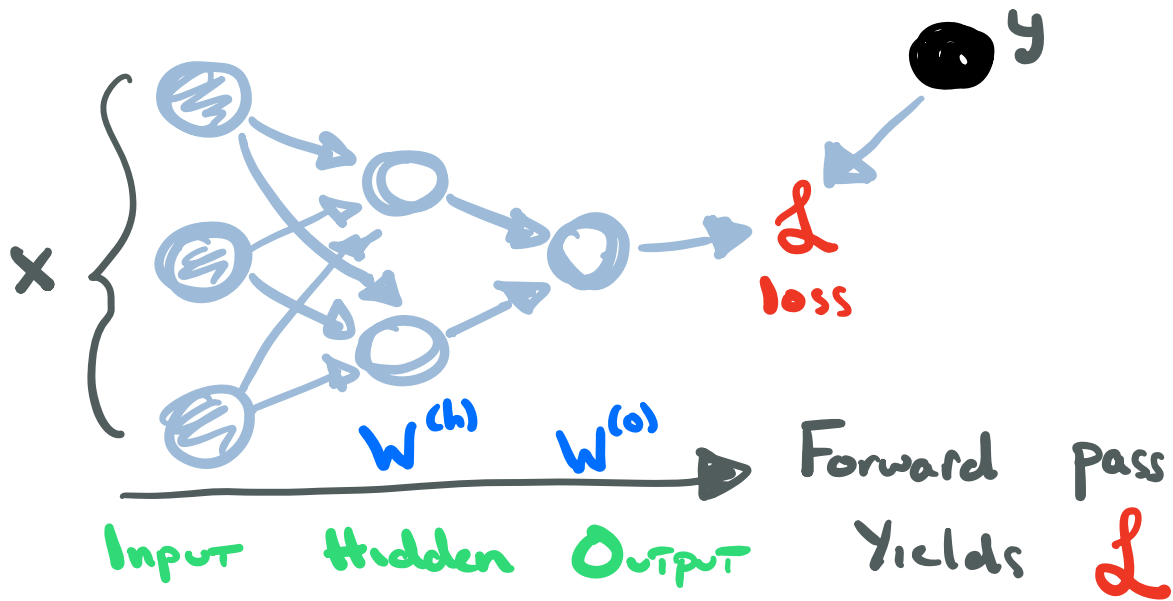


DS 4440 Backpropagation (1)

Last Time Computation graphs comprise layers to form networks

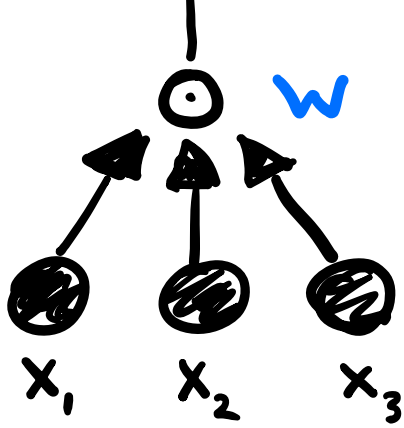


Today: How to adjust parameters in computation graphs to minimize \mathcal{L}

Backprop. This extends SGD (L1; HW2).

Consider the graph for logistic regression.





$$\nabla_w \mathcal{L} = (y - \sigma(\hat{w} \cdot x))x$$

for one instance (x, y) .

The Recipe

- Build graph (model) w/ parameters W
- Define (differentiable) \mathcal{L}
- Use $\nabla_w \mathcal{L}$ to find \hat{W} via SGD

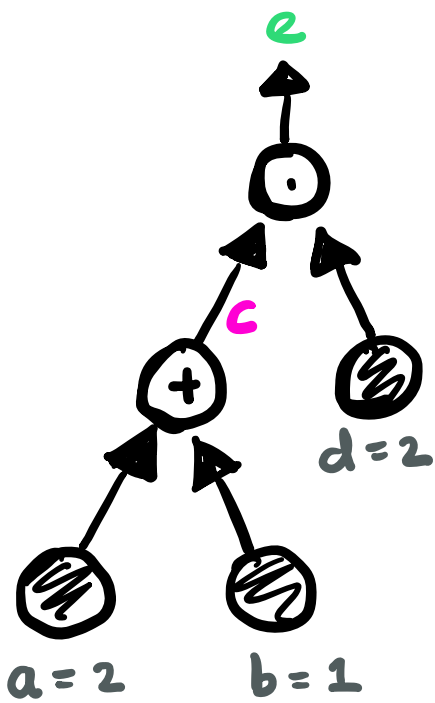
Problem Finding ∇_w is painful to do manually for big models.

Solution Use the computation graph to auto-diff via backprop.

Consider a simple scalar example

$$e = (a + b) \cdot d$$

How is e affected by a ?



Chain Rule!

$$\frac{d}{dx} f(g(x)) = \frac{d}{du} f(u) \cdot \frac{du}{dx}$$

$$\frac{\partial e}{\partial a} = \frac{\partial e}{\partial c} \cdot \frac{\partial c}{\partial a} = 2 \cdot 1 = 2$$

$$\frac{d}{dx} u$$

$$\frac{\partial}{\partial a} (a+b) = 1$$

$$\frac{\partial}{\partial c} c \cdot d = d = 2$$

$$\frac{\partial e}{\partial b} = \frac{\partial e}{\partial c} \frac{\partial c}{\partial b}$$

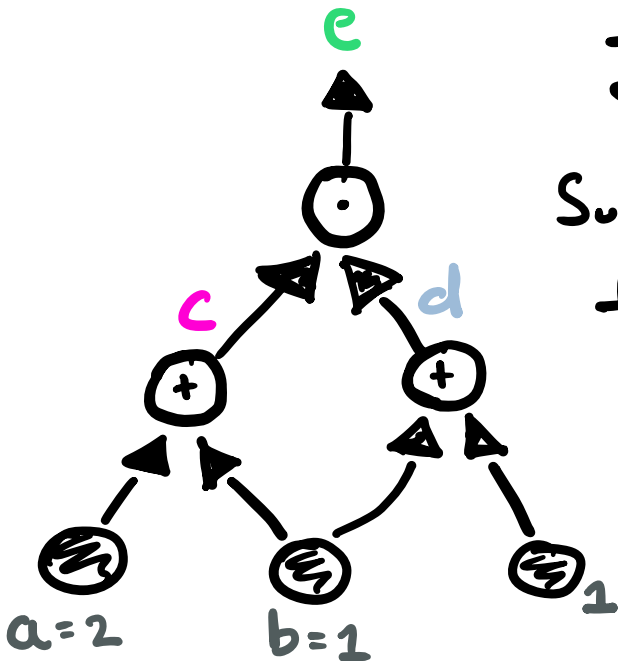
Re-use!

A Slightly Trickier example

$$e = (a+b)(b+1)$$

$$\frac{\partial}{\partial b} e = ?$$

Sum over multiple paths from $b \rightarrow e$



$$= \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial e}{\partial d} \frac{\partial d}{\partial b}$$

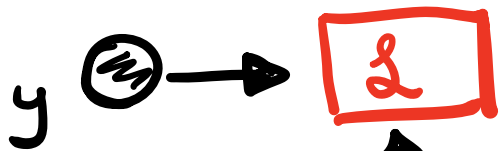
$$= (b+2) + (a+b)$$

$$= 5$$

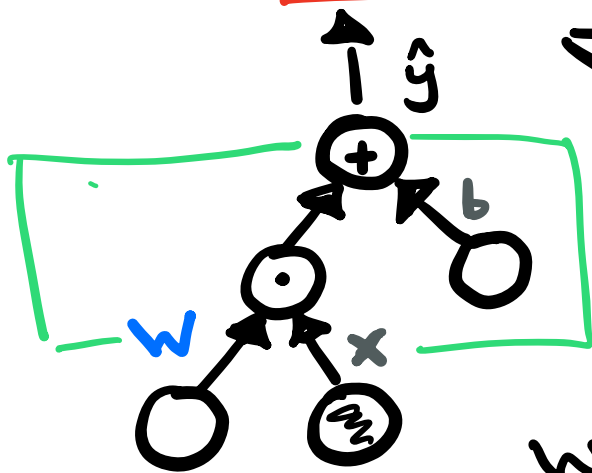
Let's add a **Loss**. Assume regression.

$$\hat{y} = \hat{w}x + \hat{b}$$

$$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2 = (y - [\hat{w}x + \hat{b}])^2$$



We previously derived ∇_w for this **single layer** model.

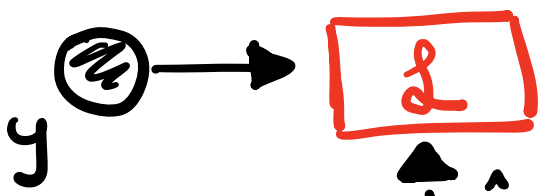


More interesting when we introduce

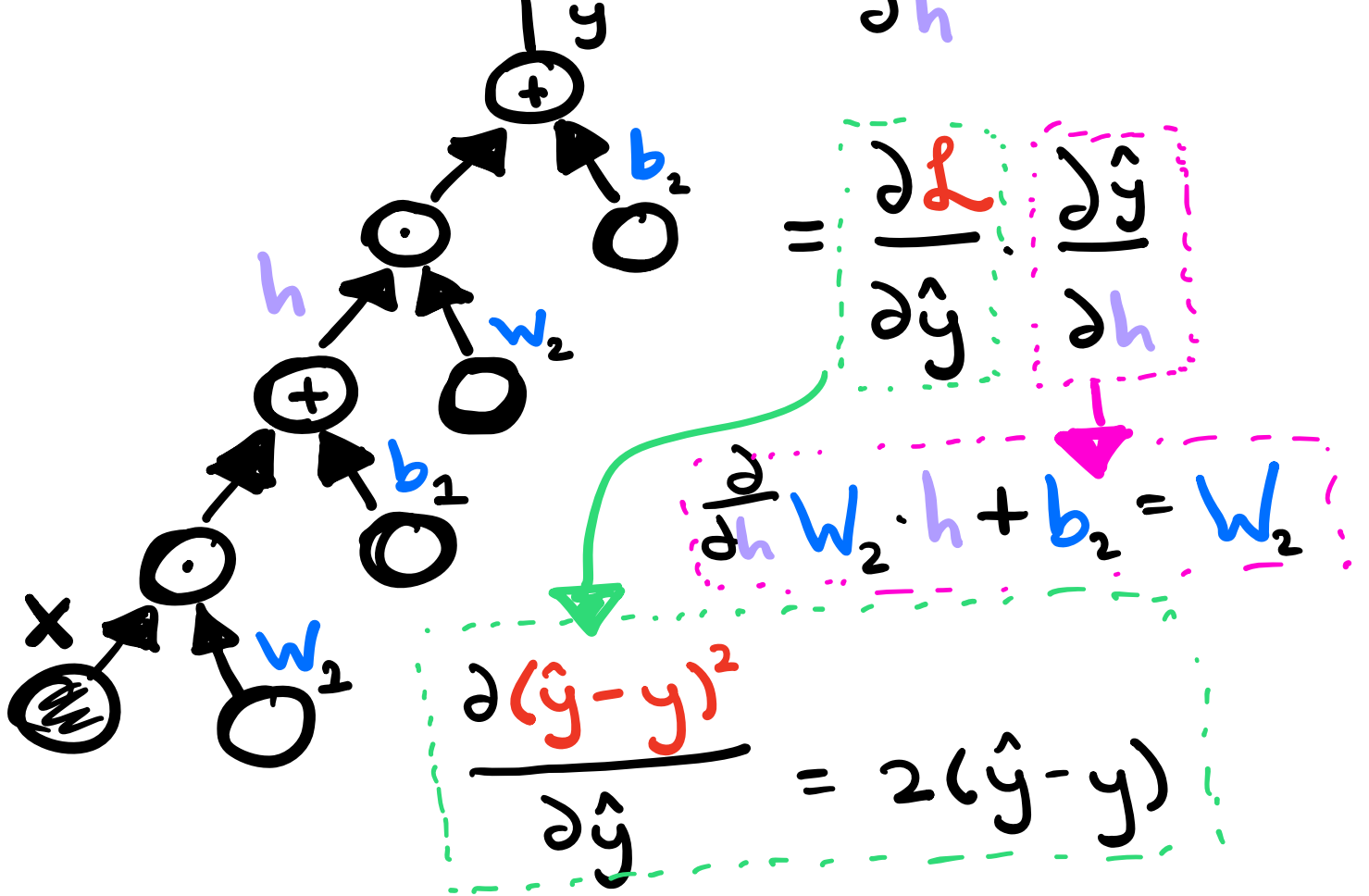
hidden layers.

(See in-class exercise!)

$$\hat{y} = w_2 \cdot \underbrace{(w_1 \cdot x + b_1)}_h + b_2$$



$$\frac{\partial \mathcal{L}}{\partial w} = ?$$



$$= 2(\hat{y} - y) w_2.$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial h}{\partial w_2} \cdot \frac{\partial \mathcal{L}}{\partial h} = x \cdot 2(\hat{y} - y) w_2$$

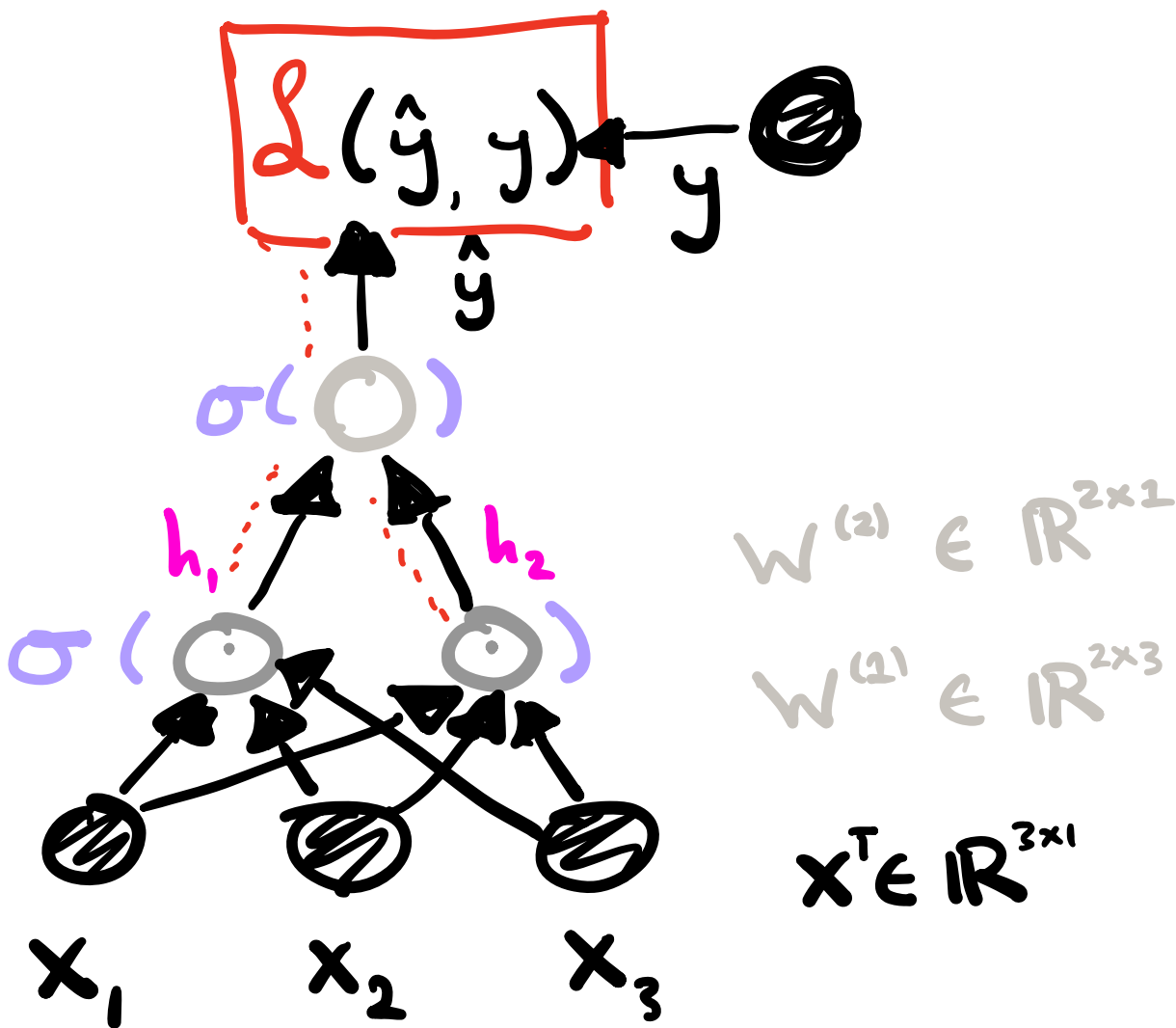
$$\frac{\partial}{\partial w_1} w_1 \cdot x + b_1 = x$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{\partial h}{\partial b_1} \cdot \frac{\partial \mathcal{L}}{\partial h} = 2(\hat{y} - y) w_2$$

$$\frac{\partial}{\partial 1} = 1$$

(See Calc!

How about This NN?



$\nabla_{W^{(2)}} \mathcal{L}$ We have already seen

$$\nabla_{W^{(1)}} \mathcal{L} = ?$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial h} \cdot \frac{\partial h}{\partial W^{(1)}} \quad | \quad \text{Chain rule!}$$

[1]
Scalar

$$\begin{bmatrix} \frac{\partial \hat{y}}{\partial h_1} \\ \frac{\partial \hat{y}}{\partial h_2} \end{bmatrix}$$

δ_1 δ_2

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

$\nabla_{W^{(1)}} h \approx x$ | Ignoring element wise σ

Local error at h_1 and h_2 , after

- with $\frac{\partial \mathcal{L}}{\partial \hat{y}}$.

$$\begin{bmatrix} \frac{\partial h_1}{\partial W^{(1)}} \\ \frac{\partial h_2}{\partial W^{(1)}} \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 \\ x_1 & x_2 & x_3 \end{bmatrix}$$

$$= \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \begin{bmatrix} \frac{\partial \hat{y}}{\partial h_1} x_1 & \frac{\partial \hat{y}}{\partial h_1} x_2 & \frac{\partial \hat{y}}{\partial h_1} x_3 \\ \frac{\partial \hat{y}}{\partial h_2} x_1 & \frac{\partial \hat{y}}{\partial h_2} x_2 & \frac{\partial \hat{y}}{\partial h_2} x_3 \end{bmatrix}$$

$$= \begin{bmatrix} \delta_1 \vec{x} \\ \delta_2 \vec{x} \end{bmatrix}$$

More Next Time!