**Exercise 3**

Consider a sigmoid neuron with 1D input $x$, weight $w$, bias $b$ and output $y = \sigma(wx + b)$. The target is the variable $z$. Consider the cost function $J(w, b) = \frac{1}{2}(y - z)^2$.

- Find $\nabla J(w, b)$ and show that $\|\nabla J\| < \frac{1}{4}\sqrt{1 + x^2}(1 + |z|)$.

- Write the gradient descent iteration for the sequence $(w_n, b_n)$.

· Gradient $\nabla J (w, b)$

define $u = wx + b$      $y = \sigma(u)$      $J = \frac{1}{2}(y - z)^2$

Using chain rule:

$$\begin{cases} \frac{\partial J}{\partial w} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial w} \\ \frac{\partial J}{\partial b} = \frac{\partial J}{\partial y} \cdot \frac{\partial y}{\partial u} \cdot \frac{\partial u}{\partial b} \end{cases}$$

$\hookrightarrow$ compute each derivative

$$\begin{cases} \frac{\partial J}{\partial y} = y - z \\ \frac{\partial y}{\partial u} = \sigma(u)(1 - \sigma(u)) = y(1 - u) \\ \frac{\partial u}{\partial w} = x, \quad \frac{\partial u}{\partial b} = 1 \end{cases}$$      (derivative of sigmoid)

$\longrightarrow$ chain rule

$$\begin{cases} \frac{\partial J}{\partial w} = (y - z) \cdot y(1 - y)x \\ \frac{\partial J}{\partial b} = (y - z) \cdot y(1 - y) \end{cases}$$

· Bound on $\|\nabla J\| = \sqrt{\left(\frac{\partial J}{\partial w}\right)^2 + \left(\frac{\partial J}{\partial b}\right)^2}$

with $0 \le y \le 1 \longrightarrow 0 \le y(1 - y) \le \frac{1}{4}$      (maximum when $y = 0.5 \longrightarrow y(1-y) = 0.25$)

denote $A = (y - z) \, y (1 - y)$

$\quad\quad \hookrightarrow \left(\frac{\partial J}{\partial w}\right)^2 + \left(\frac{\partial J}{\partial b}\right)^2 = A^2(x^2 + 1)$

using the bound $|A| = |y - z| \cdot y(1 - y) \le (|y| + |z|) \cdot \frac{1}{4} \le \frac{1 + |z|}{4}$

hence $\|\nabla J\| \le \sqrt{(x^2 + 1)} \cdot \left(\frac{1 + |z|}{4}\right) = \frac{1}{4}\sqrt{1 + x^2}(1 + |z|)$

- Gradient descent iteration

Take $\eta > 0$

$$w^{i+1} = w^i - \eta \frac{\partial J}{\partial w} = w^i - \eta (y-z) y (1-y) x$$

$$b^{i+1} = b^i - \eta \frac{\partial J}{\partial w} = b^i - \eta (y-z) y (1-y)$$

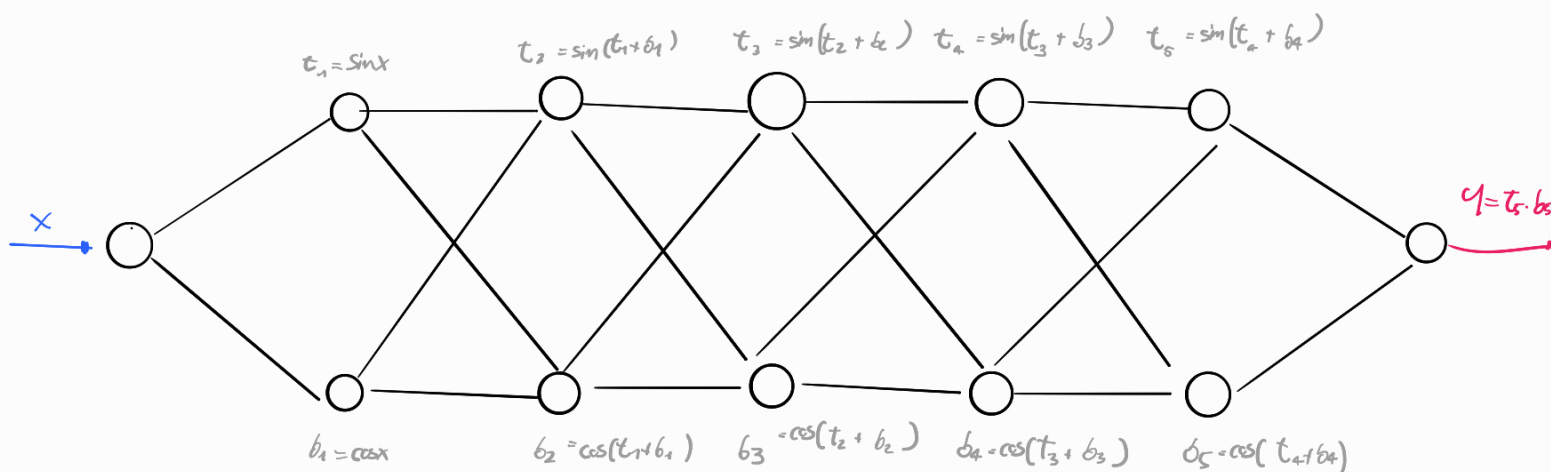$$\left( \text{with } y = \sigma(w^i x + b^i) \right)$$

**Exercise 3**

Consider the following computational graph:



Figure 1: Computational Graph

The upper node in each layer computes $\sin(x+y)$ and the lower computes $\cos(x+y)$ with respect to its 2 inputs. For the first hidden layer, there is only a single input $x$, and therefore the values $\sin(x)$ and $\cos(x)$ are computed. The final node computes the product of the two inputs. The single input is denoted by $x$ (value in radiants). Compute the numerical value of the partial derivative of the output with respect to $x$ for $x = 1$ using the backpropagation algorithm. Explain clearly each step you have performed.



$$t_1 = \sin x$$
$$t_2 = \sin(t_1 + b_1)$$
$$t_3 = \sin(t_2 + b_2)$$
$$t_4 = \sin(t_3 + b_3)$$
$$t_5 = \sin(t_4 + b_4)$$

$$y = t_5 \cdot b_5$$

$$b_1 = \cos x$$
$$b_2 = \cos(t_1 + b_1)$$
$$b_3 = \cos(t_2 + b_2)$$
$$b_4 = \cos(t_3 + b_3)$$
$$b_5 = \cos(t_4 + b_4)$$

$x = 1$

**Forward Pass**

$h_1 \begin{cases} t_1 = \sin(1) = 0,84 \\ b_1 = \cos(1) = 0.54 \end{cases}$

$h_2 \begin{cases} t_1 + b_1 = 1,38 \\ t_2 = \sin(1,38) = 0,98 \\ b_2 = \cos(1.38) = 0,18 \end{cases}$

$h_3 \begin{cases} t_2 + b_2 = 1.16 \\ t_3 = \sin(1.16) = 0,91 \\ b_3 = \cos(1.16) = 0,40 \end{cases}$

$h_4 \begin{cases} t_3 + b_3 = 1,31 \\ t_4 = \sin(1.31) = 0,92 \\ b_4 = \cos(1.31) = 0,25 \end{cases}$

$h_5 \begin{cases} t_4 + b_4 = 1,22 \\ t_5 = \sin(1.22) = 0,94 \\ b_5 = \cos(1.22) = 0,34 \end{cases}$

output $y = t_5 \cdot b_5 = 0,32$

**Backward Pass** $\left( \dfrac{d \sin x}{dx} = \cos x, \quad \dfrac{d \cos x}{dx} = -\sin x \right)$ $\qquad \dot{t} = \dfrac{\partial y}{\partial t}$ $\qquad$ obs: $T$ and $b$ have same contribute on derivative (symmetry!)

$h_5 \begin{cases} \dot{t_5} = b_5 = 0,34 \\ \dot{b_5} = t_5 = 0,94 \end{cases}$

$h_4 \begin{cases} \dot{t_4} = \dot{t_5} \cdot \cos(t_4 + b_4) + \dot{b_5}(-\sin(t_4 + b_4)) = \ldots = -0,72 \\ \dot{b_4} = -0,77 \end{cases}$

$h_3$ 
$$\begin{cases} \dot{t_3} = \dot{t_4}\cos(t_3+b_3) + \dot{b_4}\left(-\sin(t_3+b_3)\right) & = \dots \quad 0.54 \\ \dot{b_3} = 0.54 \end{cases}$$

$h_2$ 
$$\begin{cases} \dot{t_2} = \dot{t_3}\cos(t_2+b_2) \cdot \dot{b_3}\left(-\sin(t_2+b_2)\right) = \dots \quad -0.28 \\ \dot{b_2} = -0.28 \end{cases}$$

$h_1$ 
$$\begin{cases} \dot{t_1} = \dot{t_2}\cos(t_1+b_1) \cdot \dot{b_2}\left(-\sin(t_1+b_1)\right) = 0.22 \\ \dot{\delta_1} = 0.22 \end{cases}$$

input $\quad \dot{x} = \dot{t_1}\cos(1) + \dot{b_1}\left(-\sin(1)\right) = -0.02$

**Exercise 3**

Consider the two classes of patterns that are shown in the following figure where Class I represents vertical lines and Class II represents horizontal lines.



Class I

Class II

1. Are these categories linearly separable ?

2. Design a multilayer network to distinguish these categories.

---

- Encode input as $[x_1, x_2, x_3, x_4]$

$$\begin{matrix} x_1 & x_2 & x_3 & x_4 \end{matrix}$$
$$[1 \quad 0 \quad 1 \quad 0] \longrightarrow \text{class 1}$$
$$[0 \quad 1 \quad 0 \quad 1] \longrightarrow \text{class 1}$$
$$[1 \quad 1 \quad 0 \quad 0] \longrightarrow \text{class 2}$$
$$[0 \quad 0 \quad 1 \quad 1] \longrightarrow \text{class 2}$$

↳ Linear decision boundary in the form $W^T x + b > 0$    (e.g $> 0 \rightarrow$ class 1)

$$[1 \quad 0 \quad 1 \quad 0] \longrightarrow W_1 + W_3 \quad > 0 \qquad (1)$$
$$[0 \quad 1 \quad 0 \quad 1] \longrightarrow W_2 + W_4 \quad > 0 \qquad (2)$$
$$[1 \quad 1 \quad 0 \quad 0] \longrightarrow W_1 + W_2 \quad < 0 \qquad (3)$$
$$[0 \quad 0 \quad 1 \quad 1] \longrightarrow W_3 + W_4 \quad < 0 \qquad (4)$$

→ It raises inconsistencies :

From (1) and (3) $\longrightarrow \begin{cases} W_1 + W_3 > 0 \\ -W_1 > W_2 \end{cases} \xrightarrow{\text{sum}} W_3 > W_2$   ▽!

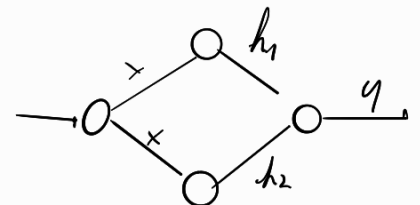From (2) and (4) $\longrightarrow \begin{cases} W_2 + W_4 > 0 \\ -W_4 > W_3 \end{cases} \xrightarrow{\text{sum}} W_2 > W_3$

- Multi-layer Perceptron

Vertical detector : $h_1 = \sigma(x_1 + x_3 - x_2 - x_4)$

horizontal detector: $h_2 = \sigma(x_1 + x_2 - x_3 - x_4)$

output : $y = \sigma(w_1 h_1 + w_2 h_2 + b)$

loss : Cross-entropy

## Exercise 3

Suppose that the output $\hat{y}_k$ of a given unit in a neural network is given by the softmax function *i.e.*:

$$\hat{y}_k = \frac{\exp(a_k)}{\sum_j \exp(a_j)}. \tag{2}$$

- Show that the output of the softmax function does not change if you shift, in all components, the activations $a_j$ by some constant $c$.

- Explain why the shift $c = -\max_j(a_j)$ can be useful.

- Invariance to shift

    suppose shift all activation $a_j$ by constant $c$    :    $a_j' = a_j + c$

    $$\hat{y}_k' = \frac{e^{a_k + c}}{\sum_j e^{a_j + c}} = \frac{e^{a_k} e^c}{\sum_j e^{a_j} e^c} = \frac{e^{a_k} e^c}{e^c \sum_j e^{a_j}} = \frac{e^{a_k}}{\sum_j e^{a_j}} = \hat{y}_k$$

- Choosing    $a_j' = a_j - \max_j(a_j)$

    cbs :    $a_j' \leq 0$    it helps numerical stability!

    - when some $a_j$ are large $e^{a_j}$ can overflow

    - if maximum is $0$ we ensure all components to be $\leq 1$
      and avoid overflow!

**Exercise 3**

Show that a multi-layer neural network with linear activation function $s(x) = x$ is equivalent to a single layer linear network. Assume that in each layer the inputs follow a Normal distribution with mean zero and small variance, *i.e.* $\sigma \ll 1$. For which of the activation functions $s(x) = 1/(1 + \exp(-x))$, $s(x) = \tanh(x)$, $s(x) = \text{relu}(x)$ and $s(x) = \text{selu}(x)$ is a deep network equivalent to a linear network for the given distribution ? The selu function is given by:

$$\text{selu}(x) = \begin{cases} \lambda x \text{ if } x > 0 \\ \alpha(\exp(x) - 1) \text{ otherwise} \end{cases}, \tag{4}$$

where $\lambda \approx 1.0507$ and $\alpha \approx 1.75814$. (*Hint:* consider the case $\sigma \to 0$ using a Taylor expansion around 0.)

---

• Deep Linear = single linear

MLP (with $L$ layers) has about: $y = W^L \cdot W^{L-1} \cdot \ldots W^2 \cdot W^1 \cdot x$

We can just truncate $W^L \cdot W^{L-1} \ldots W^2 \cdot W^1 = W_{eff}$ ⟶ $y = W_{eff} x$

(single neuron!)

• for which activation functions is a deep approximately linear when $x \sim N(0, \sigma^2)$ with $\sigma^2 \ll 1$ ?

1) Sigmoid $s(x) = \dfrac{1}{1+2^{-x}}$

Taylor expansion in $\varphi$: $s(x) \approx \dfrac{1}{2} + \dfrac{x}{4} - \dfrac{x^3}{48} + \ldots$

For small $\sigma$ $s(x) \approx \dfrac{1}{2} + \dfrac{x}{4}$     Not linear: constant bias breaks linearity over multiple layers ✗

2) $\tanh(x)$

Taylor expansion $s(x) = x - \dfrac{x^3}{3} + \dfrac{2x^5}{15} \ldots$

For small $\sigma$ $s(x) = x$      Linear ✓

3) $Relu(x) = \max(0, x)$

$x < 0$ ⟶ output 0

$x > 0$ ⟶ output $x$

  ↳ normal distribution around $\emptyset$ means that whole input are mapped into $\emptyset$

Non linearity!

a) $\text{Selu}(x) = \begin{cases} \lambda x & x > 0 \\ \lambda \alpha (e^x - 1) & x \leq 0 \end{cases}$
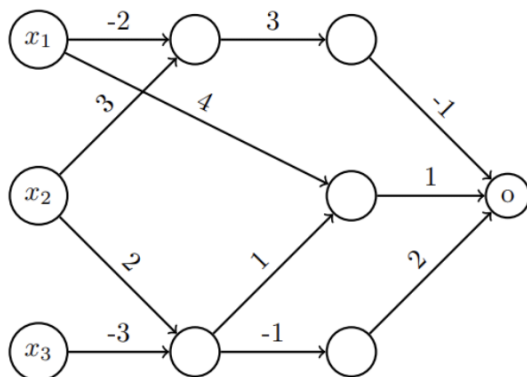
taylor expansion $\begin{cases} \lambda x & x > 0 \\ \lambda \alpha (x + \frac{x^2}{2} + \cdots) \approx \lambda \alpha x & x \leq 0 \end{cases}$

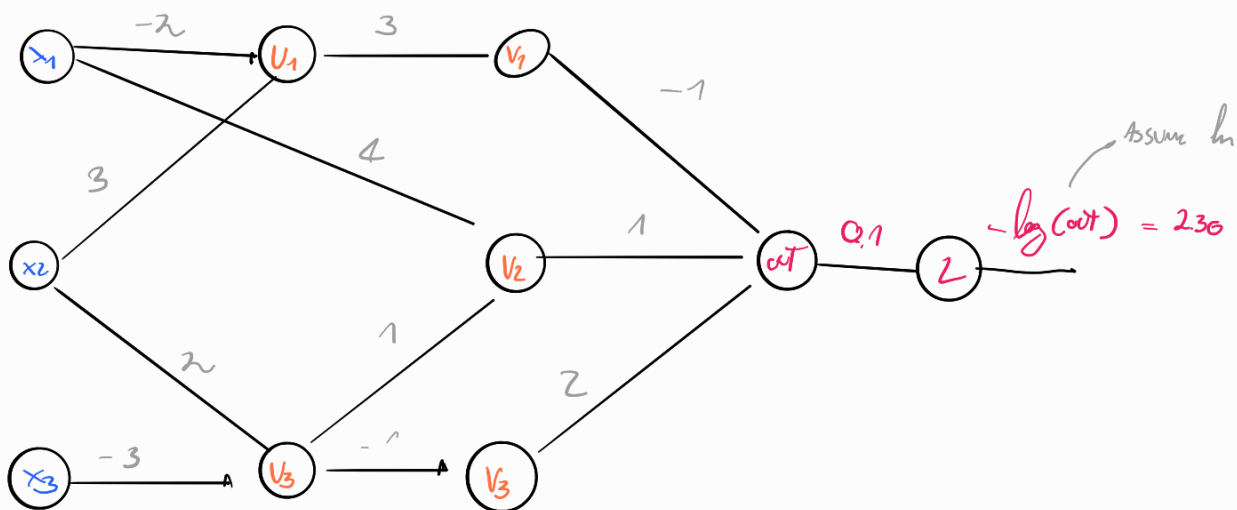$\llcorner$ this is piecewise linear (but not globally) slope differs between $x > 0$ and $x < 0$

## Exercise 3

Consider the following network where on each edge $(i, j)$ the value of $\dfrac{\partial y(j)}{\partial y(i)}$ is given; $y(k)$ denotes the activation of node $k$.



The output $o$ is equal to 0.1 and the loss function is $L = -log(o)$. Compute the value of $\dfrac{\partial L}{\partial x_i}$ for each input $x_i$ using the backpropagation method.



Assume $ln$

$-log(out) = 2.30$

$$\frac{\partial L}{\partial at} = -\frac{1}{out} = -10$$

$$\frac{\partial L}{\partial v_1} = \frac{\partial L}{\partial out} \cdot \frac{\partial at}{\partial v_1} = -10 \cdot -1 = 10$$

$$\frac{\partial L}{\partial v_2} = \frac{\partial L}{\partial out} \cdot \frac{\partial out}{\partial v_2} = -10 \cdot 1 = -10$$

$$\frac{\partial L}{\partial v_3} = \frac{\partial L}{\partial out} \cdot \frac{\partial out}{\partial v_3} = -10 \cdot 2 = -20$$

$$\frac{\partial L}{\partial u_1} = \frac{\partial L}{\partial v_1} \cdot \frac{\partial v_1}{\partial u_1} = 10 \cdot 3 = 30$$

$$\frac{\partial L}{\partial u_3} = \frac{\partial L}{\partial v_2} \cdot \frac{\partial v_2}{\partial u_3} + \frac{\partial L}{\partial v_3} \cdot \frac{\partial v_3}{\partial u_3} = (-10 \cdot 1) + (-20 \cdot -1) = 10$$

$$\frac{\partial L}{\partial x_1} = \frac{\partial L}{\partial U_1} \cdot \frac{\partial U_1}{\partial x_1} + \frac{\partial L}{\partial V_2} \cdot \frac{\partial V_1}{\partial U_1} = (30 \cdot -2) + (-10 \cdot 4) = -60 - 40 = -100$$

$$\frac{\partial L}{\partial x_2} = \frac{\partial L}{\partial U_1} \cdot \frac{\partial U_1}{\partial x_2} + \frac{\partial L}{\partial V_3} \cdot \frac{\partial U_3}{\partial x_2} = (30 \cdot 3) + (10 \cdot 2) = 90 + 20 = 110$$

$$\frac{\partial L}{\partial x_3} = \frac{\partial L}{\partial V_3} \cdot \frac{\partial U_3}{\partial x_3} = 10 \cdot -3 = -30$$
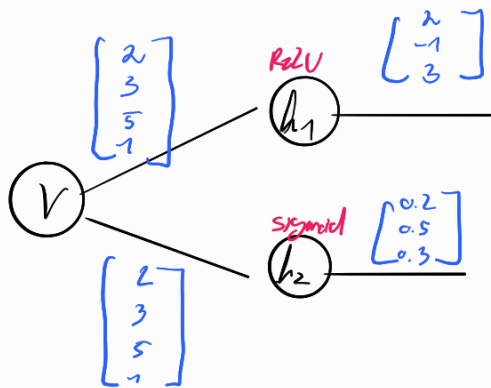
### Exercise 3

Consider a neural network in which a vectored node $v$ feeds into two distinct vectored nodes $h_1$ and $h_2$ computing different functions. The functions computed at the nodes are $h_1 = ReLU(W_1 v)$ and $h_2 = sigmoid(W_2 v)$. We do not know anything about the values of the variables in other parts of the network, but we know that $h_1 = [2, -1, 3]^T$ and $h_2 = [0.2, 0.5, 0.3]^T$, that are connected to the node $v = [2, 3, 5, 1]^T$. Furthermore the loss gradients are $\dfrac{\partial L}{\partial h_1} = [-2, 1, 4]^T$ and $\dfrac{\partial L}{\partial h_2} = [1, 3, 2]^T$, respectively. Show that the backpropagated loss gradient $\dfrac{\partial L}{\partial v}$ can be computed in terms of $W_1$ and $W_2$ as follows:

$$\frac{\partial L}{\partial v} = W_1^T \begin{bmatrix} -2 \\ 0 \\ 4 \end{bmatrix} + W_2^T \begin{bmatrix} 0.16 \\ 0.75 \\ -0.42 \end{bmatrix} \tag{5}$$

What are the sizes of $W_1, W_2$ and $\dfrac{\partial L}{\partial v}$ ?

Remember that $ReLU(x) = \max(0, x)$ and $sigmoid(x) = \dfrac{\exp(x)}{(\exp(x) + 1)}$.

$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix} \quad \overset{\text{ReLU}}{\underset{h_1}{\longrightarrow}} \quad \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} \qquad \frac{\partial L}{\partial h_1} = \begin{bmatrix} -2 \\ 1 \\ 4 \end{bmatrix}$$

$$V \qquad \overset{\text{sigmoid}}{\underset{h_2}{\longrightarrow}} \begin{bmatrix} 0.2 \\ 0.5 \\ 0.3 \end{bmatrix} \qquad \frac{\partial L}{\partial h_2} = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix}$$

$$\begin{bmatrix} 2 \\ 3 \\ 5 \\ 1 \end{bmatrix}$$

• With backpropagation rule

$$\frac{\partial L}{\partial V} = \frac{\partial L}{\partial h_1} \cdot \frac{\partial h_1}{\partial V} + \frac{\partial L}{\partial h_2} \cdot \frac{\partial h_2}{\partial V}$$

$$= W_1^T \left( \frac{\partial L}{\partial h_1} \odot \frac{\partial h_1}{\partial z_1} \right) + W_2^T \left( \frac{\partial L}{\partial h_2} \odot \frac{\partial h_2}{\partial z_2} \right)$$

with $z_1 = W_1 V$ , $h_1 = ReLU(z_1) \implies \frac{\partial h_1}{\partial z_1} = ReLU'(z_1)$

$z_2 = W_2 V$ , $h_2 = sigmoid(z_2) \overset{\sigma(z_2)}{\longrightarrow} \frac{\partial h_2}{\partial z_2} = \sigma(z_2)(1 - \sigma(z_2))$

$\implies \frac{\partial h_1}{\partial z_1}$

$$ReLU' = \begin{cases} 1 & z_i > 0 \\ \emptyset & z_i \leq 0 \end{cases}$$

from $h_1 = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} \longrightarrow z_1 = W_1 V = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix} \longrightarrow ReLU'(z_1) = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

$\longrightarrow \frac{\partial L}{\partial h_1} \odot ReLU'(z_1) = \begin{bmatrix} -2 \\ 1 \\ 4 \end{bmatrix} \odot \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -2 \\ 0 \\ 4 \end{bmatrix}$

$$\Rightarrow \frac{\partial h_2}{\partial z_2} \qquad \sigma'(z) = \sigma(z)(1 - \sigma(z))$$

$$\text{from } h_2 = \begin{bmatrix} 0.2 \\ 0.5 \\ 0.3 \end{bmatrix} \rightsquigarrow z_2 = W_2 V \begin{bmatrix} 0.2 \\ 0.5 \\ 0.3 \end{bmatrix} \longrightarrow \sigma'(z_2) = \begin{bmatrix} 0.2 \cdot 0.8 \\ 0.5 \cdot 0.5 \\ 0.3 \cdot 0.7 \end{bmatrix} = \begin{bmatrix} 0.16 \\ 0.25 \\ 0.21 \end{bmatrix}$$
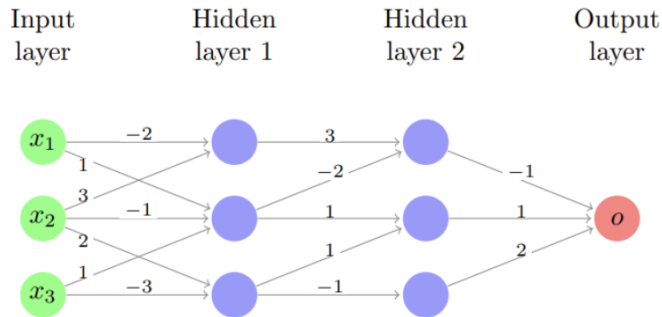
$$\Rightarrow \frac{\partial L}{\partial h_2} \odot \sigma'(z_2) = \begin{bmatrix} 1 \\ 3 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 0.16 \\ 0.25 \\ 0.21 \end{bmatrix} = \begin{bmatrix} 0.16 \\ 0.25 \\ 0.42 \end{bmatrix}$$

$$\text{Sizes:} \qquad V \in \mathbb{R}^a \qquad h_1, h_2 \in \mathbb{R}^3 \qquad W_1, W_2 \in \mathbb{R}^{3 \times a} \qquad \frac{\partial L}{\partial V} \in \mathbb{R}^a$$
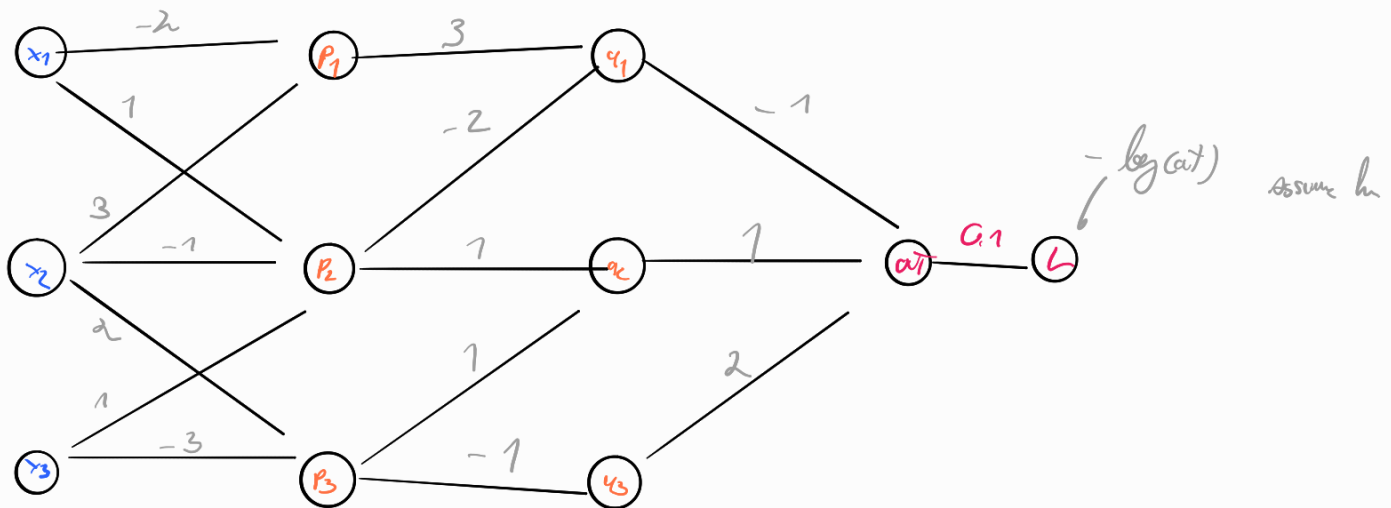
### Exercise 3

Consider the following network where on each edge $(i, j)$ the value of $\frac{\partial y(j)}{\partial y(i)}$ is given; $y(k)$ denotes the activation of node $k$.

| Input layer | Hidden layer 1 | Hidden layer 2 | Output layer |
|---|---|---|---|



The output $o$ is equal to 0.1 and the loss function is $L = -log(o)$. Compute the value of $\frac{\partial L}{\partial x_i}$ for each input $x_i$ using the backpropagation method.



$-log(o_t)$    Assume ln

**Backpropagation**

$$\frac{\partial L}{\partial o_t} = -\frac{1}{o_t} = \frac{1}{o_1} = -10$$

$$\frac{\partial L}{\partial q_1} = \frac{\partial L}{\partial o_t} \cdot \frac{\partial o_t}{\partial q_1} = -10 \cdot -1 = +10$$

$$\frac{\partial L}{\partial q_2} = \frac{\partial L}{\partial o_{ut}} \cdot \frac{\partial o_{ut}}{\partial q_2} = -10 \cdot 1 = -10$$

$$\frac{\partial L}{\partial q_3} = \frac{\partial L}{\partial o_{uT}} \cdot \frac{\partial o_{ut}}{\partial q_3} = -10 \cdot 2 = -20$$

$$\frac{\partial L}{\partial P_1} = \frac{\partial L}{\partial q_1} \cdot \frac{\partial q_1}{\partial P_1} = 10 \cdot 3 = 30$$

$$\frac{\partial L}{\partial P_2} = \frac{\partial L}{\partial q_1} \cdot \frac{\partial q_1}{\partial P_2} + \frac{\partial L}{\partial q_2} \cdot \frac{\partial q_2}{\partial P_2} = (10 \cdot -2) + (-10 \cdot 1) = -30$$