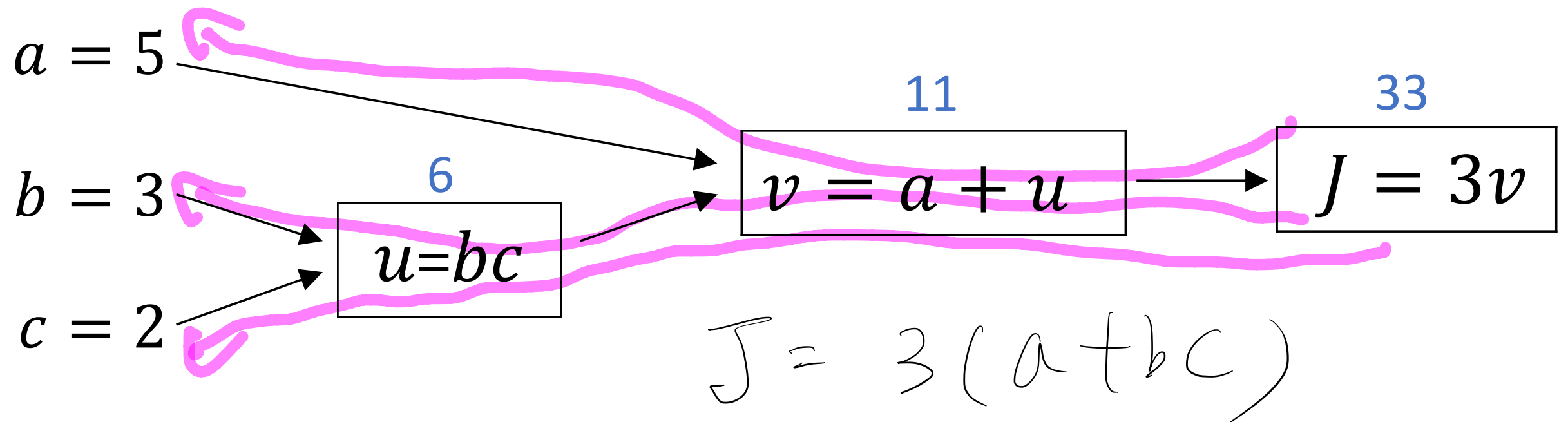
# Computational Graph

$$J(a, b, c) = \underline{3(a + bc)}$$



| input                           |                                 |                                 |       |
|---------------------------------|---------------------------------|---------------------------------|-------|
| $\frac{\partial J}{\partial a}$ | $\frac{\partial J}{\partial v}$ | $\frac{\partial v}{\partial a}$ | $= 3$ |
| $\frac{\partial J}{\partial b}$ | $\frac{\partial J}{\partial v}$ | $\frac{\partial v}{\partial u}$ | $= 3$ |
| $\frac{\partial J}{\partial c}$ | $\frac{\partial J}{\partial v}$ | $\frac{\partial v}{\partial c}$ | $= 3$ |

# Computing derivatives

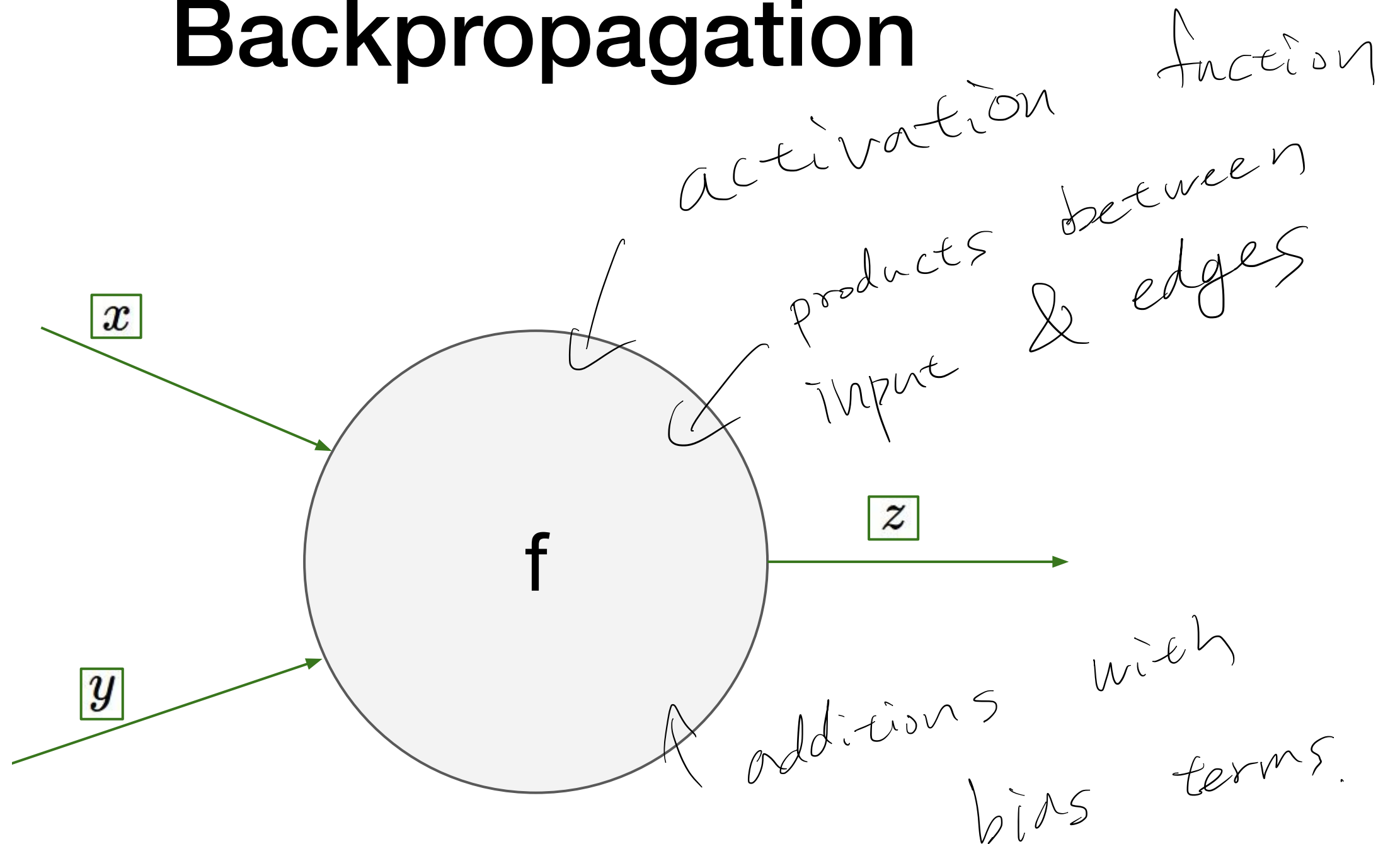


$$\frac{\partial J}{\partial a} = 11$$

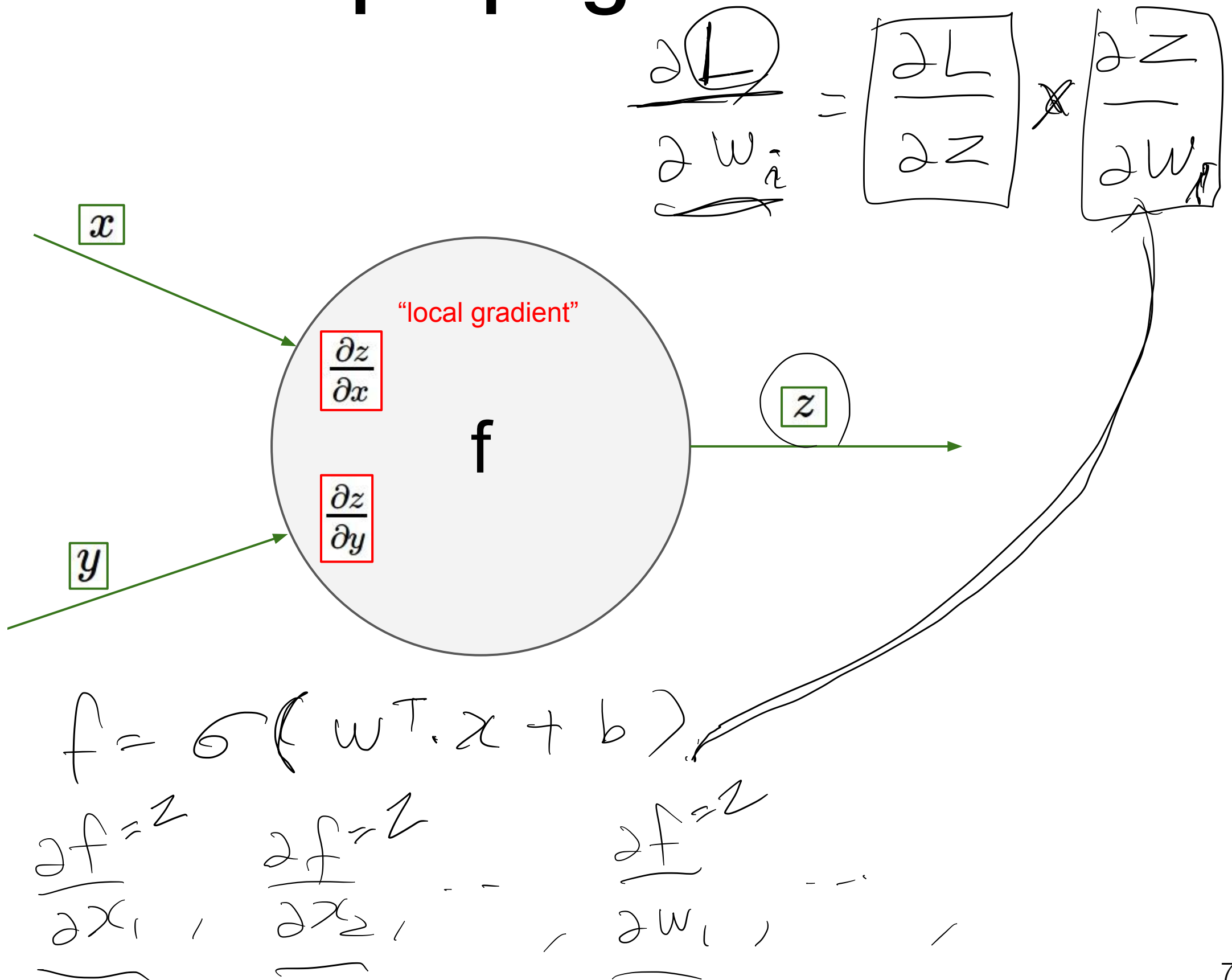
$$\frac{\partial J}{\partial b} = 6$$

$$\frac{\partial J}{\partial c} = 6$$

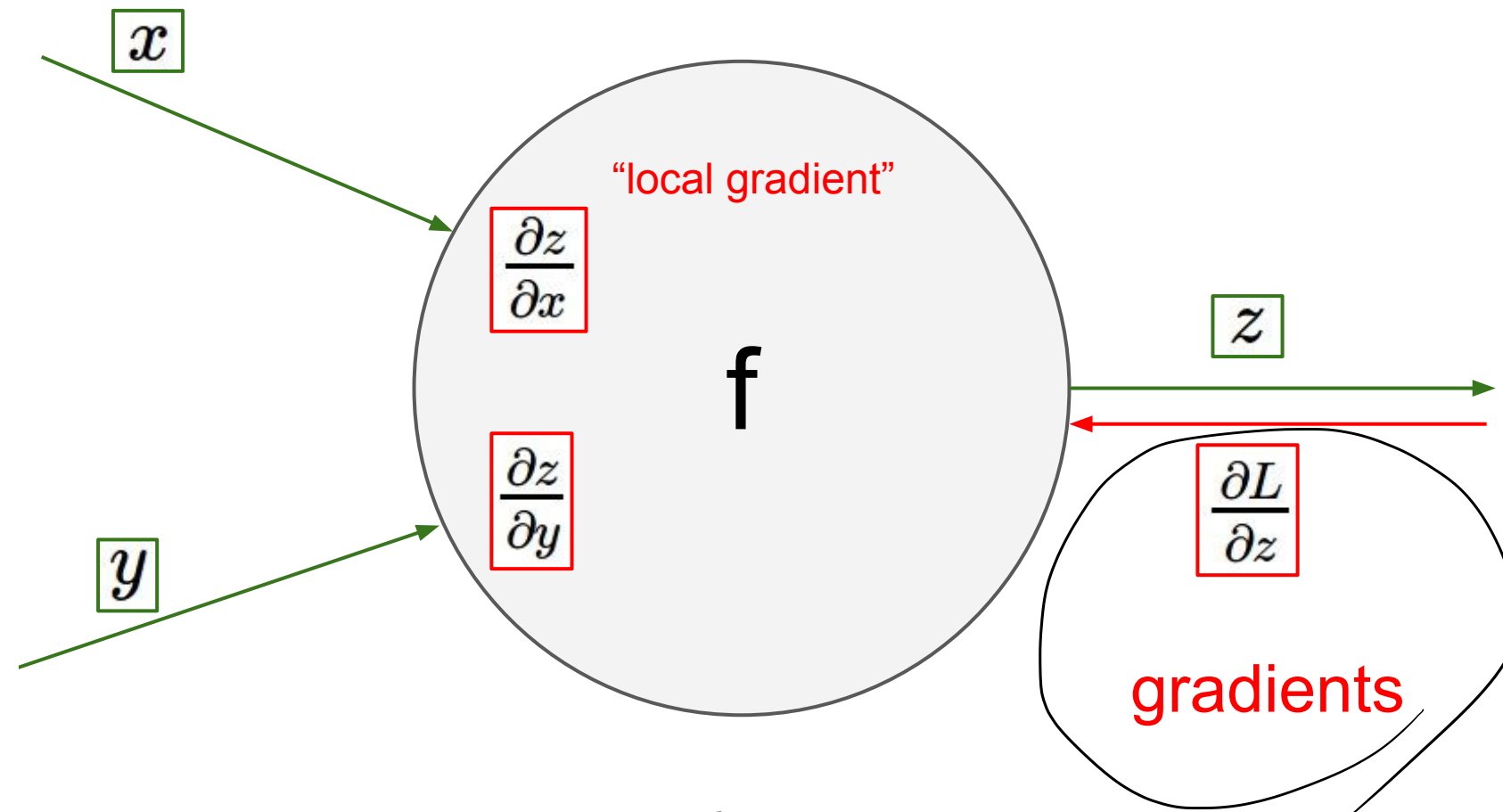
# Backpropagation



# Backpropagation



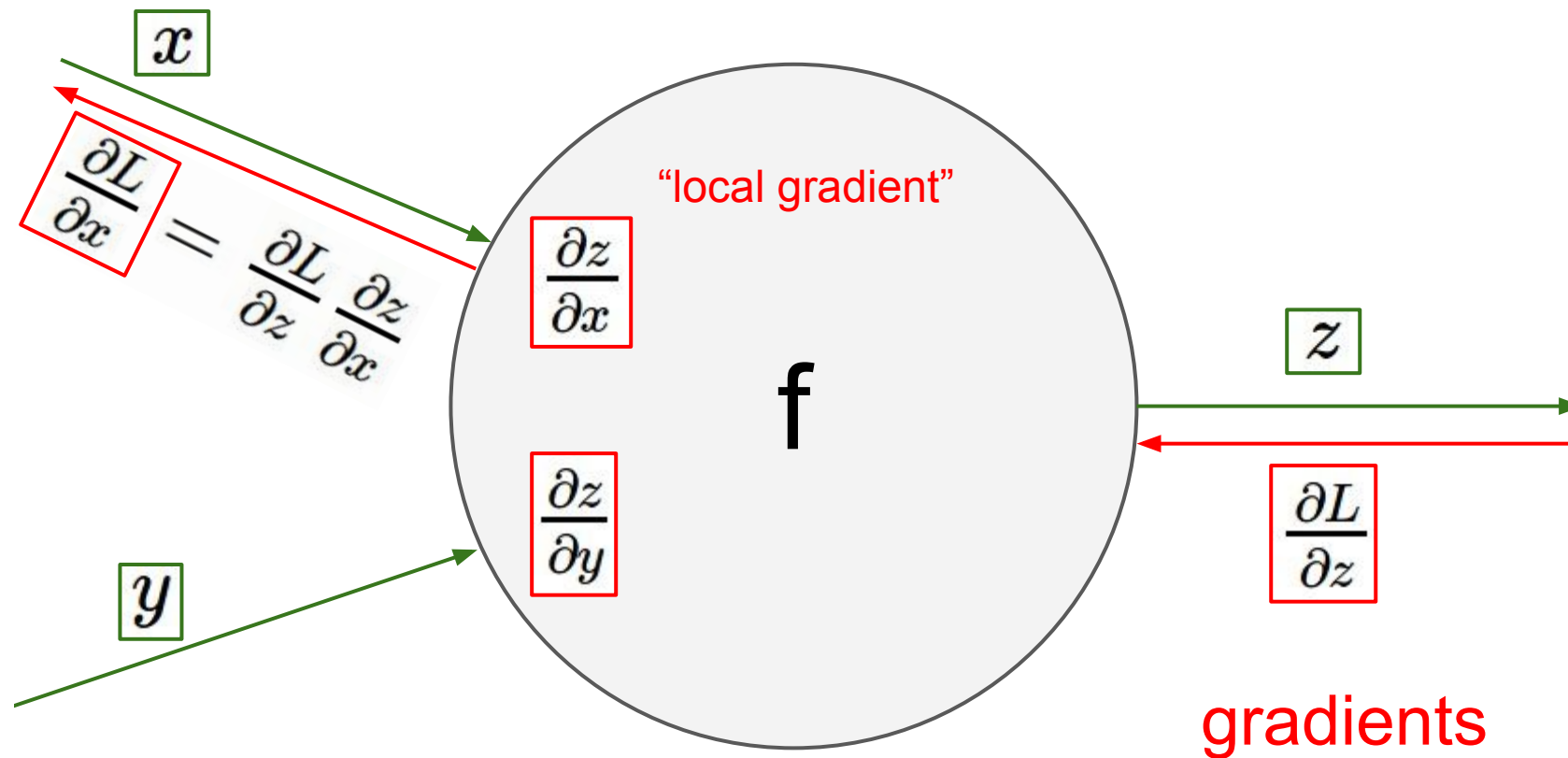
# Backpropagation



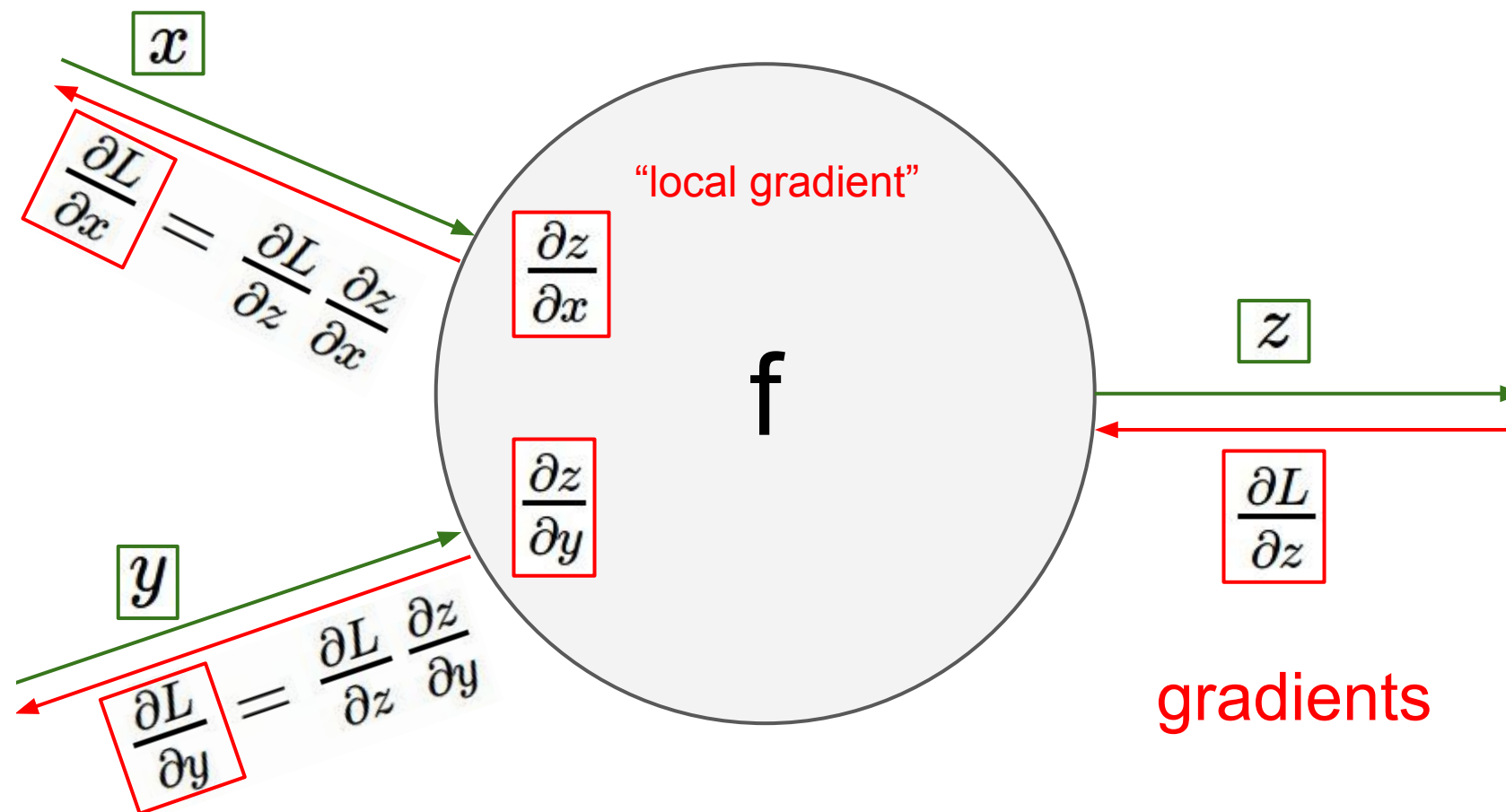
$$\frac{\partial L}{\partial z} \times \frac{\partial z}{\partial w_i}$$



# Backpropagation

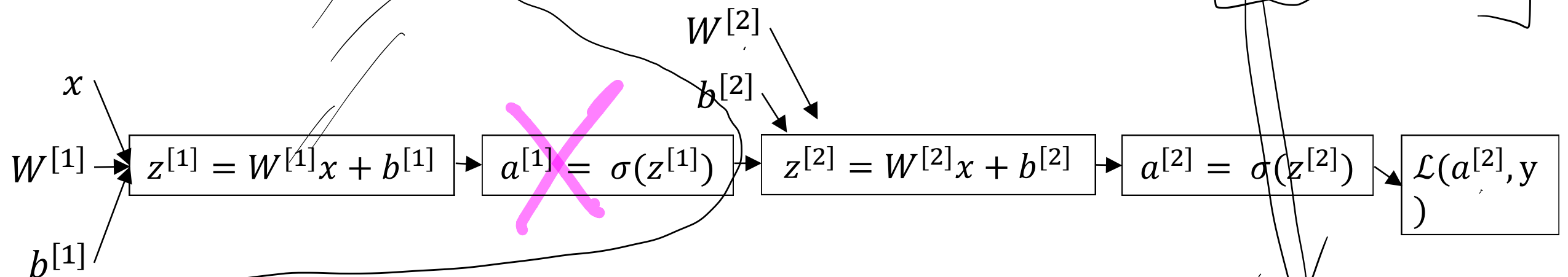


# Backpropagation



# Neural Net gradients

$$\nabla_{w^{[2]}} L$$



$$\begin{bmatrix} \sigma'(z_1^{[1]}) \\ \sigma'(z_2^{[1]}) \\ \vdots \end{bmatrix}$$

$$\nabla_{w^{[2]}} L$$

$$= \underbrace{\left( \nabla_{z^{[2]}}(w^{[2]}) \right)}_{\text{Jacobian of } z^{[2]} \text{ w.r.t } w^{[2]}} \nabla_{z^{[2]}} L$$

$$= \nabla_{z^{[2]}}(w^{[2]}) \underbrace{\left( \nabla_{a^{[2]}}(z^{[2]}) \right)}_{\text{Jacobian of } a^{[2]} \text{ w.r.t } z^{[2]}} \nabla_{a^{[2]}} L$$

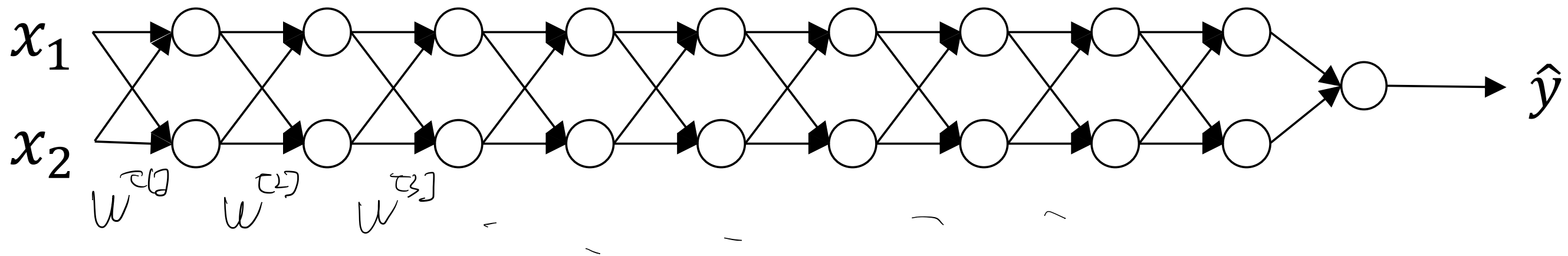
$$\begin{bmatrix} \frac{\partial L}{\partial w_1^{[2]}} \\ \frac{\partial L}{\partial w_2^{[2]}} \end{bmatrix}$$

"

$$\begin{bmatrix} \frac{\partial z_1^{[2]}}{\partial w_1^{[2]}} & \frac{\partial z_2^{[2]}}{\partial w_1^{[2]}} \\ \frac{\partial z_1^{[2]}}{\partial w_2^{[2]}} & \frac{\partial z_2^{[2]}}{\partial w_2^{[2]}} \end{bmatrix}$$

$$\begin{bmatrix} \frac{\partial L}{\partial z_1^{[2]}} \\ \frac{\partial L}{\partial z_2^{[2]}} \end{bmatrix}$$

# Vanishing/exploding gradients



$$\nabla_{w^{\tau_1}} L$$

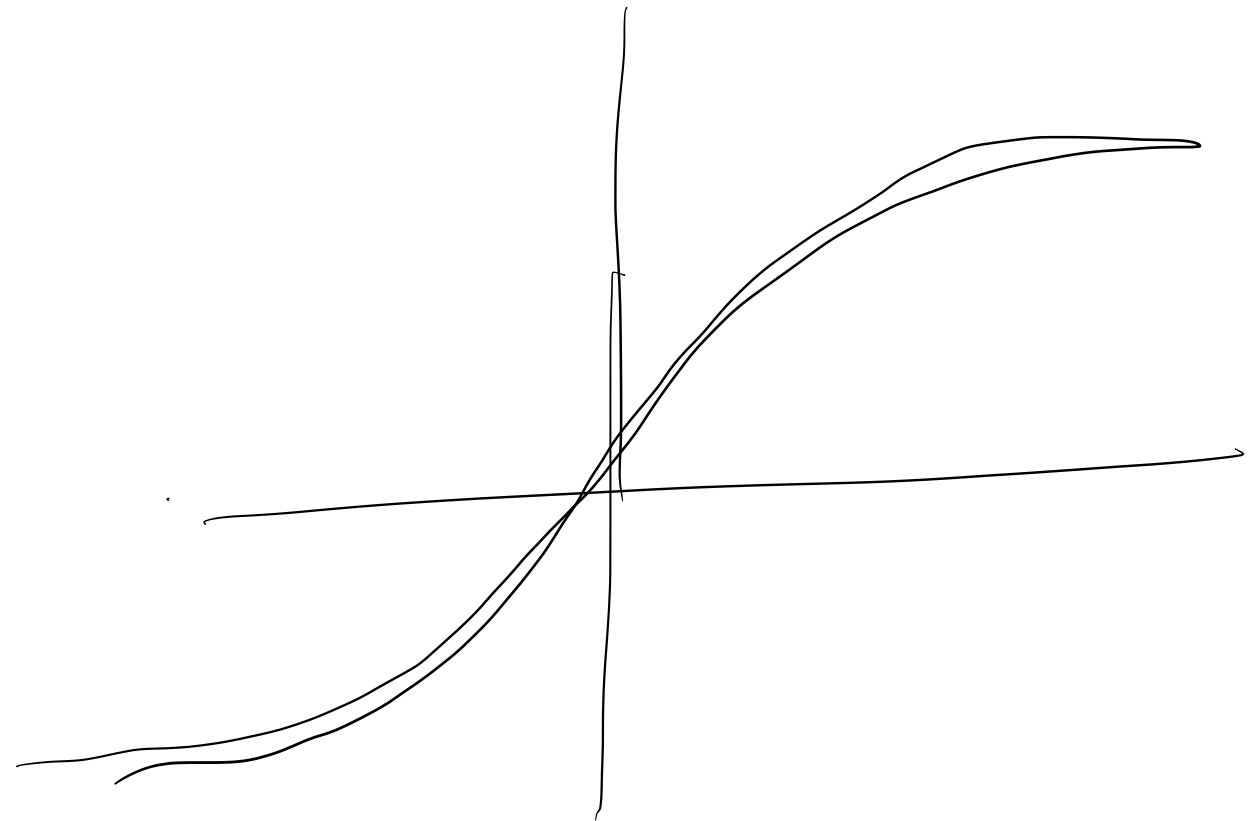
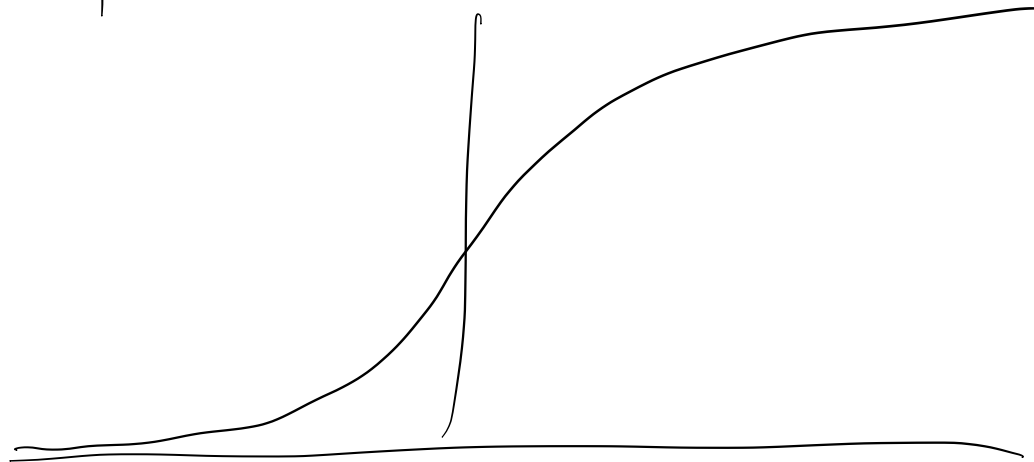
$$A = X_1 X_2 X_3 X_4 \dots X_{10}$$

$$\underbrace{\quad}^{<1} \quad \underbrace{\quad}^{<1} \quad \underbrace{\quad}^{<1}$$

$$\|A\|_2 \leq \underbrace{\|X_1\|_2}_{>1} \underbrace{\|X_2\|_2}_{>1} \dots \underbrace{\|X_{10}\|_2}_{>1}$$

# Derivatives of activation function

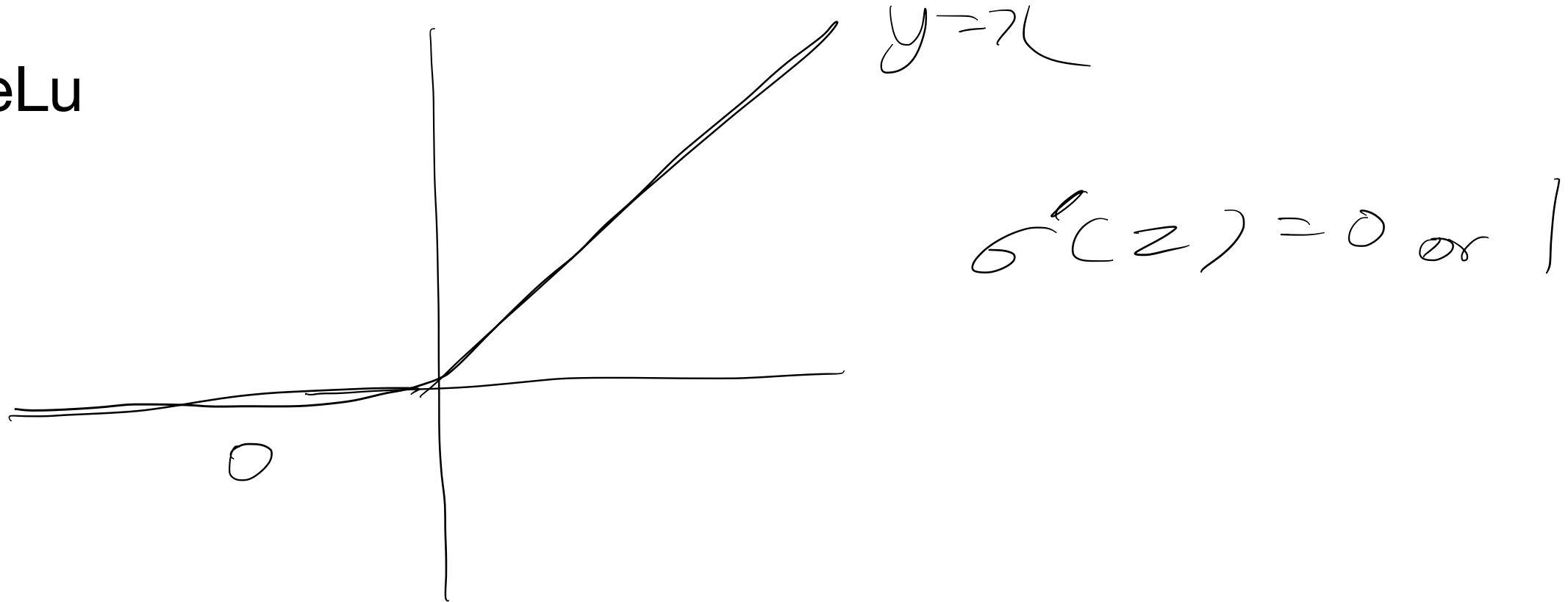
- Sigmoid and tanh



$$J_{a^{[L]}}(z^{[L]}) = \begin{bmatrix} \sigma'(z_1^{[L]}) & 0 \\ 0 & \sigma'(z_n^{[L]}) \end{bmatrix}$$

# Derivatives of activation function

- ReLu



$$\left[ \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right] \cdot \left[ \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \right] = 0$$