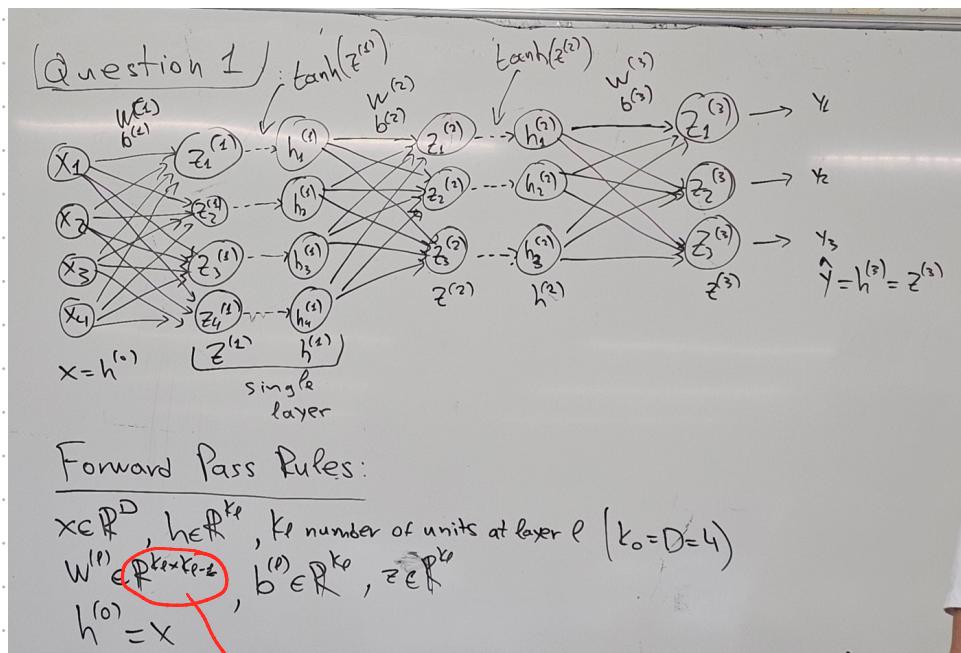


## QUESTION 1

4433 (including input & output)

All except output include hyperbolic tangent activation functions

- Initialize all weights & bias to 0.1
- Do SGD using squared error loss
- learning rate  $\eta = 0.1$



$$z \text{ preactivation: } z^{(l)} = W^{(l)} h^{(l-1)}(x) + b^{(l)}$$

$$h \text{ representation: } h^{(l)} = g(z^{(l)}), \quad g: \text{activation function}$$

$$\text{output: } \hat{y} = h^{(l=L)}, L \text{ is the last layer}$$

$$\text{loss: } L(\hat{y}, y)$$

weight init:

$$W^{(1)}_{4 \times 4} = \begin{bmatrix} w_{11} & \dots & w_{14} \\ \vdots & \ddots & \vdots \\ w_{41} & \dots & w_{44} \end{bmatrix} = \begin{bmatrix} 0.1 & \dots & 0.1 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.1 \end{bmatrix}$$

$$b^{(1)}_{4 \times 1} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}, \quad \text{And also } g = \tanh$$

$$W^{(2)}_{3 \times 4} = \begin{bmatrix} 0.1 & \dots \\ \dots & \dots \\ \dots & \dots \end{bmatrix}, \quad b^{(2)}_{3 \times 1} = \begin{bmatrix} 0.1 \\ 0.1 \\ 0.1 \end{bmatrix}$$

$$W^{(3)}_{3 \times 3} = \begin{bmatrix} \dots \\ \dots \\ \dots \end{bmatrix}, \quad b^{(3)}_{3 \times 1} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}$$

$$x = [1 \ 0 \ 1 \ 0]^T, y = [0 \ 1 \ 0]^T$$

Run the forward pass:

$$h^{(0)} = x = [1 \ 0 \ 1 \ 0]^T$$

$$z^{(1)} = W^{(1)} h^{(0)}(x) + b^{(1)} = [0.3 \ 0.3 \ 0.3 \ 0.3]^T$$

$$h^{(1)} = \tanh(z^{(1)}) = [0.2913 \ \dots \ 0.2913]^T$$

$$z^{(2)} = W^{(2)} h^{(1)} + b^{(2)} = [0.2165 \ \dots \ 0.2165]^T$$

$$h^{(2)} = \tanh(z^{(2)}) = [\dots]^T$$

$$z^{(3)} = W^{(3)} h^{(2)} + b^{(3)} = [0.16936 \ \dots \ 0.16936]^T$$

$$h^{(3)} = z^{(3)} = y$$

$$\text{Loss: } L(\hat{y}, y) = \frac{1}{2} \sum_{i=1}^3 (\hat{y}_i - y_i)^2 = \dots = 0.376$$

sum of squares

in mean squared error  
we would normalize  
 $\frac{1}{N}$  but here no!

Backpropagation (Recap):  $\frac{\partial L}{\partial W^{(l)}} , \frac{\partial L}{\partial b^{(l)}} , \frac{\partial L}{\partial h^{(l)}} , \frac{\partial L}{\partial z^{(l)}}$

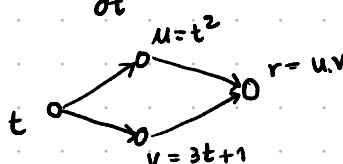
(slide 53)

Recap of the chain rule:

Example: If a function  $r(t)$  can be written in terms of intermediate results  $g_i(t)$  then we have:

$$\frac{\partial r(t)}{\partial t} = \sum_i \frac{\partial r}{\partial g_i} \frac{\partial g_i(t)}{\partial t}$$

so for example:

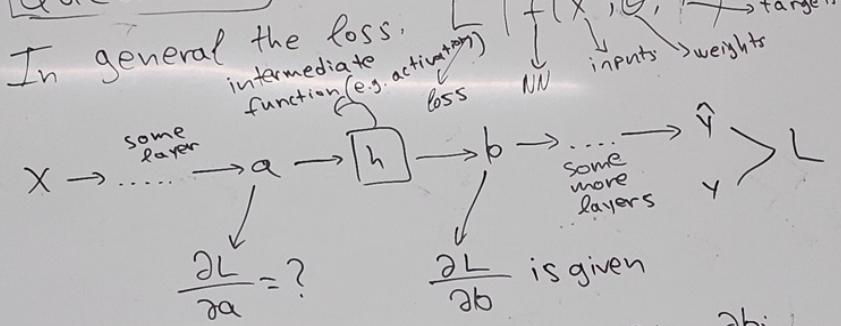


$$\frac{\partial r(t)}{\partial t} = \frac{\partial r(u)}{\partial u} \frac{\partial u(t)}{\partial t} + \frac{\partial r(v)}{\partial v} \frac{\partial v(t)}{\partial t} = \dots = 9t^2 + 2t$$

In general the loss:  $L(f(x; \theta), y)$

↓ ↓ ↓ targets  
loss NN inputs weights

### Question 1)



Use of the chain-rule:  $\frac{\partial L}{\partial a_i} = \sum_j \frac{\partial L}{\partial b_j} \frac{\partial b_j}{\partial a_i}$

Back-Prop:  $\frac{\partial L}{\partial \hat{y}}, \frac{\partial L}{\partial h^{(k)}}, \frac{\partial L}{\partial z^{(l)}}, \frac{\partial L}{\partial w^{(l)}}, \frac{\partial L}{\partial b^{(l)}}$

$$\cancel{\text{slide 54}} \cdot \frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} \frac{1}{2} (\hat{y} - y)^2 = \hat{y} - y \left( = \frac{\partial L}{\partial h^{(k)}} \right)$$

$$\frac{\partial L}{\partial z^{(l)}} = ? : h^{(l)} = g(z^{(l)}) \xrightarrow{\text{chain rule}} \frac{\partial L}{\partial z^{(l)}} = \frac{\partial L}{\partial h^{(l)}} \circ \frac{\partial h^{(l)}}{\partial z^{(l)}} \downarrow g'$$

For  $l=3$ :  $g^l$  is the identity

$$\text{so } \frac{\partial L}{\partial z^{(3)}} = \frac{\partial L}{\partial h^{(3)}}$$

For  $l \neq 3$ :  $g = \tanh$ , so  $g^l = 1 - \tanh^2$

$$h^{(l)} = \tanh(z^{(l)})$$

$$\begin{aligned} \text{so } \frac{\partial L}{\partial z^{(l)}} &= \frac{\partial L}{\partial h^{(l)}} \circ (1 - \tanh^2(z^{(l)})) \\ &= \frac{\partial L}{\partial h^{(l)}} \circ (1 - h^{(l)2}) \end{aligned}$$

$$\bullet \frac{\partial L}{\partial b^{(l)}} = \frac{\partial L}{\partial z^{(l)}} \quad (\text{slide 61})$$

$$\bullet \frac{\partial L}{\partial w^{(l)}} = \frac{\partial L}{\partial z^{(l)}} h^{(l-1)\top} \quad (\text{slide 61})$$

$$\bullet \frac{\partial L}{\partial h^{(l-1)}} = w^{(l)\top} \frac{\partial L}{\partial z^{(l)}}$$

### Question 1

$$\frac{\partial L}{\partial \hat{y}} = \frac{\partial L}{\partial h^{(3)}} = \hat{y} - y = \frac{\partial L}{\partial z^{(2)}}$$

$$\frac{\partial L}{\partial w^{(2)}} = \frac{\partial L}{\partial z^{(2)}} h^{(2)\top} = (\hat{y} - y) h^{(2)\top} = \dots$$

$$\frac{\partial L}{\partial b^{(2)}} = \frac{\partial L}{\partial z^{(2)}} = \hat{y} - y$$

$$\frac{\partial L}{\partial h^{(2)}} = w^{(2)\top} \frac{\partial L}{\partial z^{(2)}} = w^{(2)\top} (\hat{y} - y)$$

$$\frac{\partial L}{\partial z^{(2)}} = \frac{\partial L}{\partial h^{(2)}} \odot \frac{\partial h^{(2)}}{\partial z^{(2)}} = w^{(2)\top} (\hat{y} - y) \odot (1 - h^{(2)2})$$

⚠ Don't forget to update the weights

$$w^{(l)} = w^{(l)} - n \frac{\partial L}{\partial w^{(l)}}$$

$$b^{(l)} = b^{(l)} - n \frac{\partial L}{\partial b^{(l)}}$$

the grad of the sum  
is the sum of the  
grads

$$\Delta L = \Delta L_1 + \Delta L_2$$

$$\Delta L = \Delta L_1 + \Delta L_2$$

2) We know how to compute the grad for 1 example  
what if we are given more?

$$L(\theta) = \sum_i L(f(x_i; \theta), y_i)$$

$$\rightarrow \nabla_{\theta} L(\theta) = \nabla_{\theta} \sum_i L(f(x_i; \theta), y_i) = \\ = \sum_i \nabla_{\theta} L(f(x_i; \theta), y_i)$$

$$L(\theta) = \frac{1}{N} \sum_i L_i$$

$$\nabla_{\theta} L(\theta) = \frac{1}{N} \sum_i \nabla_{\theta} L_i$$

### QUESTION 2

- softmax activation
- error f. is cross-entropy

Same as question 1 but:

$$\hat{y} = h^{(3)} = \text{softmax}(z^{(3)})$$

$$\text{cross-entropy: } L(\hat{y}, y) = - \sum_k y_k \log \hat{y}_k$$

Tip: you can skip the computation of  $\frac{\partial L}{\partial h^{(3)}}$

its easier to compute directly  $\frac{\partial L}{\partial z^{(3)}}$

$$\begin{aligned}
 L(z, y) &= -\sum_k y_k \log \hat{y}_k = \\
 &= -\sum_k y_k \xrightarrow{\text{loop}} \text{softmax}(z_k) = \\
 &= -\sum_k y_k \log \left[ \frac{\exp(z_k)}{\sum_j \exp(z_j)} \right] \\
 &= -\sum_k y_k \left[ z_k - \log \sum_j \exp(z_j) \right] = \\
 &= \sum_k y_k \left[ \log \left( \sum_j \exp(z_j) \right) - z_k \right]
 \end{aligned}$$

↗ because derivative  
 w.r.t respect to  $z_i$   
 will have zero in  
 all the other  
 variables

$$\begin{aligned}
 \frac{\partial L}{\partial z_i} &= \frac{\partial}{\partial z_i} \left[ \sum_k y_k \log \sum_j \exp(z_j) - \sum_k y_k z_k \right] \\
 &= \frac{\partial}{\partial z_i} \sum_k y_k \log \left( \sum_j \exp(z_j) \right) - y_i \\
 &\quad \left( \begin{array}{l} \dots y_{i+1} + y_{i+2} + \dots \\ \dots y_1 z_1 + \dots + y_n z_n \end{array} \right) \\
 &\quad \frac{\partial}{\partial z_i} (y_i z_i) = y_i \\
 &= \sum_k y_k \frac{\partial}{\partial z_i} \log \sum_j \exp(z_j) - y_i \\
 &= \frac{\partial}{\partial z_i} \log \left( \sum_j \exp(z_j) \right) - y_i \\
 &= \frac{1}{\sum_j \exp(z_j)} \cdot \frac{\partial}{\partial z_i} \sum_j \exp(z_j) - y_i = \\
 &= \frac{\exp(z_i)}{\sum_j \exp(z_j)} - y_i = \\
 &= \text{softmax}(z_i) - y_i \quad (\text{slide 48})
 \end{aligned}$$

$$\text{So } \frac{\partial L}{\partial z^{(3)}} = \text{softmax}(z^{(3)}) - y = \hat{y} - y$$

$$\nabla_{z^{(3)}} L = \text{softmax}(z^{(3)}) - y$$

And we continue by applying exactly the same rules as in question 1

2) As in Q1 the derivative of the sum is the sum of the derivatives ...

$$\begin{aligned}
 \frac{\partial L}{\partial w} &\xrightarrow{\text{matrix IR}^{k \times k-1}} \\
 w^{(l)} &= w^{(l)} - \eta \frac{\partial L}{\partial w}
 \end{aligned}$$

