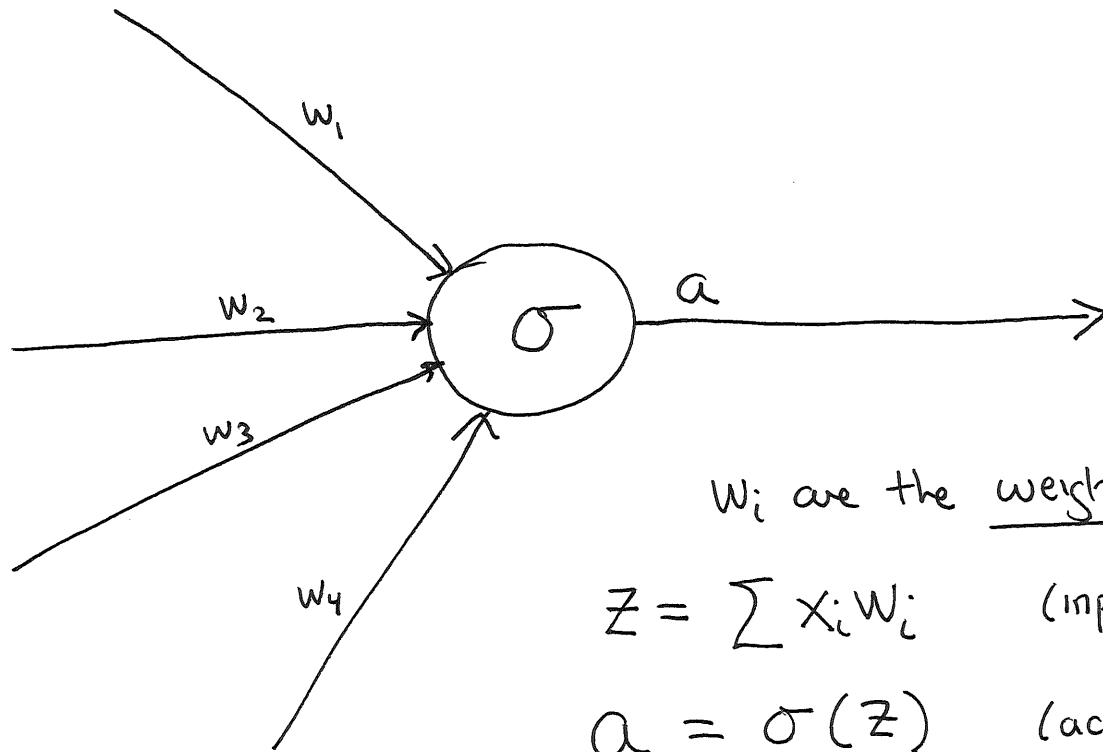


# A Neuron



$w_i$  are the weights

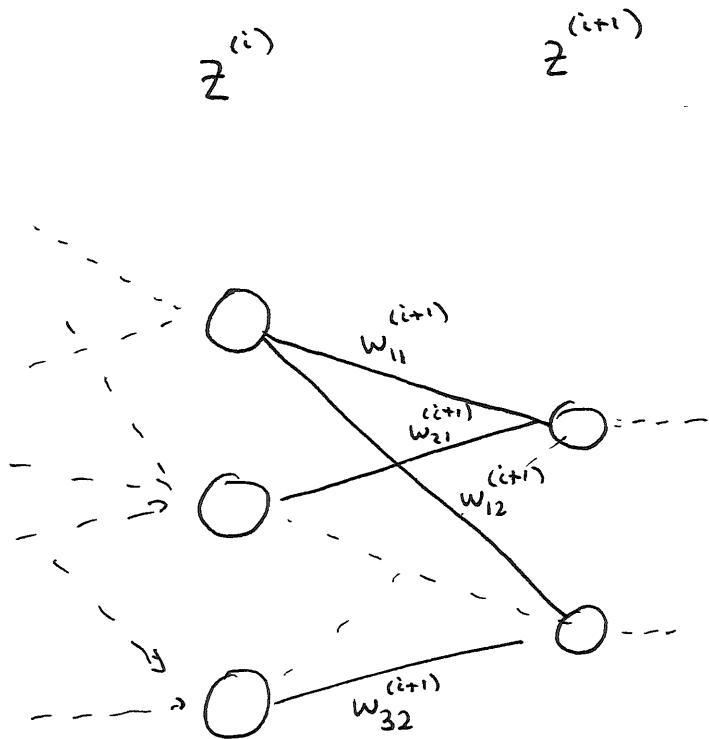
$$z = \sum x_i w_i \quad (\text{input})$$

$$a = \sigma(z) \quad (\text{activation})$$

$$\sigma(t) = \frac{1}{1+e^{-t}} \quad \text{OR} \quad \sigma(t) = \begin{cases} t & \text{if } t \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{OR} \quad \sigma(t) = \tanh(t)$$

$$\text{OR} \quad \sigma(t) = t$$

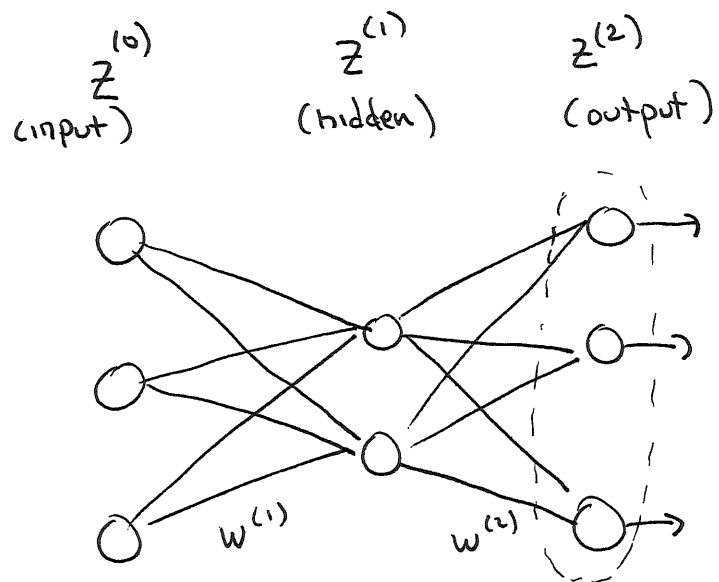
Layers



$$z_l^{(i+1)} = \sum_{j=1}^{n^{(i)}} a_j^{(i)} w_{jl}^{(i+1)}$$

$$a_l^{(i+1)} = \sigma(z_l^{(i+1)})$$

With all weights specified, a network defines a function  
 called the "feed-forward" function



$z^{(0)}$  1x3 matrix

$w^{(1)}$  3x2 matrix

$w^{(2)}$  2x3 matrix

$$z^{(0)} \rightarrow z^{(0)} w^{(1)} \rightarrow \sigma(z^{(0)} w^{(1)}) \rightarrow \sigma(z^{(0)} w^{(1)}) w^{(2)}$$

(elementwise)

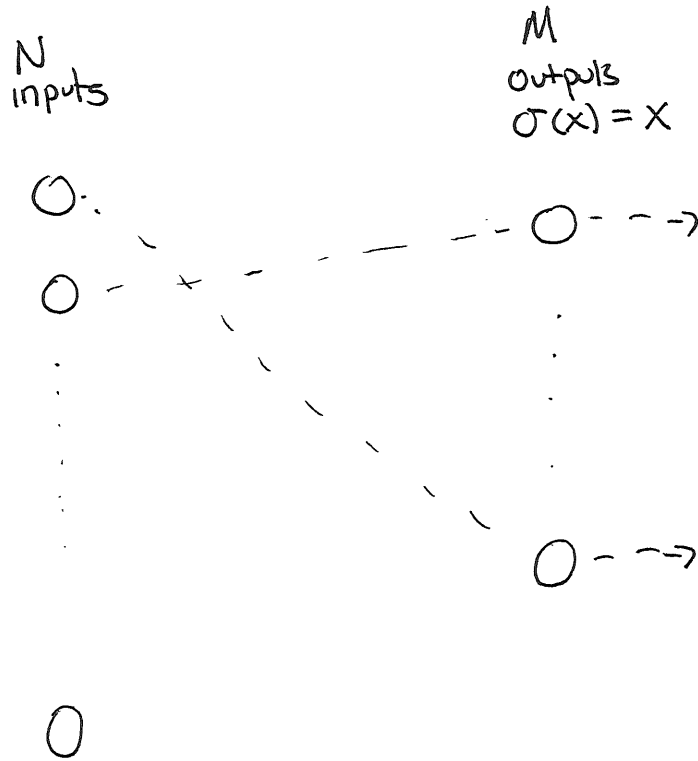
↓

$$\mathbb{R} F(\sigma(z^{(0)} w^{(1)}) w^{(2)})$$

↗  
might not be elementwise

# Linear Regression

Input vectors  $X$  in  $N$  dimensions  
Response vectors  $Y$  in  $M$  dimensions



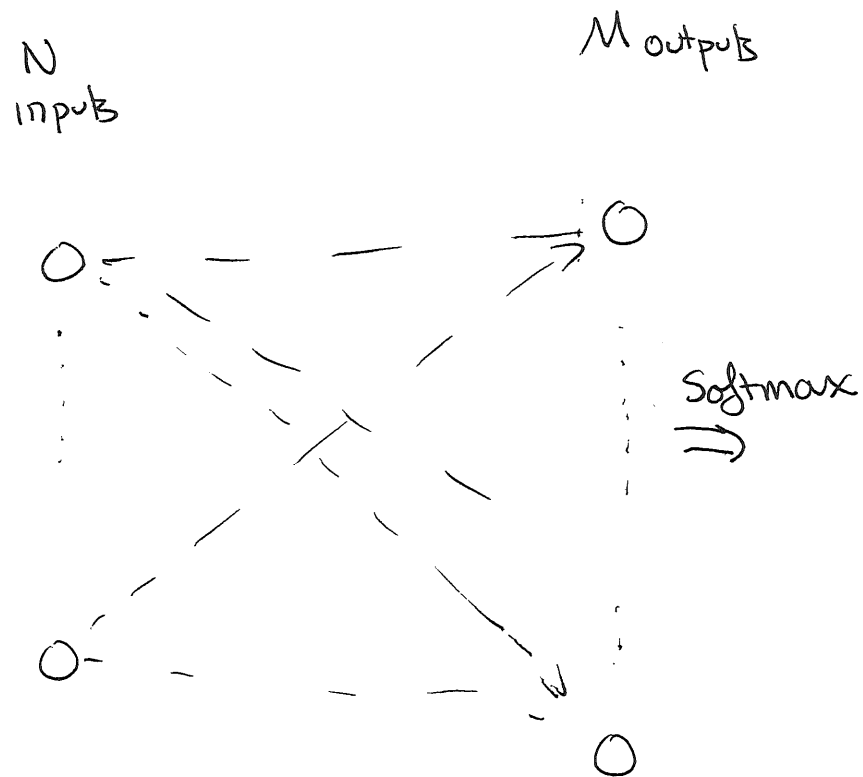
Weights

$$W_{ij}^{(1)}$$

$$1 \leq i \leq N$$
$$1 \leq j \leq M$$

$$\text{output} = X W^{(1)}$$

# Logistic Regression



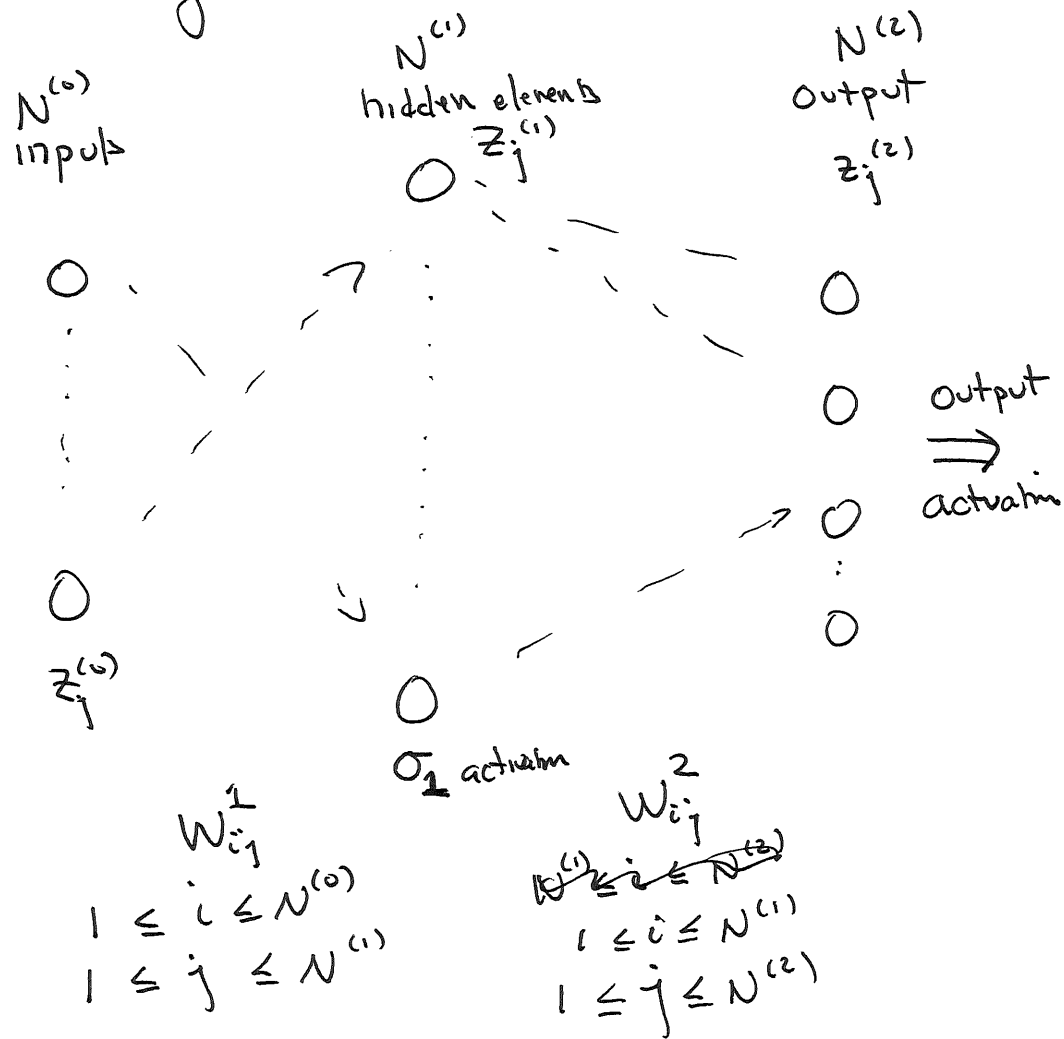
Weights  $W_{ij}^{(1)}$

$$1 \leq i \leq N$$

$$1 \leq j \leq M$$

$$\text{output} = \text{Softmax}(XW^{(1)})$$

# A "hidden" Layer



Training: Given a network ~~with~~  $F_W$  depending on the weights

Ingredients: data  $(x^{[i]}, y^{[i]})$   $i = 1, \dots, M$

A loss function

$$L_W = \frac{1}{M} \sum_{i=1}^M \mathcal{L}(y^{[i]}, F_W(x^{[i]}))$$

that compares the output  $F_W(x^{[i]})$  to  $y^{[i]}$

NOTE:  $L_W$  is a function of the WEIGHTS  
The "data"  $(x^{[i]}, y^{[i]})$  is fixed.

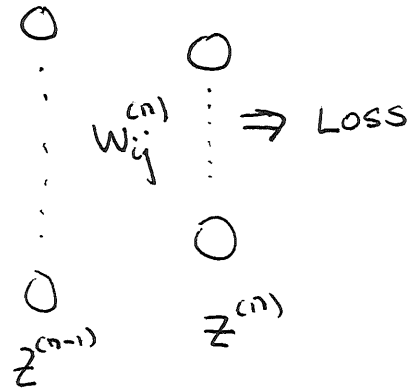
Goal is to minimize  $L$  by varying  $W$ .

Strategy: We have  $W_{ij}^{(l)}$   $l = 1, \dots, L$  for  $L$  layers

Compute  $\frac{\partial L}{\partial W_{ij}^{(l)}}$  and use gradient descent

$$W_{ij}^{(l)} = W_{ij}^{(l)} - \lambda \frac{\partial L}{\partial W_{ij}^{(l)}}$$

# Backpropagation



Step 1:

Compute  $\frac{\partial L}{\partial z_j^{(n)}}$

Examples:  $L = \sum (y_j - z_j)^2$

MSE  $L(y^{[i]}, z_j^{(n)}(x^{[i]})) = \frac{1}{2} \|y^{[i]} - z_j^{(n)}\|^2 = \frac{1}{2} \sum_{j=1}^{N_n} (y_j^{[i]} - z_j^{(n)})^2$

"cross entropy"  $L(y^{[i]}, z_j^{(n)}(x^{[i]})) = \sum_{j=1}^{N_n} y_j \log\left(\frac{e^{z_j}}{H}\right) = \sum_1^{N_n} y_j z_j - \log H$

$\frac{\partial L}{\partial z_j} = y_j - p_j$  where  $p_j = e^{z_j} / H$ .





# Earlier Layers

$$\delta_j^{(i-1)} = \sum_l \frac{\partial L}{\partial z_l^{(i)}} \frac{\partial z_l^{(i)}}{\partial z_j^{(i-1)}}$$

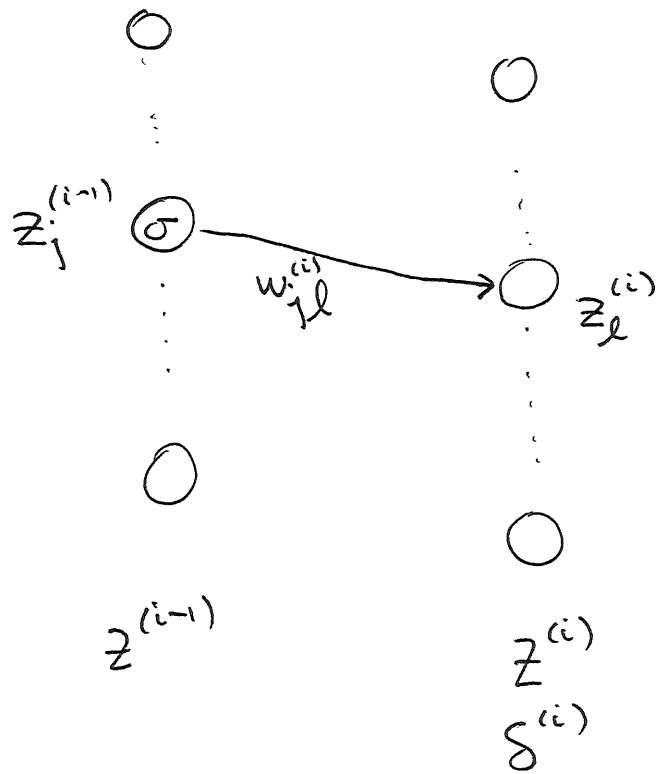
$$= \sum_l \delta_l^{(i)} \sigma'(z_j^{(i-1)}) w_{jl}^{(i)}$$

$$= \sigma'(z_j^{(i-1)}) \sum_l w_{jl}^{(i)} \delta_l^{(i)}$$

$l^{\text{th}}$  entry of  $W^{(i)} \delta^{(i)}$

So  $\delta_j^{(i-1)} = \sigma'(z_j^{(i-1)}) \sum_l w_{jl}^{(i)} \delta_l^{(i)}$

where the multiplication by  $\sigma'(z_j^{(i-1)})$  means entry by entry.



$$z_l^{(i)} = \sum_j \sigma(z_j^{(i-1)}) w_{jl}^{(i)}$$

## Backpropagation cont'd

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^{(l)}} = \sum \frac{\partial \mathcal{L}}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} = \sum \delta_j^{(l)} \sigma(z_i^{(l-1)})$$

$$\text{since } z_j^{(l)} = \sum_i \sigma(z_i^{(l-1)}) w_{ij}^{(l)}$$

In matrix terms

$$\frac{\partial \mathcal{L}}{\partial w^{(l)}}$$

is

the matrix

with entries

$$\left( \delta_j^{(l)} \sigma(z_i^{(l-1)}) \right)$$

This is the "outer product" of the vectors  $\delta^{(l)}$  and  $\sigma(z^{(l-1)})$

where  $\sigma(z^{(l-1)})$  means apply  $\sigma$  element wise.

## Backpropagation cont'd

To exploit this, during the forward pass, save the  $z^{(i)}$  and also compute  $\sigma'(z^{(i)})$

Then make a backward pass to compute the  $\delta^{(i)}$  using the weights from the forward pass. ~~to~~

Since the total  $\mathcal{L}$  is the sum of  $\mathcal{L}(y^{[i]}, F_w(x^{[i]}))$  you can accumulate the  $\delta^{(i)}$  on each pass, the  $\frac{\partial \mathcal{L}}{\partial w^{(i)}}$  one each pass

# Training Algorithm

Initialize network with random weights. Set all the  $\delta_x^{(i)}$   $\nabla W_*$  to zero

For  $x, y$  in the data:

- make a forward pass through the network, computing and saving the inputs  $z^i$  ~~and~~ at each stage

- make a backwards pass through the network

Compute  $\delta^{(i)}$  and  $\nabla W^{(i)}$  at each stage using the back propagation equations.

Accumulate  $\delta_x^{(i)} = \delta_x^{(i)} + \delta^{(i)}$  and  $\nabla W_*^{(i)} = \nabla W_*^{(i)} + \nabla W^{(i)}$

- periodically (maybe every time, maybe after  $B$  data points, maybe only at the end of the data)

Update the weights  $W^{(i)} = W^{(i)} - \lambda \nabla W^{(i)}$

Reset the accumulators to zero

This is ONE TRAINING EPOCH

Repeat until the loss stops decreasing...