

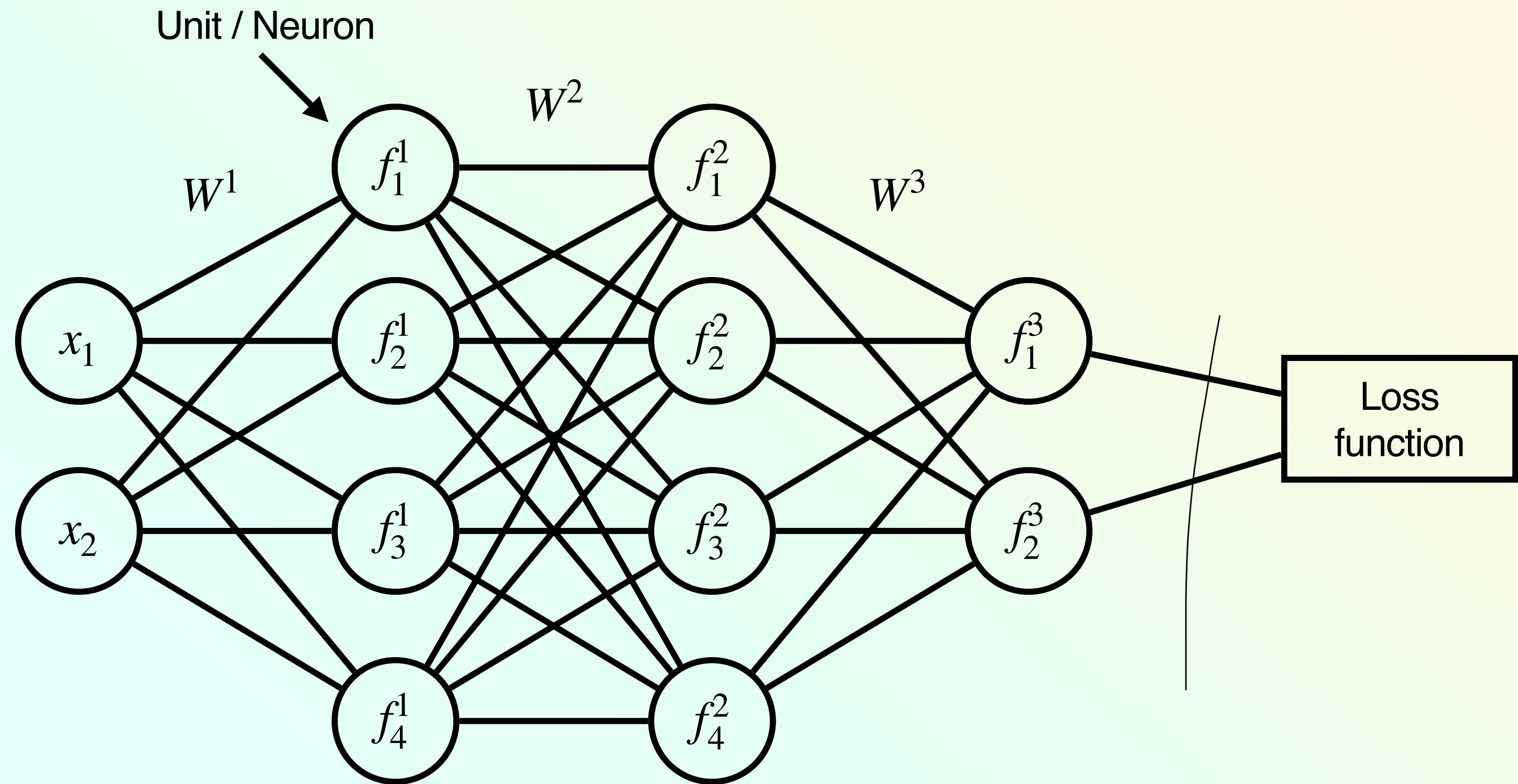
# GRADIENTS OF NEURAL NETWORKS



Input Layer

Hidden Layers

Output Layer



$$h^0$$

$$h^1 = f^1(h^0, W^1)$$

$$h^2 = f^2(h^1, W^2)$$

$$h^3 = f^3(h^2, W^3)$$



# THE STRUCTURE OF NEURAL NETWORKS

## ABSTRACT PROCESS

We can write the neural network outputs as a sequential process.

$$\begin{aligned} h^0 &= x \\ h^1 &= f^1(h^0, W^1) \\ h^2 &= f^2(h^1, W^2) \\ &\vdots \\ h^i &= f^i(h^{i-1}, W^i) \\ &\vdots \\ h^k &= f^k(h^{k-1}, W^k) \end{aligned}$$

To be concise, we can write the network output as  $h^k = f(x, \theta)$ ,  $\theta = \{W^i\}_{i=1}^k$



# LOSS FUNCTION

## EXAMPLE

$$l(\theta) = \frac{1}{2} \mathbf{E} \left[ (f(X, \theta) - Y)^2 \right] \rightarrow \begin{array}{l} \text{will apply it backward} \\ \text{for each } h \text{ with respect} \\ \text{to } w_{\sum} \text{ for each layer} \end{array}$$

We have a batch of data  $D = (x, y)$  of  $m$  samples

$x \in \mathbb{R}^{m \times n_0}$  and  $y \in \mathbb{R}^{m \times 1}$ ,  $x_i$  and  $y_i$  are the features and target for the  $i^{\text{th}}$  data point.

$$l_D(\theta) = \frac{1}{m} \frac{1}{2} \sum_{i=1}^m (f(x_i, \theta) - y_i)^2$$



# LOSS FUNCTION

## EXAMPLE

$$l(\theta) = \frac{1}{2} \mathbf{E} \left[ (f(X, \theta) - Y)^2 \right]$$

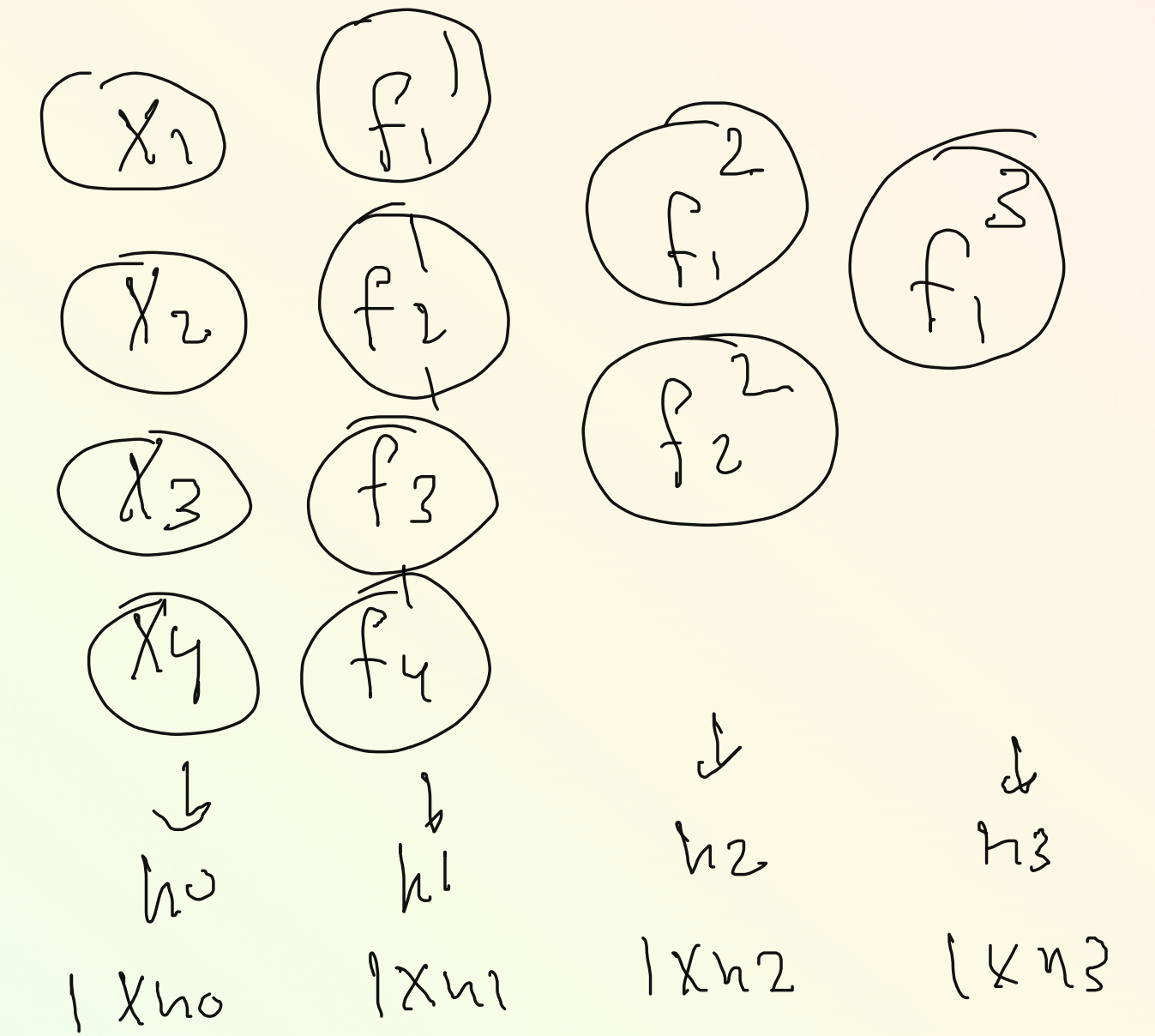
We have a batch of data  $D = (x, y)$  of  $m$  samples

$x \in \mathbb{R}^{m \times n_0}$  and  $y \in \mathbb{R}^{m \times 1}$ ,  $x_i$  and  $y_i$  are the features and target for the  $i^{\text{th}}$  data point.

$$l_D(\theta) = \frac{1}{m} \frac{1}{2} \|f_i(X, \theta) - y\|_2^2$$

$$\theta = \theta - \eta \nabla l(\theta) \quad \nabla l(\theta) = \begin{bmatrix} \frac{\partial l(\theta)}{\partial w^3} \\ \frac{\partial l(\theta)}{\partial w^2} \\ \frac{\partial l(\theta)}{\partial w^1} \end{bmatrix}$$

$l(\theta) = \frac{1}{2} \frac{1}{m} (h_{1,1}^3 - y_1)^2$





# GRADIENT DESCENT

## SET UP

$$\nabla l(\theta) = ?$$

$$\begin{aligned} \frac{\partial l(\theta)}{\partial h_{1,1}^3} &= (h_{1,1} - y_i) = \delta_1 \\ h_{1,1}^3 &= f_1^3(h^2, w^3) = h^2 w_{1,1}^3 \quad \begin{matrix} 1 \times h_2 & h_2 \times h_3 & 2 \times 1 \end{matrix} \\ \begin{bmatrix} h_{1,1}^2 & h_{1,2}^2 \end{bmatrix} & \begin{bmatrix} w_{1,1}^3 \\ w_{1,2}^3 \end{bmatrix} \\ \frac{\partial l(\theta)}{\partial w_{1,1}^3} &= \frac{\partial l(\theta)}{\partial h_{1,1}^3} \frac{\partial h_{1,1}^3}{\partial w_{1,1}^3} \end{aligned}$$

We need to know compute the loss function's gradient with respect to all the parameters in the neural network before running gradient descent.



# BACKPROP

## FORWARD PASS

Compute the outputs of each layer and the loss  $l_D(\theta)$

$$\begin{aligned}h^0 &= x \\h^1 &= f^1(h^0, W^1) \\h^2 &= f^2(h^1, W^2) \\\vdots \\h^i &= f^i(h^{i-1}, W^i) \\\vdots \\h^k &= f^k(h^{k-1}, W^k)\end{aligned}$$



# BACKPROP

## BACKWARD PASS

Using the results of the forward pass, apply the chain rule to compute the derivatives for each layer

Compute  $\frac{\partial l_D(\theta)}{\partial f(X, \theta)}$

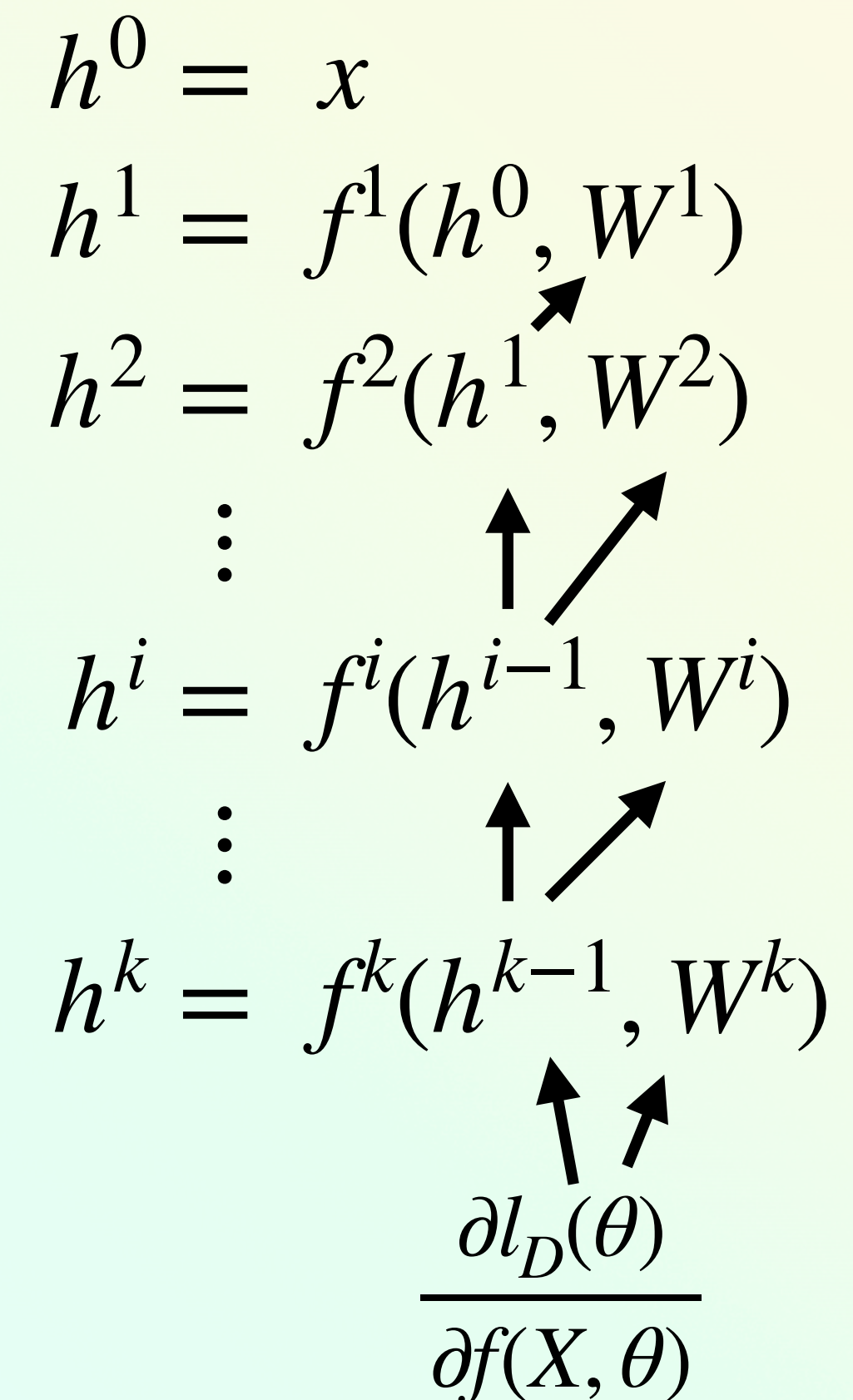
Then compute  $\frac{\partial l_D(\theta)}{\partial W^k}$  and  $\frac{\partial l_D(\theta)}{\partial h^{k-1}}$

Then compute  $\frac{\partial l_D(\theta)}{\partial W^{k-1}}$  and  $\frac{\partial l_D(\theta)}{\partial h^{k-2}}$

Repeat till  $W^1$

$$\begin{array}{l} h^0 = x \\ h^1 = f^1(h^0, W^1) \\ h^2 = f^2(h^1, W^2) \\ \vdots \\ h^i = f^i(h^{i-1}, W^i) \\ \vdots \\ h^k = f^k(h^{k-1}, W^k) \end{array}$$

$\frac{\partial l_D(\theta)}{\partial f(X, \theta)}$





# BACKPROP

## BACKWARD PASS

What are the partial derivatives for  $\frac{\partial l_D(\theta)}{\partial W^i}$  and  $\frac{\partial l_D(\theta)}{\partial h^i}$ ?



# NEXT CLASS

Next Class — Training Neural Networks