

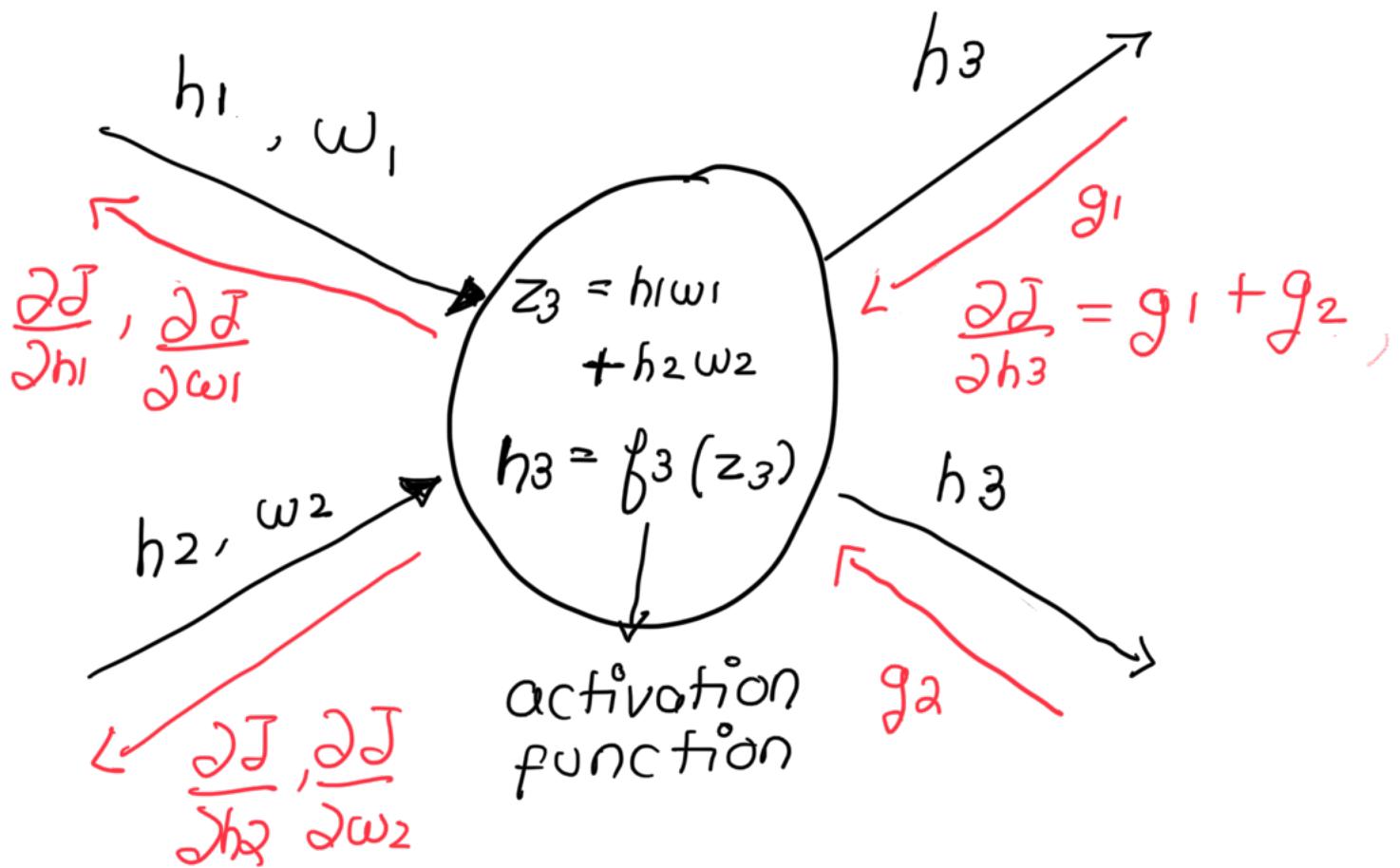
Backpropagation

The goal of backpropagation is to compute $\frac{\partial J}{\partial \omega}$ for every weight ω in the network and perform the update

$$\omega = \omega - \ell \frac{\partial J}{\partial \omega}$$

so that on the next forward pass, the loss J will reduce.

Let's zoom in and observe what BP does at a particular neuron in the network:



where

$$\frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_1}$$

$$\frac{\partial J}{\partial \omega_1} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial \omega_1}$$

$$\frac{\partial J}{\partial h_2} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_2}$$

$$\frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial w_2}$$

Also keep in mind that

$$\frac{\partial h_i}{\partial a} \quad \text{for some } i \text{ and variable } a$$

$$= \frac{\partial h_i}{\partial z_i} \cdot \frac{\partial z_i}{\partial a}$$

because of activations, but for the sake of keeping the chains shorter, I won't always expand

Now, we've setup a procedure we can follow at every neuron that computes

$$h' = f'(z')$$

f' is the activation and z' is the

$v = \omega_1 h_1 + \dots + \omega_K h_K$ ("" "
weighted sum of outputs from the
previous layer)

1. If h' reaches the final output through multiple paths, it will accumulate multiple sub gradients' g^1, \dots, g_j . Add these up, this gives
$$g_1 + \dots + g_j = \frac{\partial J}{\partial h},$$

(You can think of each g_i as one of the summands when you compute $\frac{\partial J}{\partial h}$, using the chain rule)

2. Compute the partial derivative of h' w.r.t every input $\omega_1, h_1, \dots, \omega_K, h_K$, the gradients

(we also need " $\frac{\partial J}{\partial \omega}$ " because we root h_1, \dots, h_K because they are required to calculate $\frac{\partial J}{\partial \omega}$ for weights ω in earlier layers)

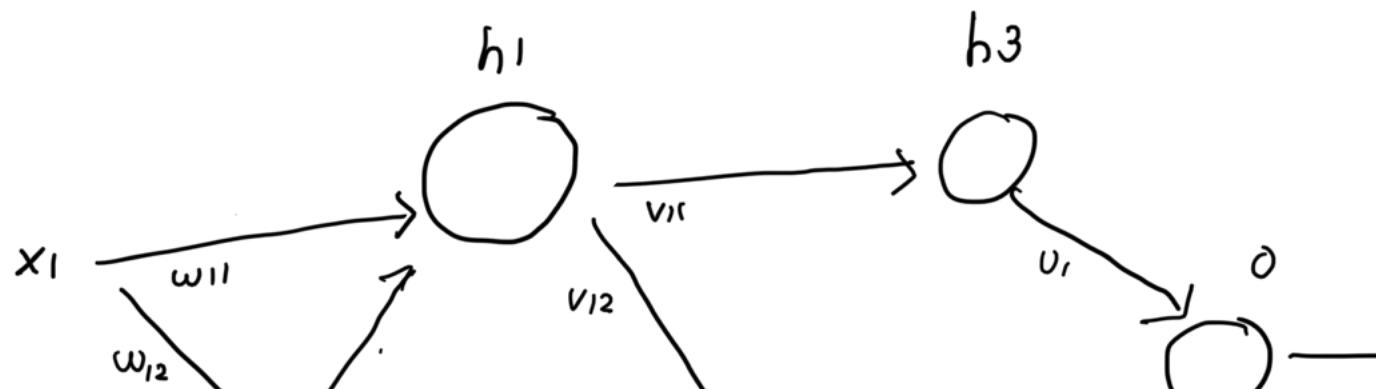
- Applying the chain rule for every input!

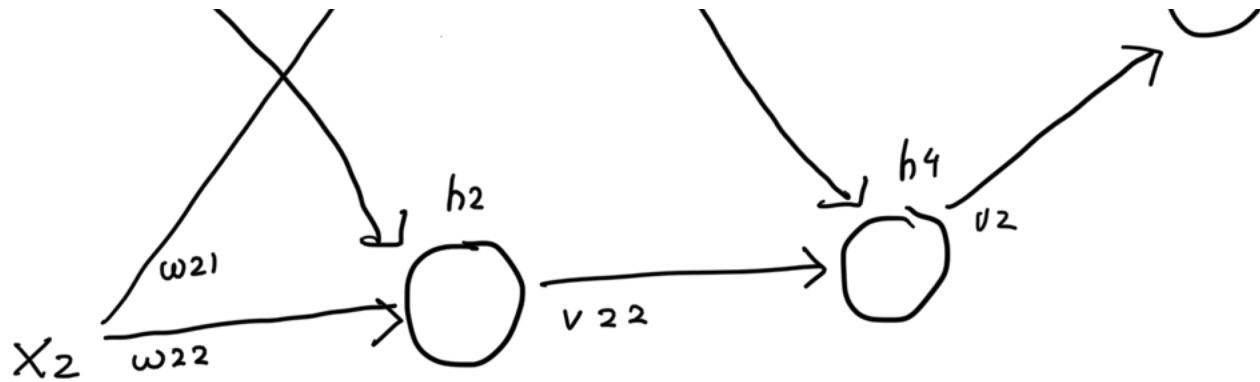
$$\text{Ex: } \underbrace{\frac{\partial J}{\partial h_i}}_{\text{outgoing gradient}} = \underbrace{\frac{\partial J}{\partial h'}}_{\text{incoming, precomputed gradient}} \circ \underbrace{\frac{\partial h'}{\partial h_i}}_{\text{local gradient computed in step 2}}$$

Notice how we are performing the chain rule incrementally, one chain 'link' at a time. This is exactly why backprop is efficient - it isn't repeatedly

- UV calculating long chains for every variable because it doesn't need to. The compositional nature of neural networks means that 'sub' chains computed at later layers can be reused to compute the longer chains at earlier layers, meaning HUGE savings on computation overall, especially for large, wide, denser nets

Let's do an example:





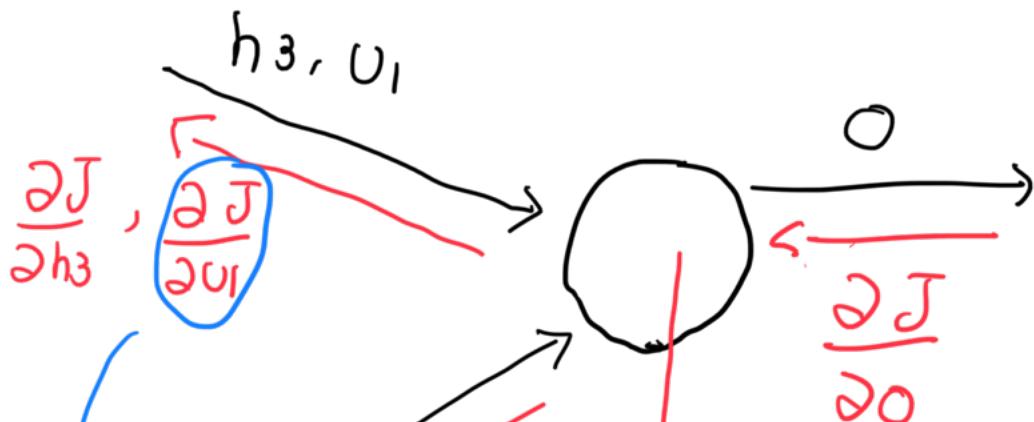
$h_i^o = f_i^o(z_i^o)$ where $z_i^o = \text{weighted sum of inputs}$

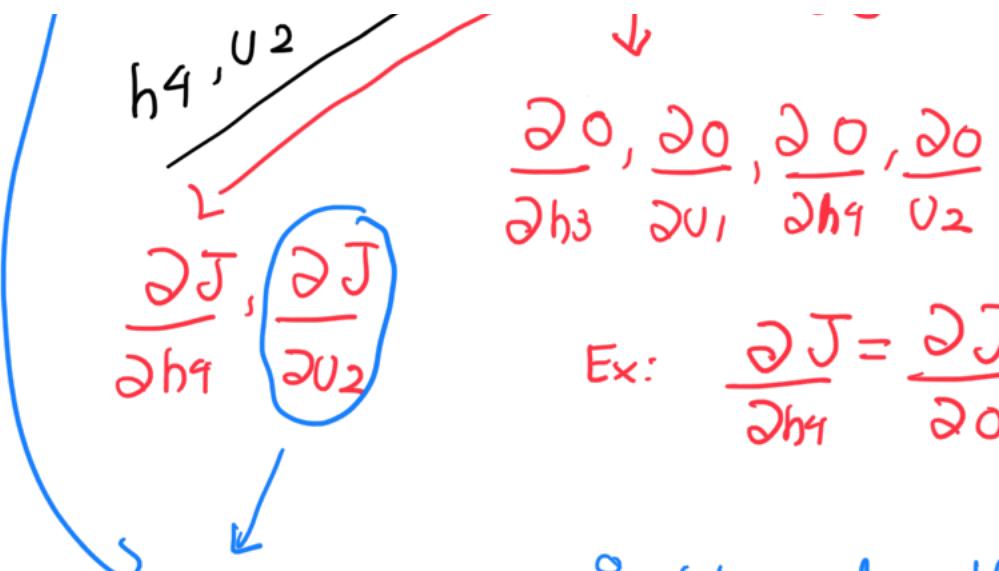
$o = f_o(z_o)$ where $z_o = v_1 h_3 + v_2 h_4$

o = predictions of the network

With o , we can compute our loss J since we have both the targets and predictions. And then as our first step in back prop, we

compute $\frac{\partial J}{\partial o}$





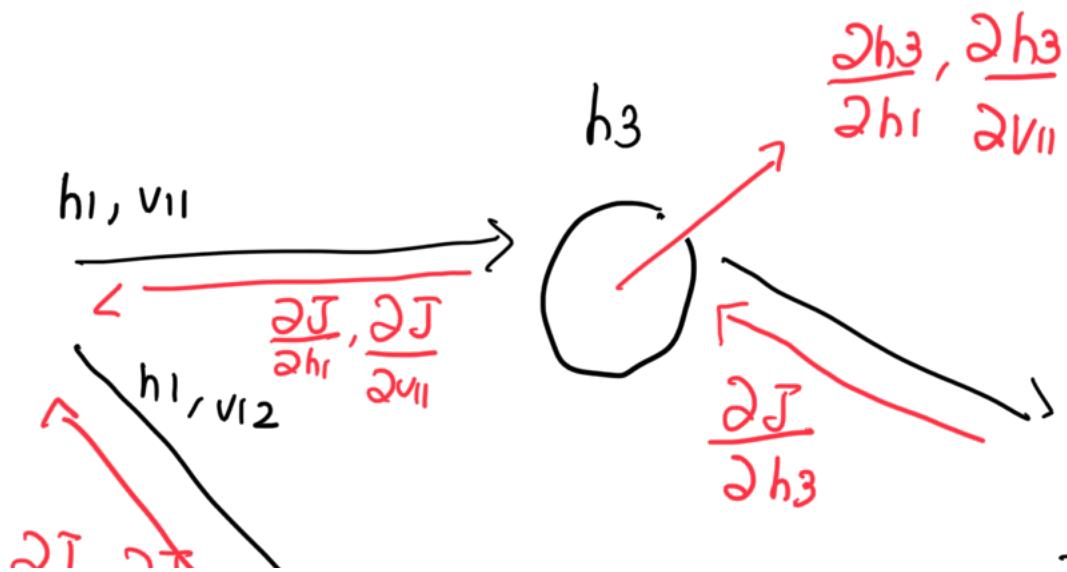
$$\text{Ex: } \frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial o} \cdot \frac{\partial o}{\partial h_1}$$

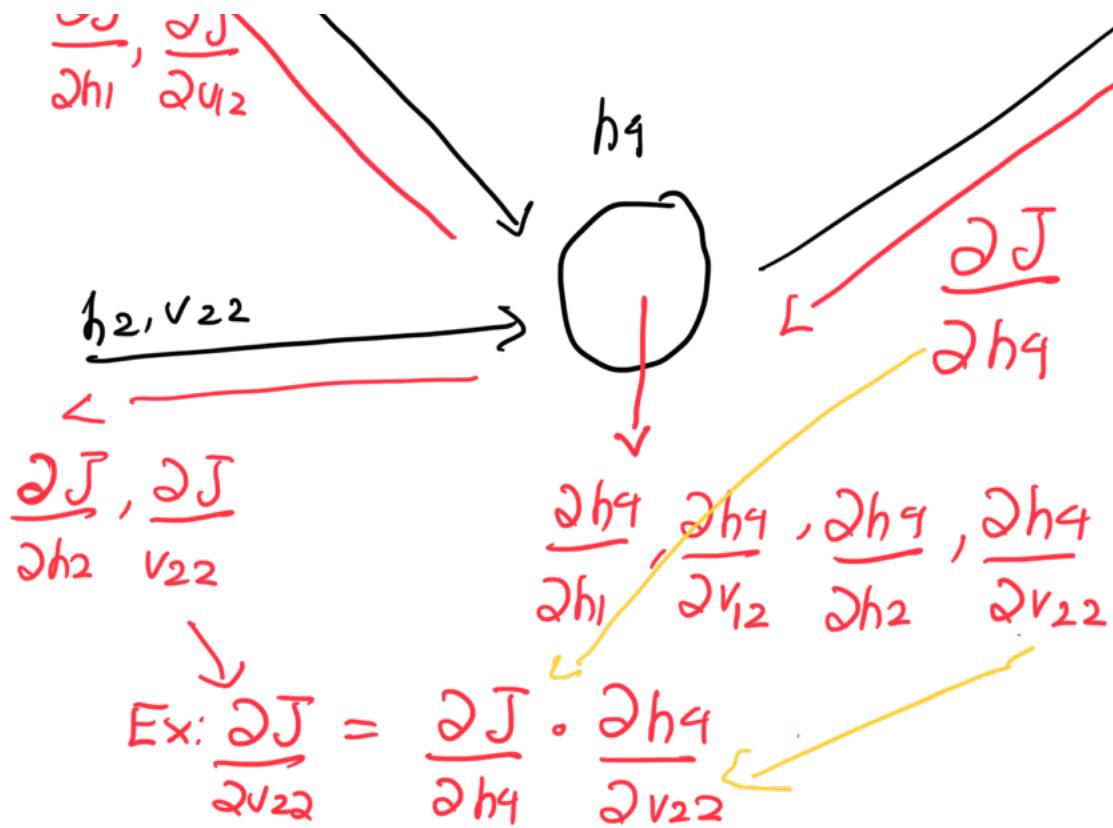
These are weights, do the gradient descent update now!

$$v_1 = v_1 - \ell \frac{\partial J}{\partial v_1}$$

$$v_2 = v_2 - \ell \frac{\partial J}{\partial v_2}$$

Now we move backward to the previous layer ...

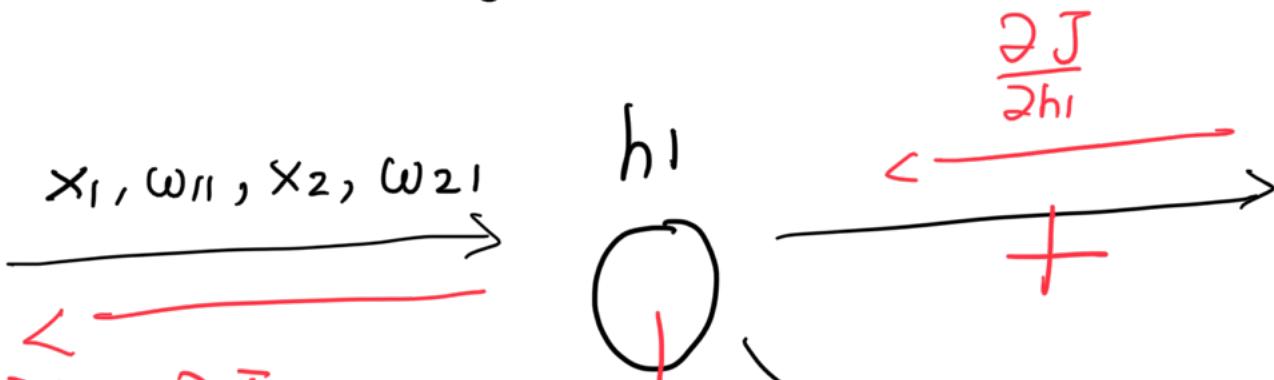




We've calculated $\frac{\partial J}{\partial v_{11}}, \frac{\partial J}{\partial v_{12}}, \frac{\partial J}{\partial v_{22}}$.

Now we can do the gradient descent update for those

Last layer! (combining input arrows)



$$\frac{\partial J}{\partial w_{11}}, \frac{\partial J}{\partial w_{21}}$$

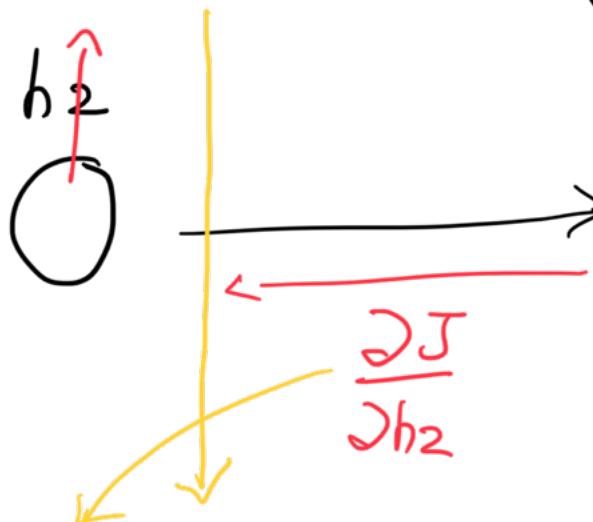
$$\frac{\partial h_1}{\partial w_{11}}, \frac{\partial h_1}{\partial w_{21}}$$

$$\frac{\partial J}{\partial h_1}$$

$$\frac{\partial h_2}{\partial w_{12}}, \frac{\partial h_2}{\partial w_{22}}$$

$$x_1, w_{12}, x_2, w_{22}$$

$$\frac{\partial J}{\partial w_{12}}, \frac{\partial J}{\partial w_{22}}$$



$$\text{Ex: } \frac{\partial J}{\partial w_{22}} = \frac{\partial J}{\partial h_2} \circ \frac{\partial h_2}{\partial w_{22}}$$

Notice a few things:

- 1) There are 2 $\frac{\partial J}{\partial h_i}$'s coming into h_1 . These are the 'sub' gradients I was referring to earlier. The actual, original $\frac{\partial J}{\partial h_i}$ would be the

$\sum \text{outputs} \frac{\partial J}{\partial h_i}$

Sum of these and what
ultimately gets back propagated
/ used. Ex: when calculating

$$\frac{\partial J}{\partial w_{ii}} = \frac{\partial J}{\partial h_i} \cdot \frac{\partial h_i}{\partial w_{ii}}$$

I didn't relabel these sub
gradients because it keeps
things cleaner and easier
to follow. But remember,
neurons with more than 1
channel of output accumulate
/ add multiple subgradients
to compute a final gradient
that is shared with earlier
layers

→ we didn't calculate $\frac{\partial J}{\partial}$

δ_j for x_1, x_2 . Because we don't have to! They would be useless - there are no earlier layers and it doesn't make sense to adjust inputs (we don't have control over that, duh)

That is one complete backward pass! We could compute $\frac{\partial J}{\partial \omega}$ for all weights ω in the network. Next, we would repeat the cycle of doing a forward pass of predicting with the updated weights, backpropagating ...

the gradients and so on till convergence (weights stabilize)

To show that backprop is correct and efficient, consider what the formula for $\frac{\partial J}{\partial w_{11}}$ would be.

w_1 influences o through 2 paths :

$$1) w_1 \rightarrow h_1 \rightarrow h_3 \rightarrow o$$

$$2) w_1 \rightarrow h_1 \rightarrow h_4 \rightarrow o$$

so we'll have 2 summands in the formula for $\frac{\partial J}{\partial w_{11}}$:

$$\underline{\frac{\partial J}{\partial w_{11}}} = \underline{\frac{\partial J}{\partial o}} \underline{\frac{\partial o}{\partial h_3}} \underline{\frac{\partial h_3}{\partial h_1}} +$$

$$\frac{\partial J}{\partial w_{11}} = \frac{\partial J}{\partial o} \frac{\partial o}{\partial h_3} \frac{\partial h_3}{\partial h_1} \frac{\partial h_1}{\partial w_{11}}$$

$$\frac{\partial J}{\partial w_{11}} = \frac{\partial J}{\partial h_4} \frac{\partial h_4}{\partial h_1} \frac{\partial h_1}{\partial w_{11}}$$

Now let's see if we get this from our backprop calculation. Continuously unavelling the chain:

$$\frac{\partial J}{\partial w_{11}} = \frac{\partial J}{\partial h_1} \cdot \frac{\partial h_1}{\partial w_{11}}$$

||

$$\left(\frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_1} + \frac{\partial J}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_1} \right) \frac{\partial h_1}{\partial w_{11}}$$

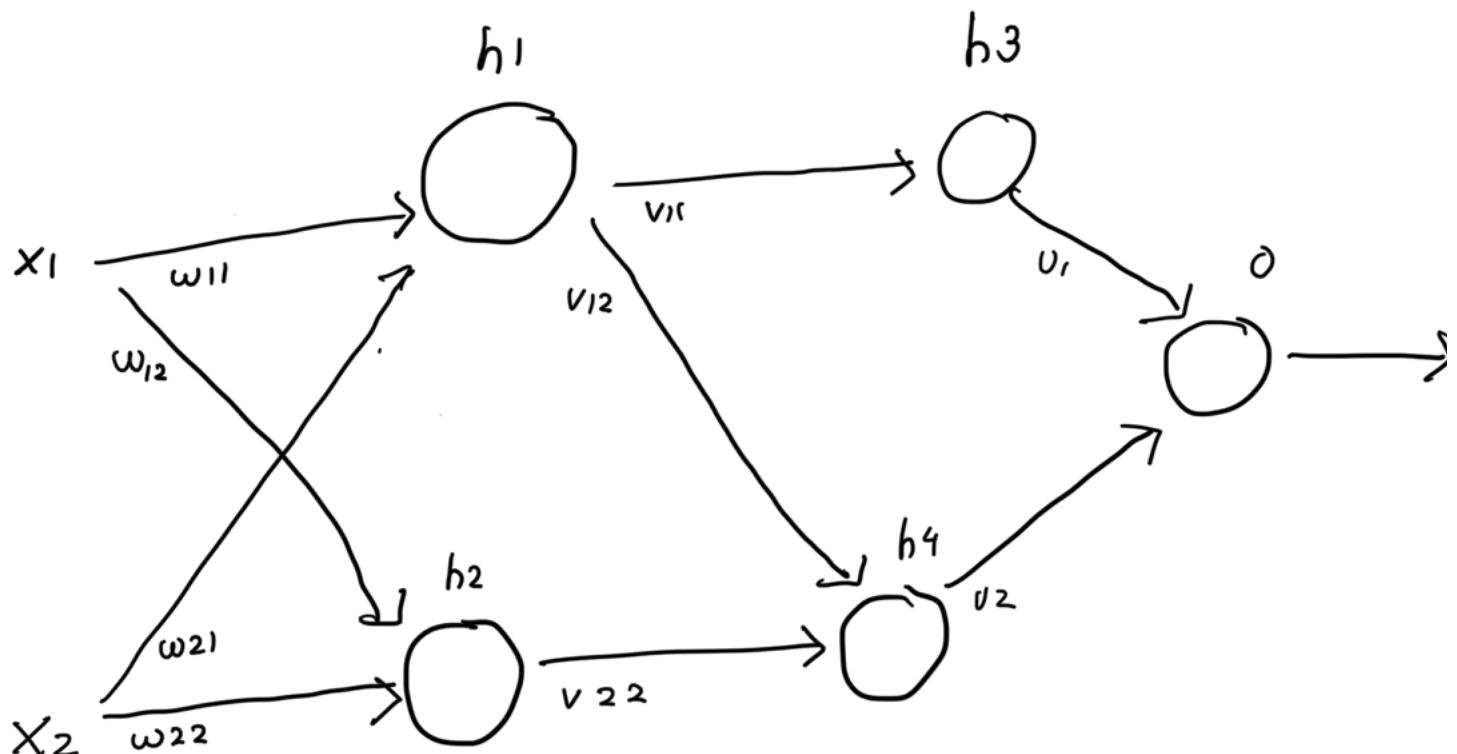
$$= \left(\underbrace{\frac{\partial J}{\partial o} \frac{\partial o}{\partial h_3} \frac{\partial h_3}{\partial h_1}}_{\text{J}} + \underbrace{\frac{\partial J}{\partial o} \frac{\partial o}{\partial h_4} \frac{\partial h_4}{\partial h_1}}_{\text{J}} \right) \frac{\partial h_1}{\partial w_{11}}$$

$$= \frac{\partial J}{\partial h_3} \frac{\partial o}{\partial h_3} \frac{\partial h_3}{\partial h_1} \frac{\partial h_1}{\partial w_{11}} + \frac{\partial J}{\partial h_4} \frac{\partial o}{\partial h_4} \frac{\partial h_4}{\partial h_1} \frac{\partial h_1}{\partial w_{11}}$$

and voila !

Now let's do an example
with actual numbers and
functions

Reminder about the network:



$$h_i^o = f_i^o(z_i^o)$$

$z_i^o = \text{weighted sum of inputs}$

Define

$$h_i = f_i(x) = x$$

$$h_2 = f_2(x) = \frac{1}{2}x^2$$

$$h_3 = f_3(x) = 2x$$

$$h_4 = f_4(x) = e^x$$

$$O = f_0(x) = \sin(x)$$

And assume all weights are initialized to 0, i.e.

$$w_{11} = w_{12} = w_{21} = w_{22} = v_{11} = v_{12} = v_{22} = v_1 = v_2 = 0$$

Let's do one cycle of a forward and backward pass with the training instance

$$(x_1, x_2, y) = (0, 1, 0)$$

Step 1: forward pass

Calculate:

$$h_1 =$$

1 -

$h_2 =$

$h_3 =$

$h_4 =$

$o =$

Step 2 : Compute loss

Let's define

$$J \text{ as } \frac{1}{2} (o - t)^2$$

where t is the target
and o is the output of
the network

Calculate

$$\frac{\partial J}{\partial o} = o - t$$

and plug in t, o

Step 3 : Back propagate

Layer 3

o

(Remember to plug in values you've already computed, including from the forward pass)

Calculate :

$$\frac{\partial o}{\partial h_3} = \frac{\partial f_o}{\partial z_o} \cdot \frac{\partial z_o}{\partial h_3} = \cos(z_o) u_1$$

$$\frac{\partial o}{\partial h_4} = \frac{\partial f_o}{\partial z_o} \cdot \frac{\partial z_o}{\partial h_4} = \cos(z_o) u_2$$

$$\frac{\partial o}{\partial u_1} = \frac{\partial f_o}{\partial z_o} \cdot \frac{\partial z_o}{\partial u_1} = \cos(z_o) h_3$$

$$\frac{\partial \underline{O}}{\partial U_2} = \frac{\partial f_0}{\partial z_0} \cdot \frac{\partial z_0}{\partial U_2} = (\cos(z_0)) h^4$$

$$\frac{\partial J}{\partial h^3} = \frac{\partial J}{\partial O} \circ \frac{\partial O}{\partial h^3} = (t - o) \cos(z_0) U_1$$

$$\frac{\partial J}{\partial h^4} = \frac{\partial J}{\partial O} \circ \frac{\partial O}{\partial h^4} = (t - o) \cos(z_0) U_2$$

$$\frac{\partial J}{\partial U_1} = \frac{\partial J}{\partial O} \circ \frac{\partial O}{\partial U_1} = (t - o) \cos(z_0) h^3$$

$$\frac{\partial J}{\partial U_2} = \frac{\partial J}{\partial O} \circ \frac{\partial O}{\partial U_2} = (t - o) \cos(z_0) h^4$$

Layer 2

h_3

Calculate:

$$\partial h_3 = \partial f_3 \cdot \partial z_3 = 2 \cdot v_{11}$$

$$\frac{\partial \underline{h_1}}{\partial h_1} \quad \frac{\partial \underline{z_3}}{\partial z_3} \quad \frac{\partial \underline{h_1}}{\partial h_1}$$

$$\frac{\partial h_3}{\partial v_{II}} = \frac{\partial f_3}{\partial z_3} \cdot \frac{\partial z_3}{\partial v_{II}} = 2 \cdot h_1$$

$$\frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial h_1} = 2(t-o) \cos(z_0) u_1 v_{II}$$

$$\frac{\partial J}{\partial v_{II}} = \frac{\partial J}{\partial h_3} \cdot \frac{\partial h_3}{\partial v_{II}} = 2 \cdot h_1 (t-o) \cos(z_0) u_1$$

h_4

Calculate:

$$\frac{\partial J}{\partial h_1} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_1} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial f_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial h_1}$$

$$\frac{\partial J}{\partial h_2} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_2} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial f_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial h_2}$$

$$\frac{\partial J}{\partial h_3} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial h_4}{\partial h_3} = \frac{\partial J}{\partial h_4} \cdot \frac{\partial f_4}{\partial z_4} \cdot \frac{\partial z_4}{\partial h_3}$$

$$\partial v_{12} \quad \partial h^4 \quad \partial v_{12} \quad \partial h^4 \quad \partial z^4 \quad \partial v_{12}$$

$$\frac{\partial J}{\partial v_{22}} = \frac{\partial J}{\partial h^4} \cdot \frac{\partial h^4}{\partial v_{22}} = \frac{\partial J}{\partial h^4} \cdot \frac{\partial h^4}{\partial z^4} \cdot \frac{\partial z^4}{\partial v_{22}}$$

Add up the $\frac{\partial J}{\partial h_i}$'s you

calculated and treat this
as your new $\frac{\partial J}{\partial h_1}$

Layer 3

h_1

Calculate

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

$\sigma \sim$

$$\underline{\partial J} =$$

$$2\omega_{12}$$

$$\underline{h_2}$$

Calculate

$$\underline{\partial J} =$$

$$2\omega_{21}$$

$$\underline{\partial J} =$$

$$2\omega_{22}$$

Finally, do all the
constraint updates using

$$\omega_{\text{new}} = \omega - \ell \frac{\partial J}{\partial \omega}$$

what are the new values
of

$$\omega_{11} =$$

$$\omega_{12} =$$

$$\omega_{21} =$$

$$\omega_{22} =$$

$$v_{11} =$$

$$v_{12} =$$

$$v_{22} =$$

$$u_1 =$$

$$u_2 =$$

v <