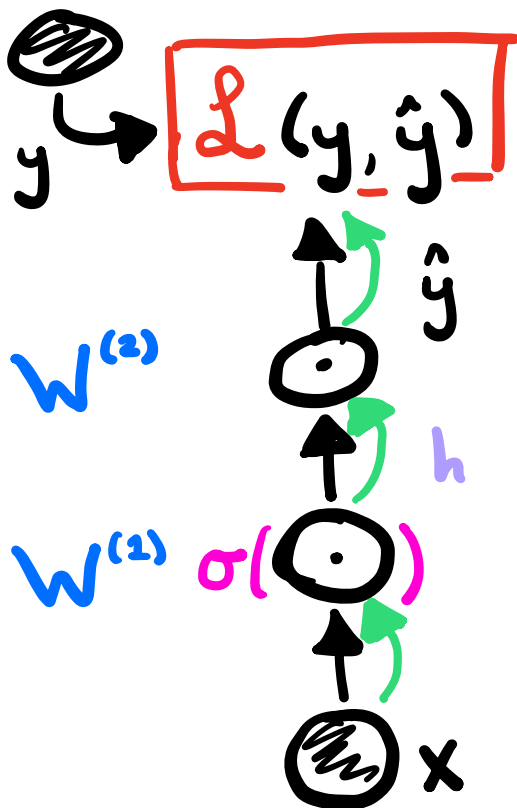


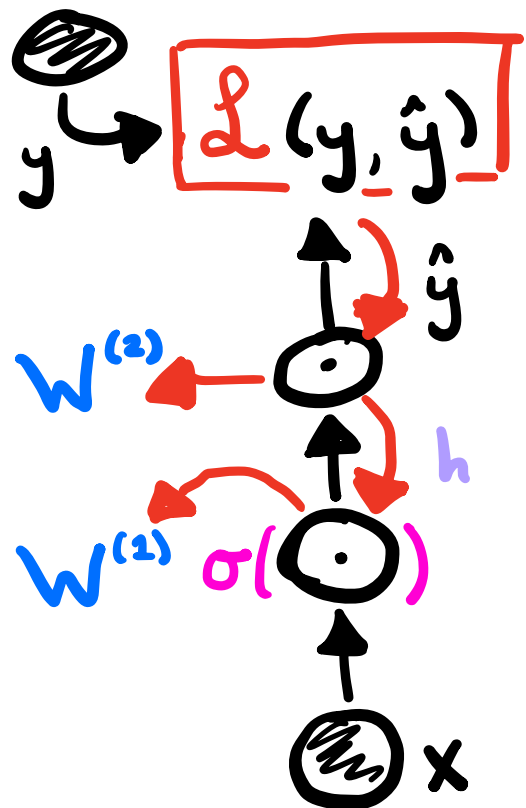
DS 4440

# Backpropagation (2)

Last time: Gradients on Computation graphs via backprop.



• forward()



• backward()

For implementation, each node/layer must know how to go forward and backward.

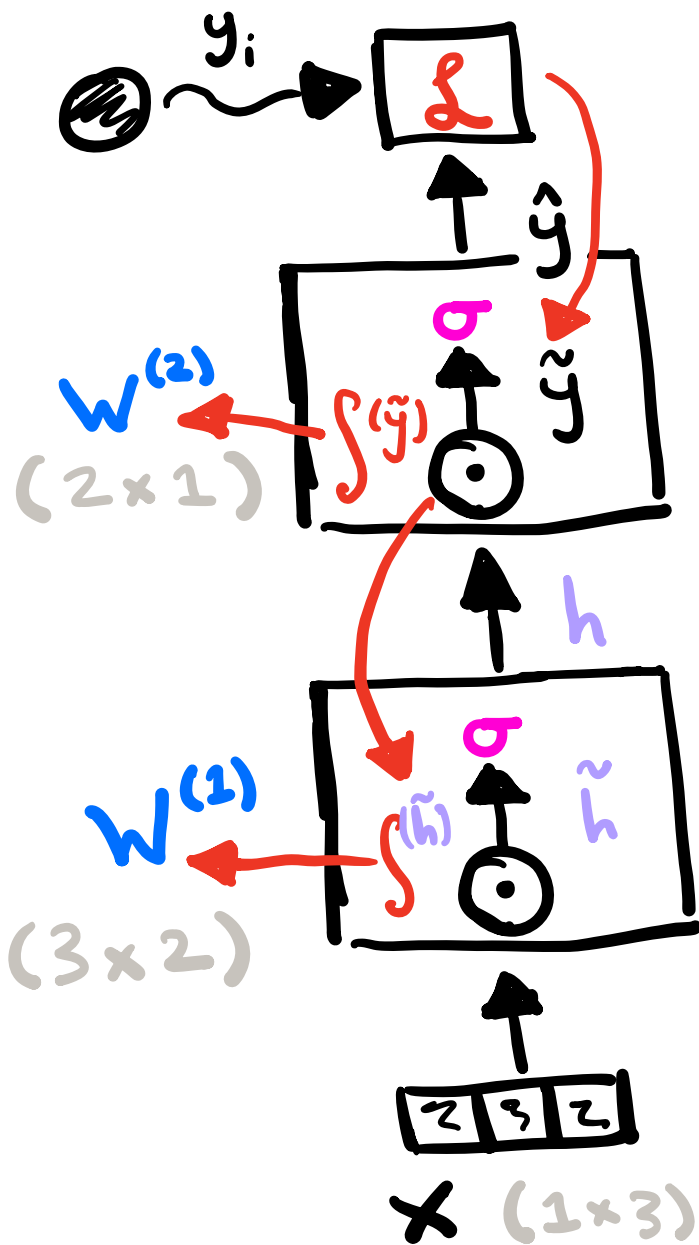
Let's consider an example

$$\mathcal{L} = \text{BCE Loss}(y, \hat{y})$$

$$= -y \log \hat{y} - (1-y) \log (1-\hat{y})$$

Probabilities

0/1



$$\frac{\partial \mathcal{L}}{\partial \tilde{y}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \tilde{y}}$$

"local error"  $\int^{(\hat{y})}$

$$\nabla_{W^{(2)}} \mathcal{L} = \nabla_{W^{(2)}} \tilde{y} \cdot \frac{\partial \mathcal{L}}{\partial \tilde{y}}$$

$$= \nabla_{W^{(2)}} \tilde{y} \int^{(\hat{y})}$$

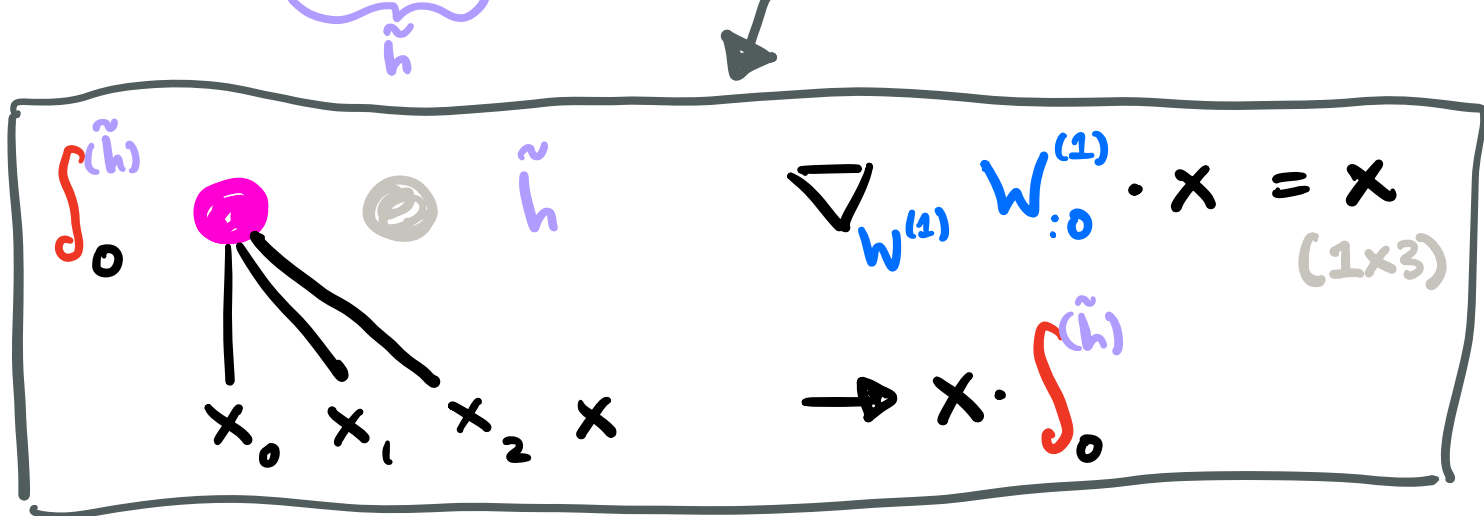
$$\nabla_{\tilde{h}} \mathcal{L} = \nabla_{\tilde{h}} \tilde{y} \int^{(\hat{y})}$$

local error act  $\tilde{h}$   $\int^{(\tilde{h})}$

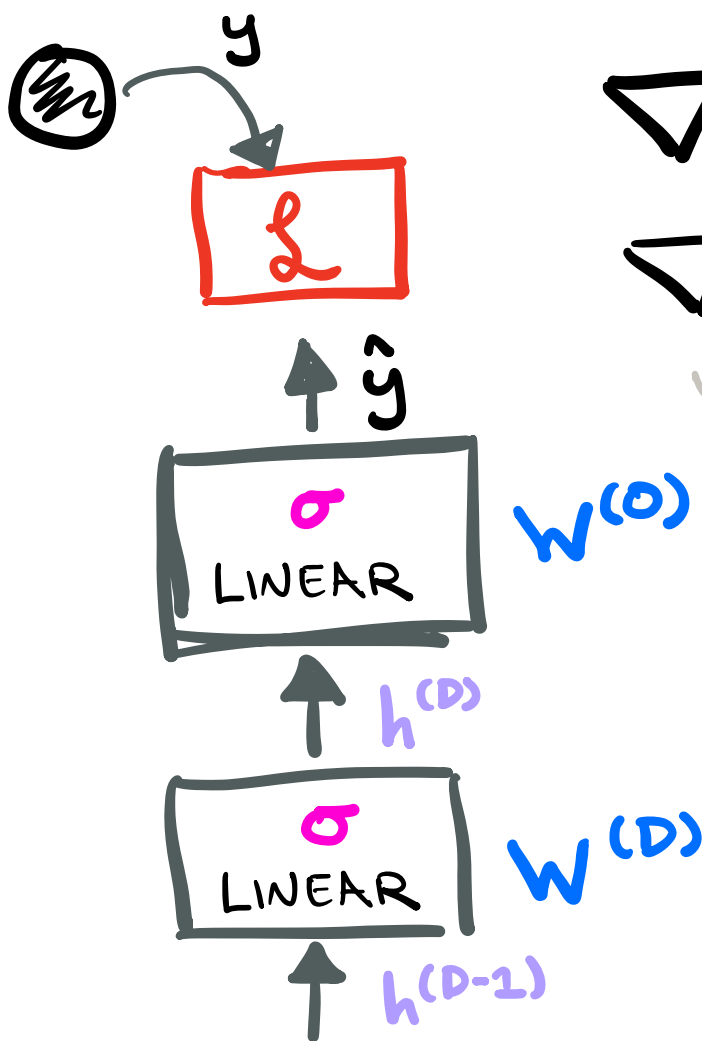
$$\nabla_{W^{(1)}} \mathcal{L} = \nabla_{W^{(1)}} \tilde{h} \cdot \int^{(\tilde{h})}$$

$$\hat{y} = \sigma(W^{(2)} \cdot h)$$

$$h = \sigma(W^{(1)} \cdot x)$$



Generalizing: Consider a feed forward network with  $D$  layers.

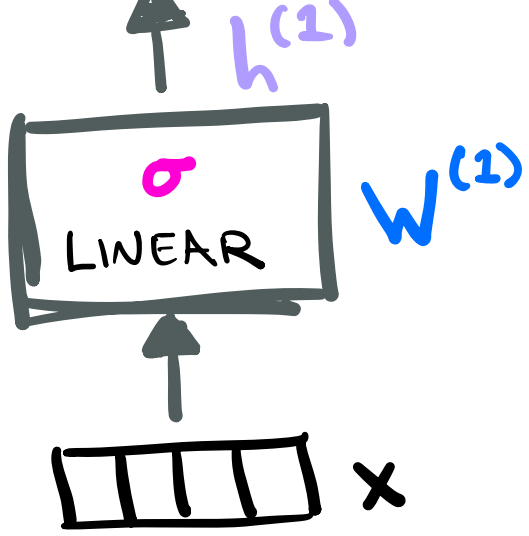


$$\nabla_{W^{(k)}} \mathcal{L} = \nabla_{W^{(k)}} \cdot \int h^{(k)}$$

$$\nabla_{h^{(k)}} \mathcal{L} = \nabla_{h^{(k)}} h^{(k+2)} \int h^{(k+2)}$$

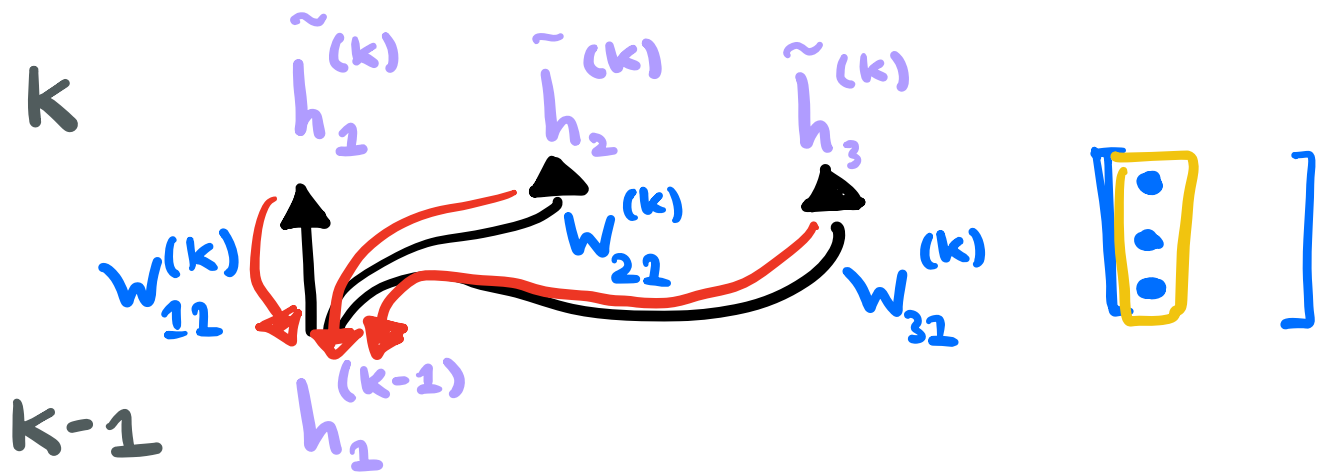
Matrices! Jacobian

$$\begin{bmatrix} \frac{\partial h_2^{(k+2)}}{\partial h_1^{(k)}} & \dots & \frac{\partial h_d^{(k+2)}}{\partial h_1^{(k)}} \\ \dots & \dots & \dots \\ \dots & \dots & \frac{\partial h_d^{(k+2)}}{\partial h_d^{(k)}} \end{bmatrix}$$



In a Simple fully Connected layer, The **local error** is INTUITIVE

$$h^{(k)} = \sigma \left( \underbrace{W^{(k)} \cdot h^{(k-1)}}_{\tilde{h}^{(k)}} \right)$$



$$\frac{\partial \tilde{h}_2^{(k)}}{\partial h_1^{(k-1)}} = \frac{\partial W_{22}^{(k)} h_1^{(k-1)}}{\partial h_1^{(k-1)}} = W_{22}^{(k)}$$

So "blame" for error at  $h_1^k$  on  $h_1^{(k-1)}$  Scales with Connection weight

$$\int_1^{(k-1)} = \nabla_{h_1^{(k-1)}} \int^k$$

$$= \sum_j \frac{\partial \tilde{h}_j^{(k)}}{\partial h_1^{(k-1)}} \left. \int^k \right\} \text{local error at } j$$

$$\frac{\partial h_1^{(k-1)}}{\partial h_1^{(k-1)}} \approx \sum_{j=1}^d W_{j1} \int_j^k$$

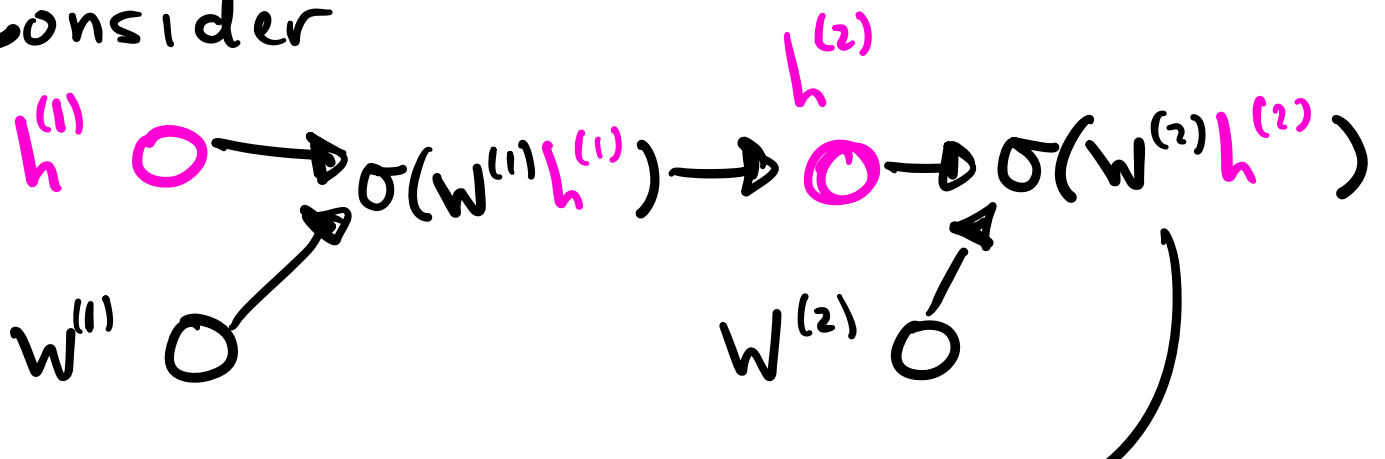
Influence of  $h_1^{(k-1)}$  on  $h_j^{(k)}$

ignoring  $\sigma$

Exercise! Let's implement a wacky custom layer (Colab)

Memory, Backprop & Detach

- What do we need to store for Backprop @ each layer?
- Consider





$$\frac{\partial}{\partial W^{(1)}} =$$

$$\frac{\partial}{\partial \hat{y}} L \cdot \frac{\partial}{\partial h^{(2)}} \hat{y} \cdot \frac{\partial}{\partial h^{(1)}} h^{(2)} \cdot \frac{\partial}{\partial W^{(1)}} h^{(1)}$$

Note that if we freeze lower layers we can discard corresponding ACTIVATIONS.