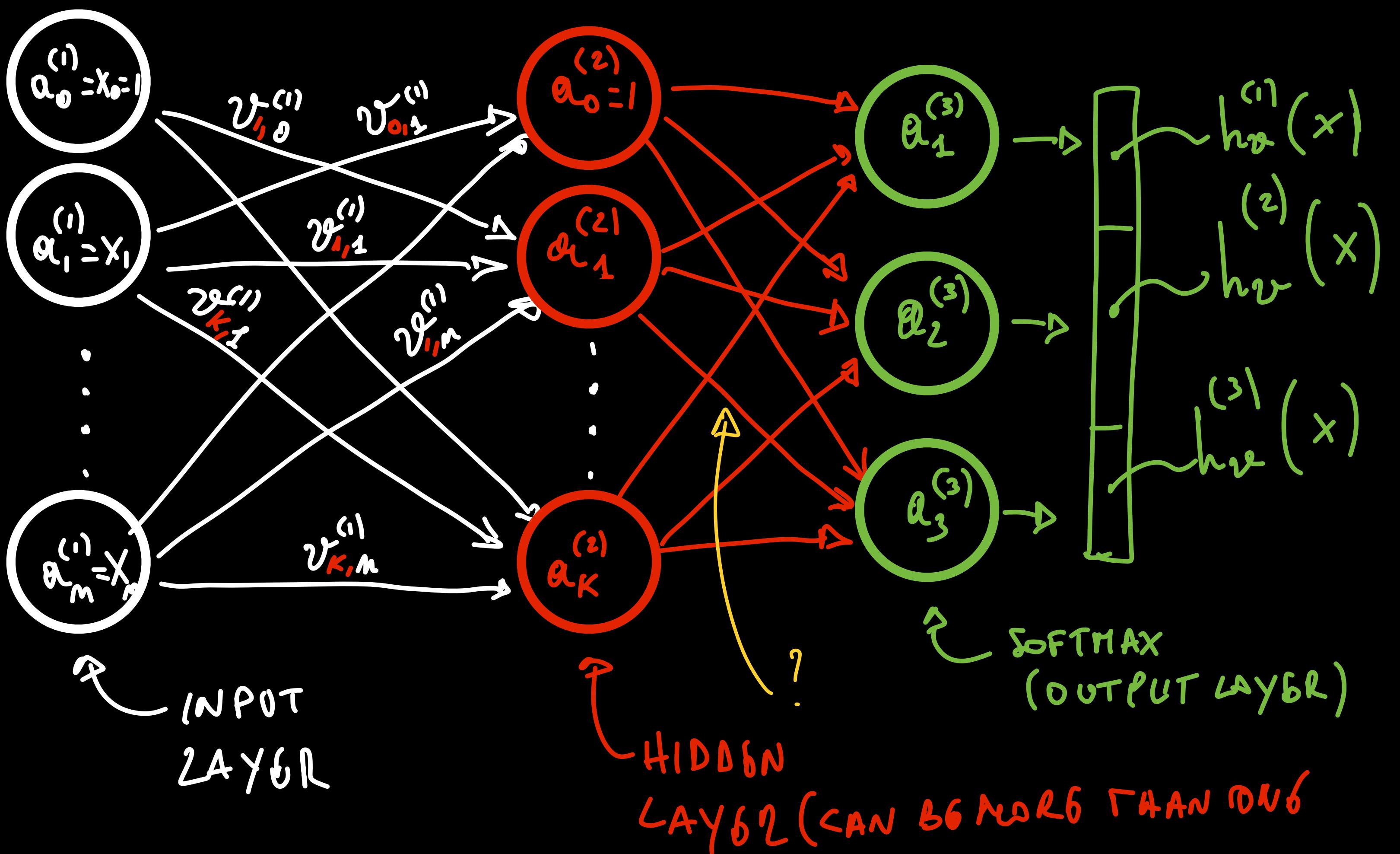


Backpropagation

(via Computational Graphs)

SEE ALSO CHAPTER 7 ("THE BACKPROPAGATION ALGORITHM")
OF THIS BOOK "NEURAL NETWORKS - A SYSTEMATIC INTRODUCTION"



AL CRESCEBISSE DELLA DIMENSIONE DEL INPUT CROSS
LA PROBABILITA' CAPOLO SPAZIO MA NON LINEARE (AD)
PERDERSI OR/AND PERDERSI NEURONI PER LAYER)

$$a_0^{(1)} = x_0 = 1 \quad a_1^{(1)} = x_1 \quad a_2^{(1)} = x_2 \quad \text{INPUT}$$

$$\mathcal{V}^{(1)} = \begin{bmatrix} v_{1,0}^{(1)} & \overset{3 \times 3}{v_{1,1}^{(1)}} & v_{1,2}^{(1)} \\ v_{2,0}^{(1)} & v_{2,1}^{(1)} & v_{2,2}^{(1)} \\ v_{3,0}^{(1)} & v_{3,1}^{(1)} & v_{3,2}^{(1)} \end{bmatrix}$$

$$\mathcal{V}^{(2)} = \begin{bmatrix} v_{1,0}^{(2)} & v_{1,1}^{(2)} & v_{1,2}^{(2)} & v_{1,3}^{(2)} \\ v_{2,0}^{(2)} & v_{2,1}^{(2)} & v_{2,2}^{(2)} & v_{2,3}^{(2)} \end{bmatrix} \quad 2 \times 4$$

WITH \mathcal{V} WE INDICATE ALL SET OF PARAMETERS

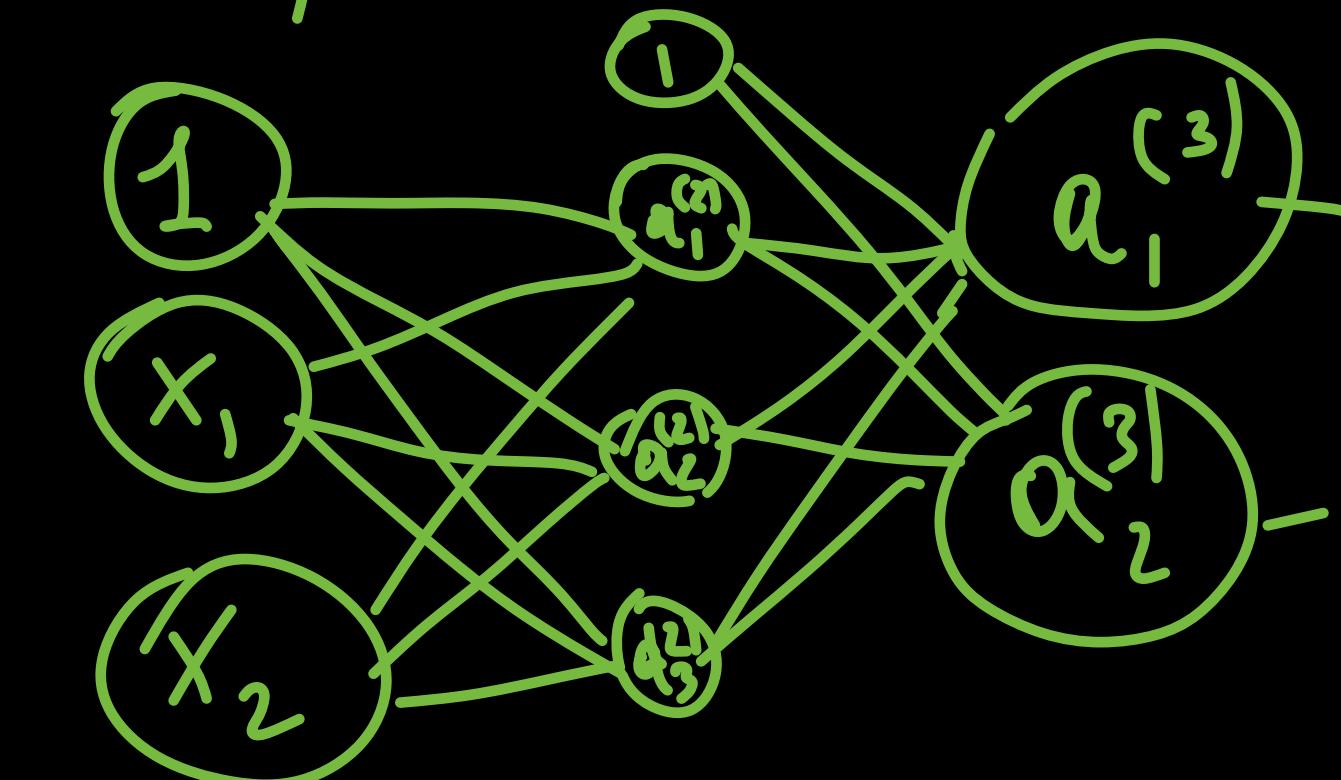
CONSIDERING THE MODEL ABOVE, THE NETWORK IS COMPOSED

BY 2 INPUT FEATURES EXTENDED WITH BIAS, JUST

ONE HIDDEN LAYER WITH 3 NEURONS + THE BIAS

NEURON OF THE HIDDEN LAYER AND AN OUTPUT LAYER

WITH 2 NEURONS



$$a_0^{(2)} = 1$$

$$x_i \downarrow$$

$$a_1^{(2)} = f\left(\sum_{i=0}^2 v_{1,i}^{(1)} a_i^{(1)}\right) = f\left(v_{1,0}^{(1)} x_0 + v_{1,1}^{(1)} x_1 + v_{1,2}^{(1)} x_2\right)$$

activation function (e.g., SIGMOID, RELU, ...)

QUANTI PARAMETRI HA
LA RETE NELLA SL1 D6
PROPSIBILITÀ?

$$a_2^{(2)} = f\left(\sum_{i=0}^2 v_{2,i}^{(1)} a_i^{(1)}\right)$$

$$a_3^{(2)} = f\left(\sum_{i=0}^2 v_{3,i}^{(1)} a_i^{(1)}\right)$$

$$a_1^{(3)} = f\left(\sum_{i=0}^3 v_{1,i}^{(2)} a_i^{(2)}\right)$$

$$a_2^{(3)} = f\left(\sum_{i=0}^3 v_{2,i}^{(2)} a_i^{(2)}\right)$$

$$h_2 = \begin{bmatrix} a_1^{(3)} & a_2^{(3)} \end{bmatrix}^T$$

SE LA RETE HA S_L UNITÀ NELLAYER L
E S_{L+1} UNITÀ NELLAYER L+1, ALLORA
 $v^{(L)}$ HA UNA DIMENSIONE DI $S_{L+1} \times S_L$
(NOTA CHE ① È CONNESSO CON I NEURONI)

UNA VOLTA FISSATI I PARAMETRI ϑ , LE FORMULE NELLA

SLI DEDICATE AGLI CODIMENTI POSSONO SERVIRE USANDO PNL
DISTRIBUZIONI LA MISTRIBURZIONE DI PROBABILITÀ SULLE
VARIABILI CLASSI DI UN NUOVO DATO DI INPUT x .

Oltre a operazioni PRATICHE IL NEURONI DI
FEEDFORWARD PROPAGATION (I DATI SI

PROPAGANO NEL GRADIVO CASO RAPPRESENTA UNA

RISATA DA SINISTRA VERSO DESTRA).

PER DETERMINARE I PARAMETRI BISOGNA USARE GRADIENTE
DESCENDENTE PER OTTENERE IL MINIMO DELLA FUNZIONE COSTO

$$J(\vartheta) = \frac{1}{2m} \sum_{i=1}^m L(\vartheta, x^{(i)}, y^{(i)}) + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{i=1}^{S_l} \sum_{j=1}^{S_{l+1}} (\vartheta_{j,i}^{(l)})^2$$

A LOSS FUNCTION (CLASSIFICATION/REGRESSION)
CROSS-ENTROPY MSE

$$h_\vartheta(x) = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_K^{(1)} \end{bmatrix} = \begin{bmatrix} P(y=1|x) \\ P(y=2|x) \\ \vdots \\ P(y=3|x) \end{bmatrix}$$

$L = \# \text{ LAYER DELLA RETE} = \text{ULTIMO LAYER}$

$$\sum_{i=1}^K a_i^{(L)} = 1$$

GRADIENT DESCENT NEEDS THE COMPUTATION
OF THE DERIVATIVES OF $J(\varphi)$

$$v_{ji}^{(e)} \leftarrow v_{ji}^{(e)} - \lambda \frac{\partial J(\varphi)}{\partial v_{ji}}$$

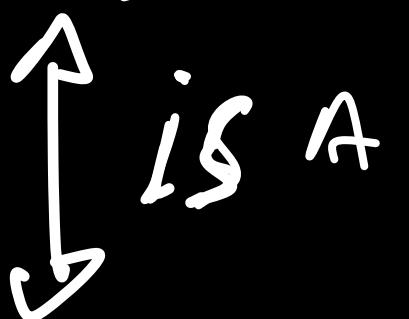
COMPUTES THE DERIVATIVES FOR EVERY PATH OF
THE NETWORK IS INEFFICIENT !!!

MANUAL
COMPUTATION

IS BORING !

WE NEED A SMART ALGORITHM \rightarrow BACKPROPAGATION

FEEDFORWARD NETWORK



COMPUTATIONAL GRAPH

NEURAL NETWORK REPRESENTS

A CHAIN OF FUNCTION COMPOSITIONS
WHICH TRANSFORM AN INPUT
TO AN OUTPUT

ACTIVATION FUNCTIONS

ARE ASSUMED SAME FOR EACH NODE

NODES ARE COMPUTING UNITS $(+ \dots * \dots)$

EDGES TRANSMIT (WEIGHT) $\circ \rightarrow 0$

NUMERICAL INFORMATION

FROM NODE TO NODE
(FORWARD)

EACH COMPUTING UNIT
IS ABLE TO COMPUTE

A SINGLE

PRIMITIVE
FUNCTION OF
ITS INPUT

$$\frac{1}{x^2}$$



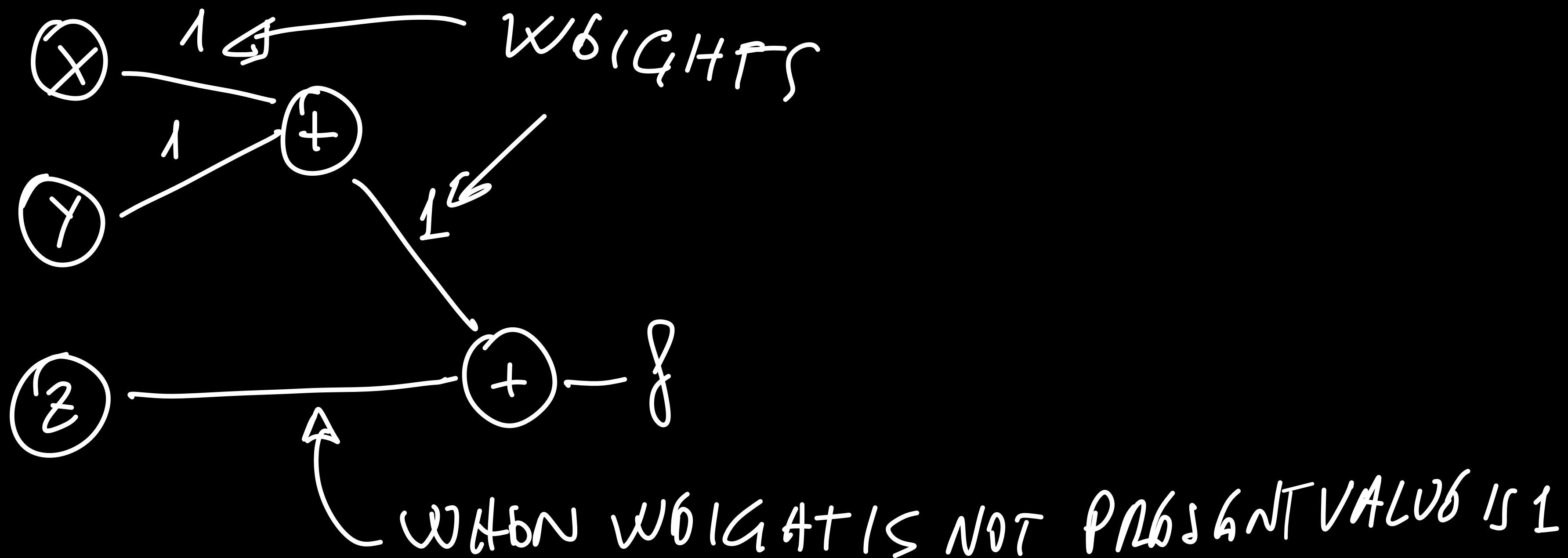
$$x^2 \rightarrow \frac{1}{x}$$

PRIMITIVE FUNCTIONS
HAVE TO BE CONTINUOUS
AND DIFFERENTIABLE

EXAMPLE OF COMPUTATIONAL GRAPH TO WHICH WE KNOW
ANALITICAL FORM IN ADVANCE:

$$f(x, y, z) = (x + y) + z$$

↑
↑
PRIMITIVE FUNCTION



LEARNING PROBLEM

- FIND THE OPTIMAL SET OF WEIGHTS SO THAT A NEURAL NETWORK h_θ APPROXIMATES A GIVEN FUNCTION F AS CLOSELY AS POSSIBLE WITH RESPECT TO A LOSS FUNCTION
- F IS NOT GIVEN EXPLICITLY, BUT ONLY IMPLICITLY THROUGH EXAMPLES (TRAINING SET)
- WEIGHTS ARE SMALL REAL NUMBERS INITIALIZED RANDOMLY AND CONSIDERING SIMPLICITY DRAWING

BACK PROPAGATION

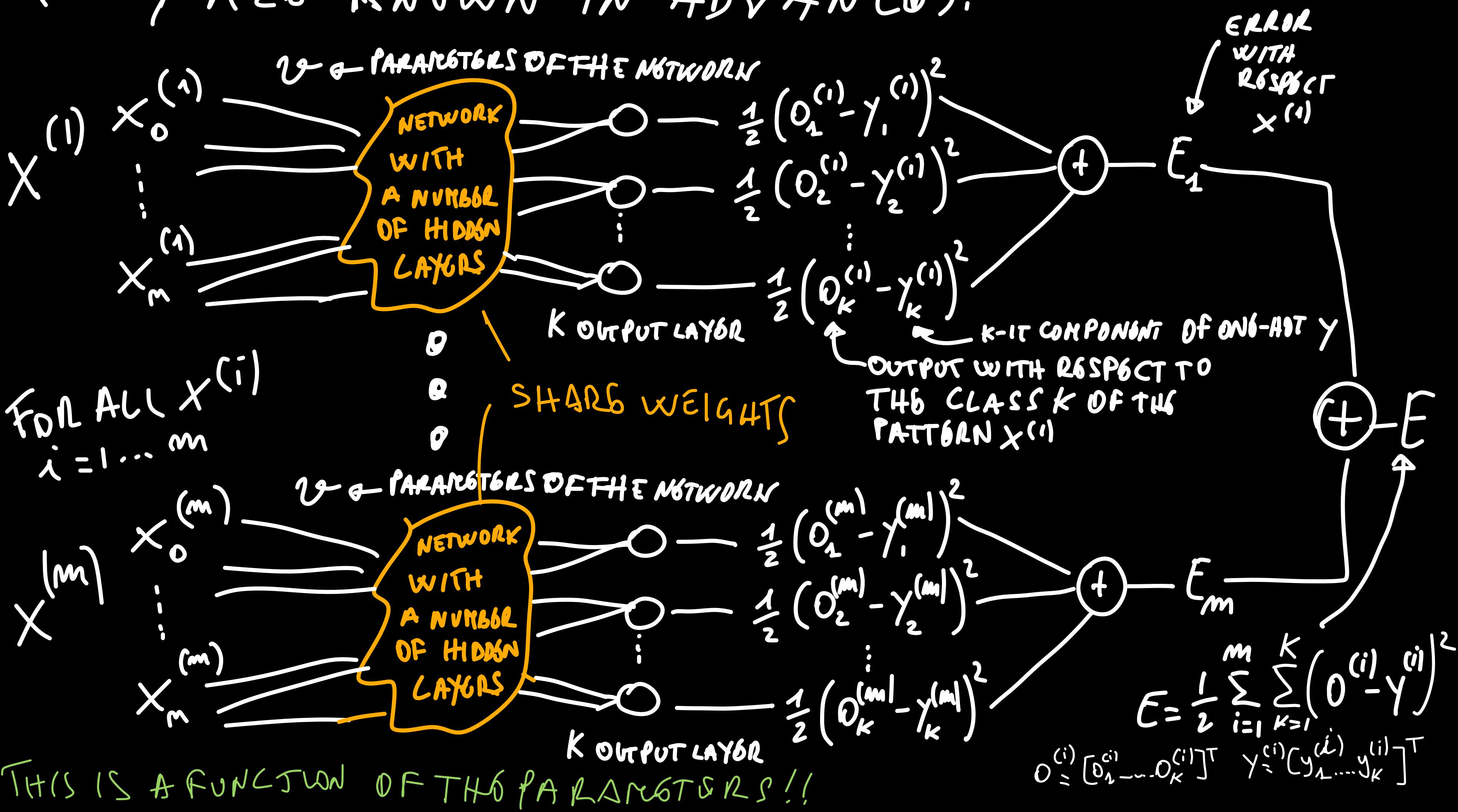
IT IS USED TO FIND THE LOCAL MINIMA OF THE LOSS FUNCTION (e.g., THE WEIGHTS RELATED TO MSE LOCAL MINIMA) CONSIDERING THE TRAINING SET:

- 1) THE NETWORK IS INITIALIZED WITH RANDOM WEIGHTS
(INITIALIZATIONS)
- 2) EVERY ^{TRAINING} INPUT IS FEED TO THE NETWORK
IS MONITORED
LEARNING
AT THE
PLOT OF
THE
LOSS
FUNCTION
(SAVE
WEIGHTS
EVERY
TIME YOU
FIND A
NEW MINIMUM.
- 3) FOR EVERY TRAINING INPUT THE GRADIENT OF THE LOSS FUNCTION IS COMPUTED WITH RESPECT TO THE PARAMETERS OF THE NETWORK
REMING
SYMMETRY
BREAKING
- 4) THE COMPUTED GRADIENTS ARE CUMULATED AND THE AVERAGE IS COMPUTED
- 5) THE GRADIENT OF THE LOSS COMPUTED AT 4) IS USED TO UPGRADE WEIGHTS WITH GRADIENT DESCENT
(ITERATIONS)
- 6) RESTART FROM 2) FOR A GIVEN NUMBER OF EPOCHS
YOU CAN
RESTART FROM
THREE

EXTENDED NETWORK (MULTIPLY REGRESSION & ANPES)

\equiv MSE LOSS

IT IS THE NETWORK WHICH INCLUDES THE LOSS FUNCTION AND IT IS EXTENDED FOR ALL THE TRAINING SAMPLES (THEY ARE KNOWN IN ADVANCE).



DESIRE DATA: Suppose we have P parameters in the net $\vartheta_1 \dots \vartheta_P$ (considering all layers). We wish to compute:

$$\nabla E = \left(\frac{\partial E}{\partial \vartheta_1}, \frac{\partial E}{\partial \vartheta_2}, \dots, \frac{\partial E}{\partial \vartheta_P} \right)$$

Specifically, we need the following value to upgrade the parameters with gradient descent:

$$\Delta \vartheta_p = -\alpha \frac{\partial E}{\partial \vartheta_p} \quad \forall p = 1 \dots P$$

LEARNING RATE

Considering the feedforward network we can only modify parameters to reduce errors (i.e., all other values are fixed)

BACKPROPAGATION EXPLOITS CHAIN RULE:

$$\text{If } F(x) = f(g(x)) \text{ then } F'(x) = g'(x)f'(g(x))$$

DERIVATIVES OF NETWORK FUNCTION

GOAL : FIND A SMART WAY TO COMPUTE GRADIENT

OBSERVATION 1: NETWORK IS A COMPOSITION OF
DRIVABLE FUNCTIONS

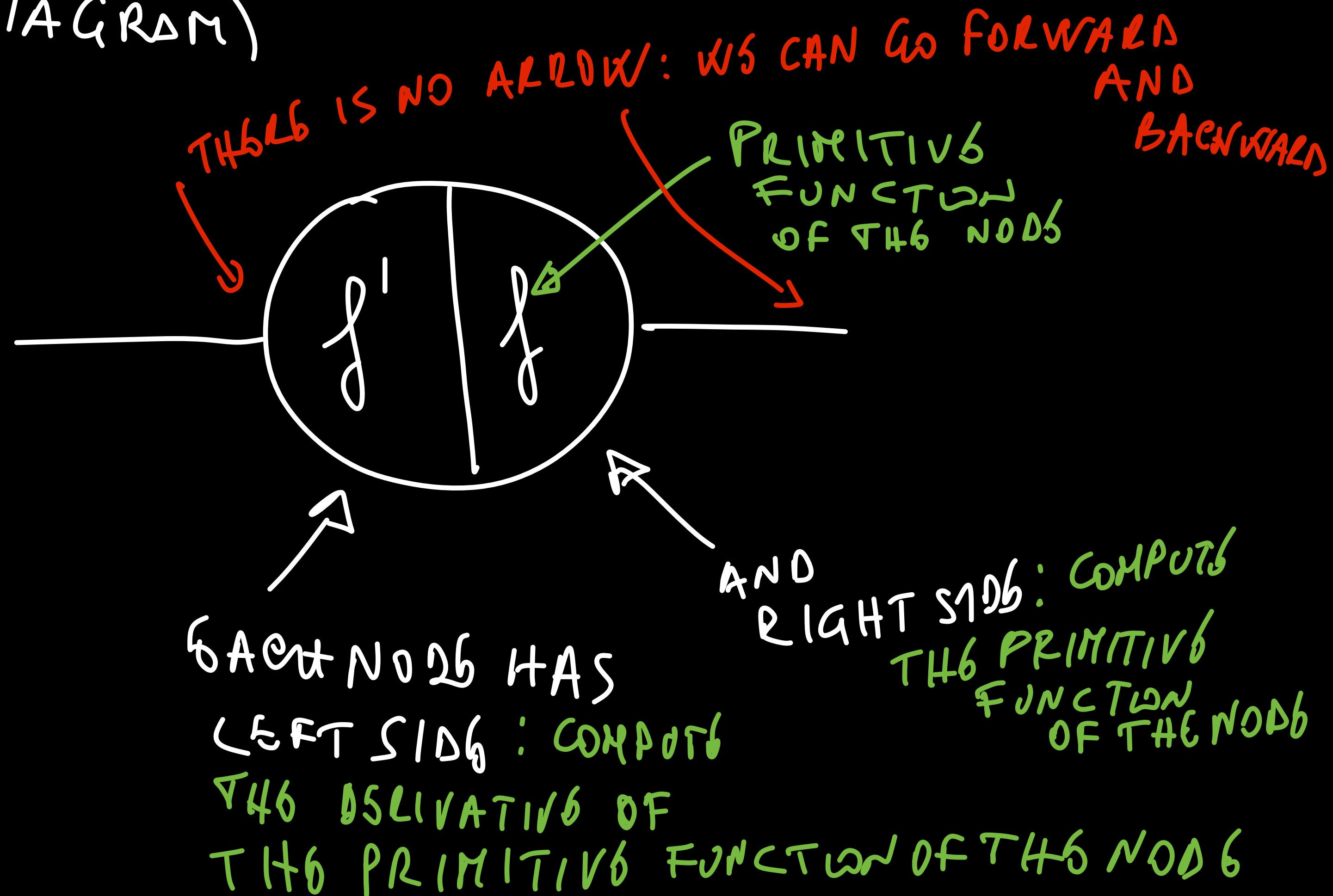
OBSERVATION 2: WE KNOW THE CHAIN RULE PROPERTY

$$f(g(x)) \xrightarrow{\text{DERIVATIVE}} g'(x) f'(g(x))$$

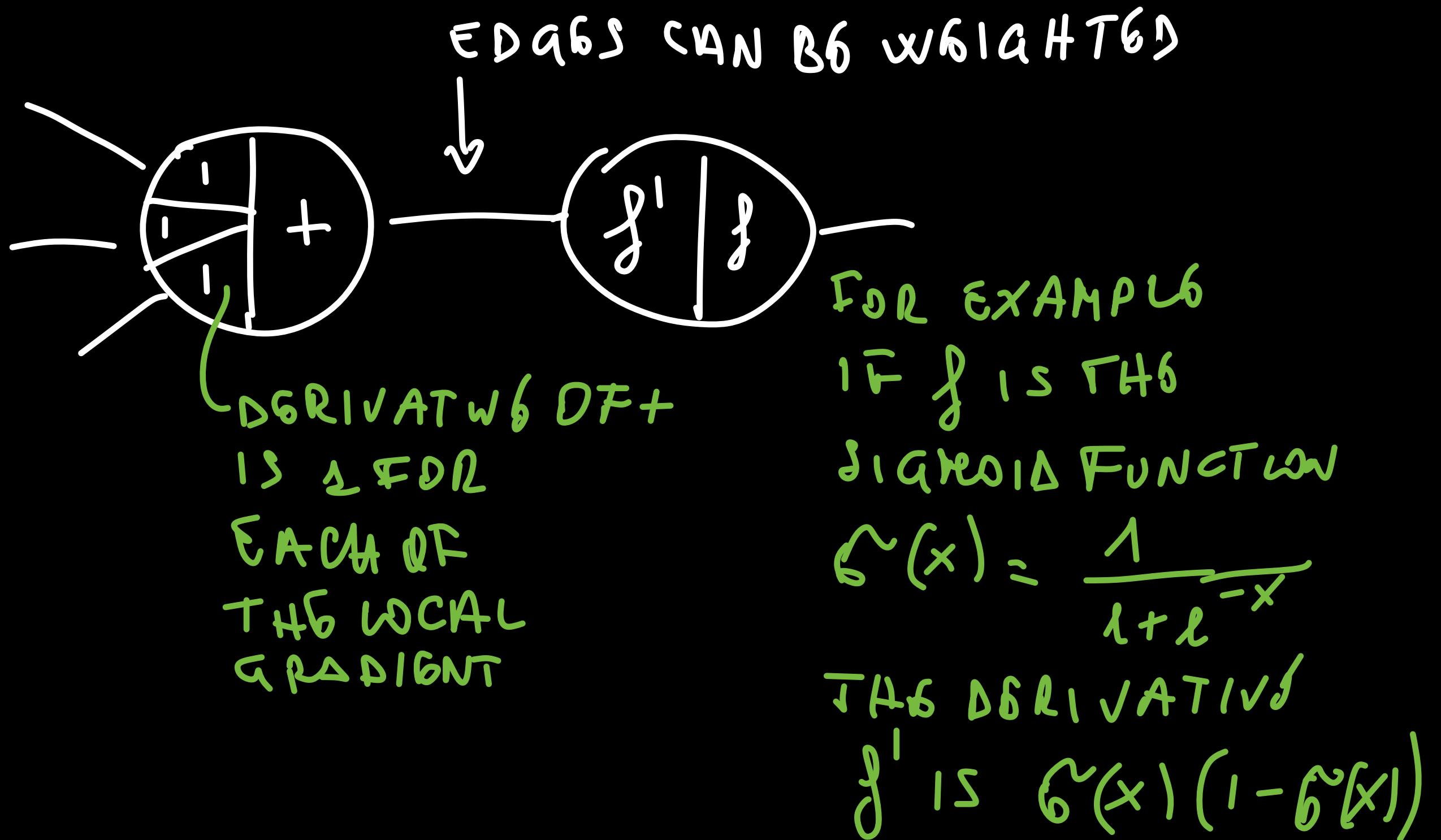
TO FIND A SMART WAY TO COMPUTE GRADIENT ON A
NETWORK WE WILL CONSIDER A VERY OLD CONCEPT...
COMPUTATIONAL GRAPHS !!

COMPOSITE STRUCTURE OF THE NODES

- WE GIVE A COMPOSITE STRUCTURE TO EACH NODE OF THE NETWORK \Rightarrow BACKPROPAGATION DIAGRAM (B-DIAGRAM)

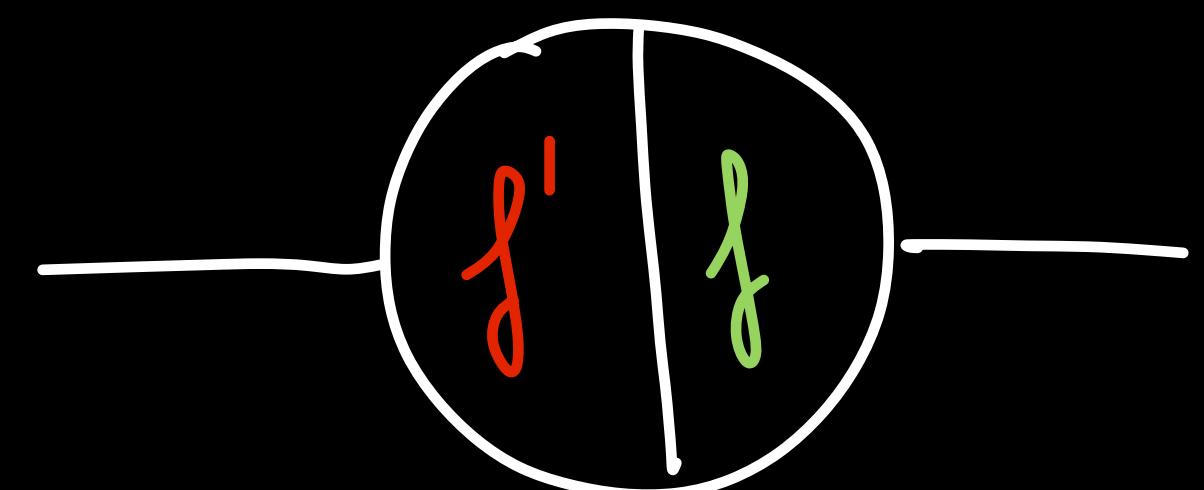


CONSIDERING PRIMITIVE FUNCTIONS, INTEGRATION FUNCTIONS
(SUM) ARE SEPARATED FROM THE ACTIVATION FUNCTIONS



THE NETWORK IS EVALUATED IN TWO STEPS FOR TRAINING:

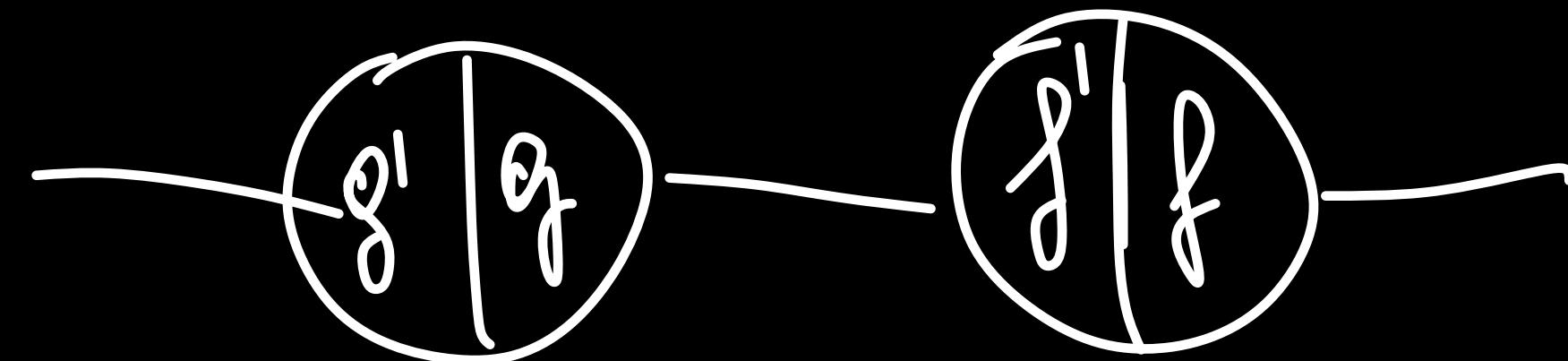
- 1) FEEDFORWARD STEP: INFORMATION GOES FROM LEFT TO RIGHT
TO EACH NODE THE PRIMITIVE FUNCTIONS ARE EVALUATED
AS WELL AS THE DERIVATIVES (OF THE LEFT SIDE).
RESULTS ARE STORED. THE RESULTS OF THE RIGHT SIDE
ARE FORWARDED TO THE CONNECTED UNITS ON THE RIGHT



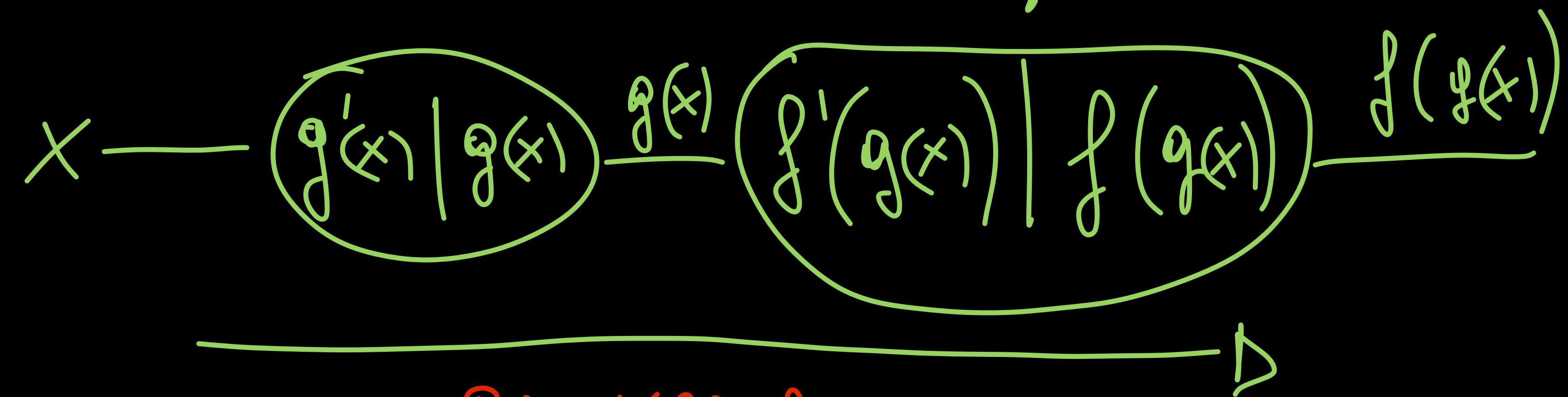
- 2) BACKWARD STEP: IT IS USED TO EVALUATE ERRORS
AND TO PROPAGATE THEM BACKWARD TO UPDATE THE
PARAMETERS. IN THIS CASE THE RESULTS OF THE LEFT
SIDE (DERIVATIVES) ARE USED. THESE ARE THEN USED
TO BE CONSIDERED (WE WILL SEE THEM IN THE FOLLOWING SLIDES).

ANY COMPUTATIONAL GRAPH (JUNCTION COMPOSITION)
CAN BE EVALUATED WITH BACKPROPAGATION TO
COMPUTE DERIVATIVES

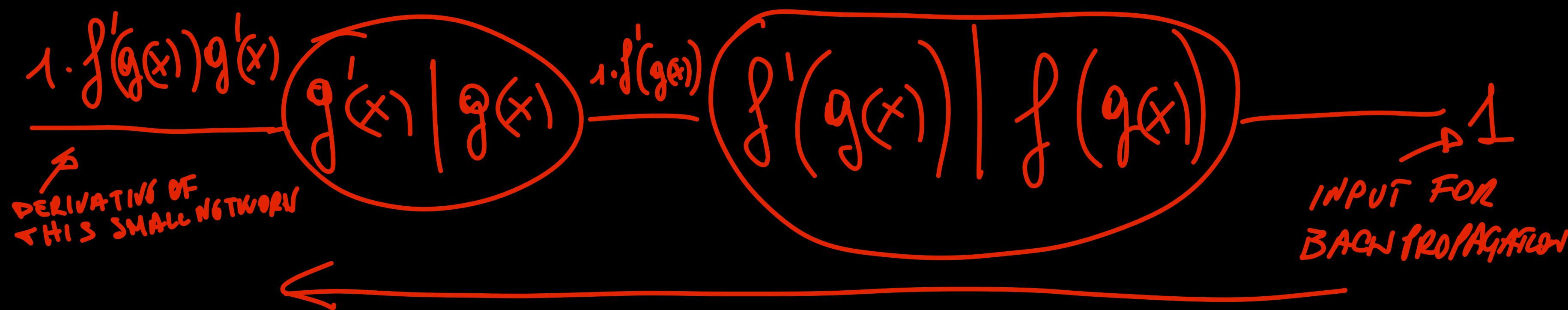
FUNCTIONAL COMPOSITION (C4S6 1)



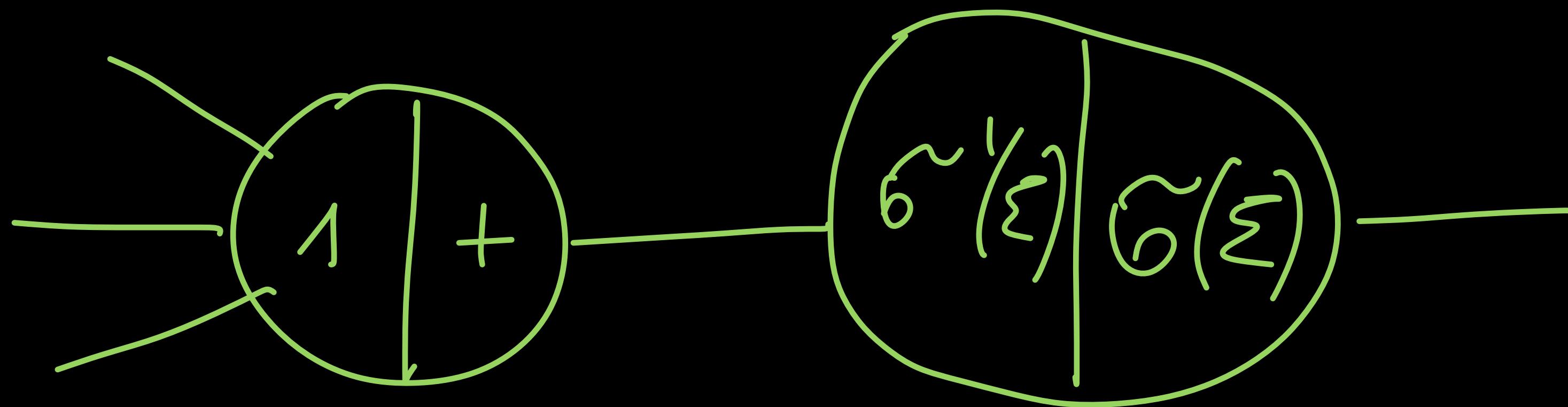
FORWARD PASS / STEP



BACK PROPAGATION



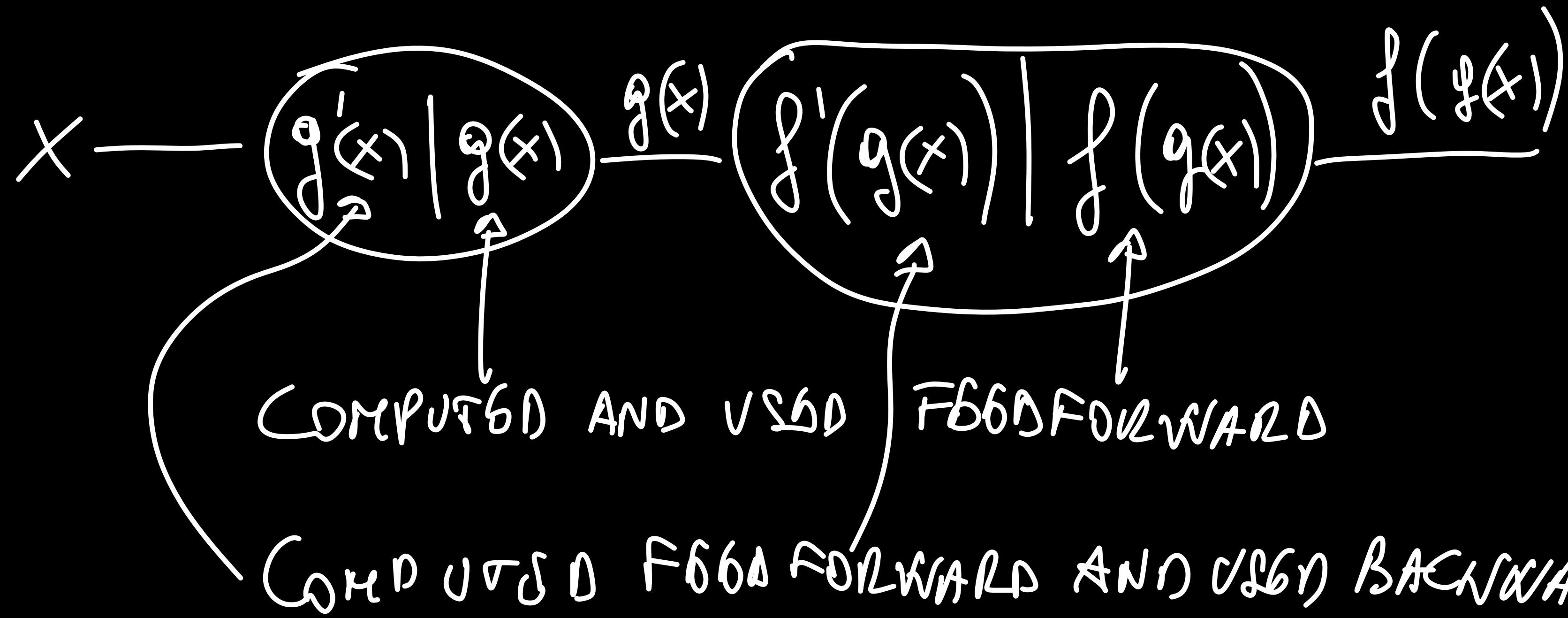
CONSIDERING SIGMOID WHICH INPUT OF IT
IS A LINEAR COMBINATION



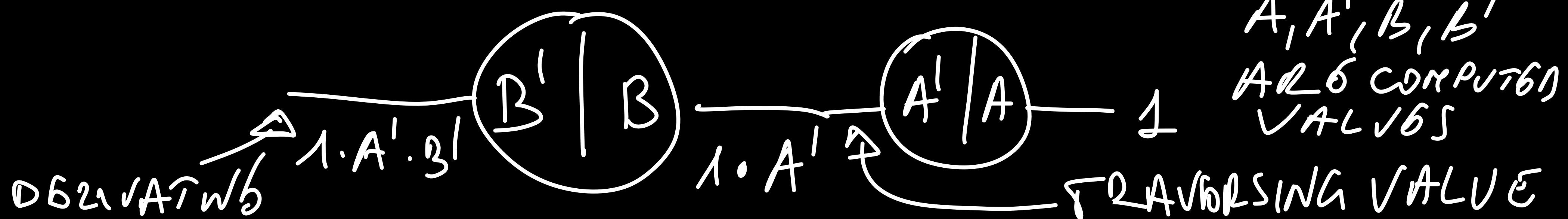
$$G(z) = \frac{1}{1 + e^{-z}}$$

IN THIS CASE z IS A
LINEAR COMBINATION \sum

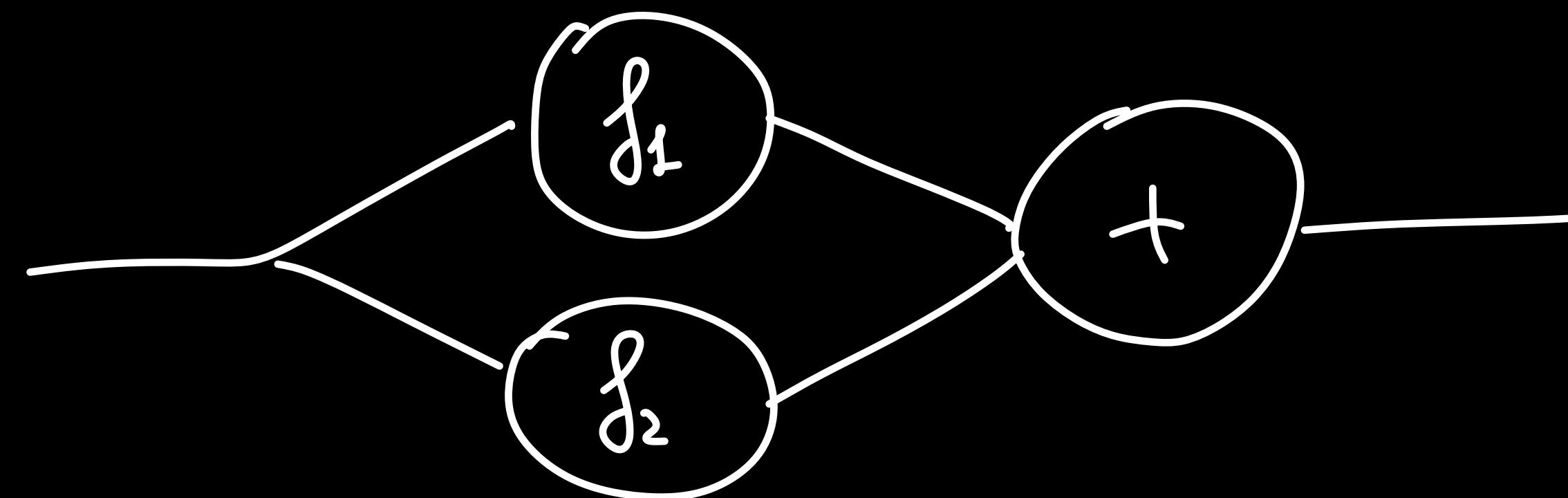
$$\tilde{G}(z) = G(z)(1 - G(z)) = \frac{1}{1 + e^{-z}} \left(1 - \frac{1}{1 + e^{-z}} \right)$$



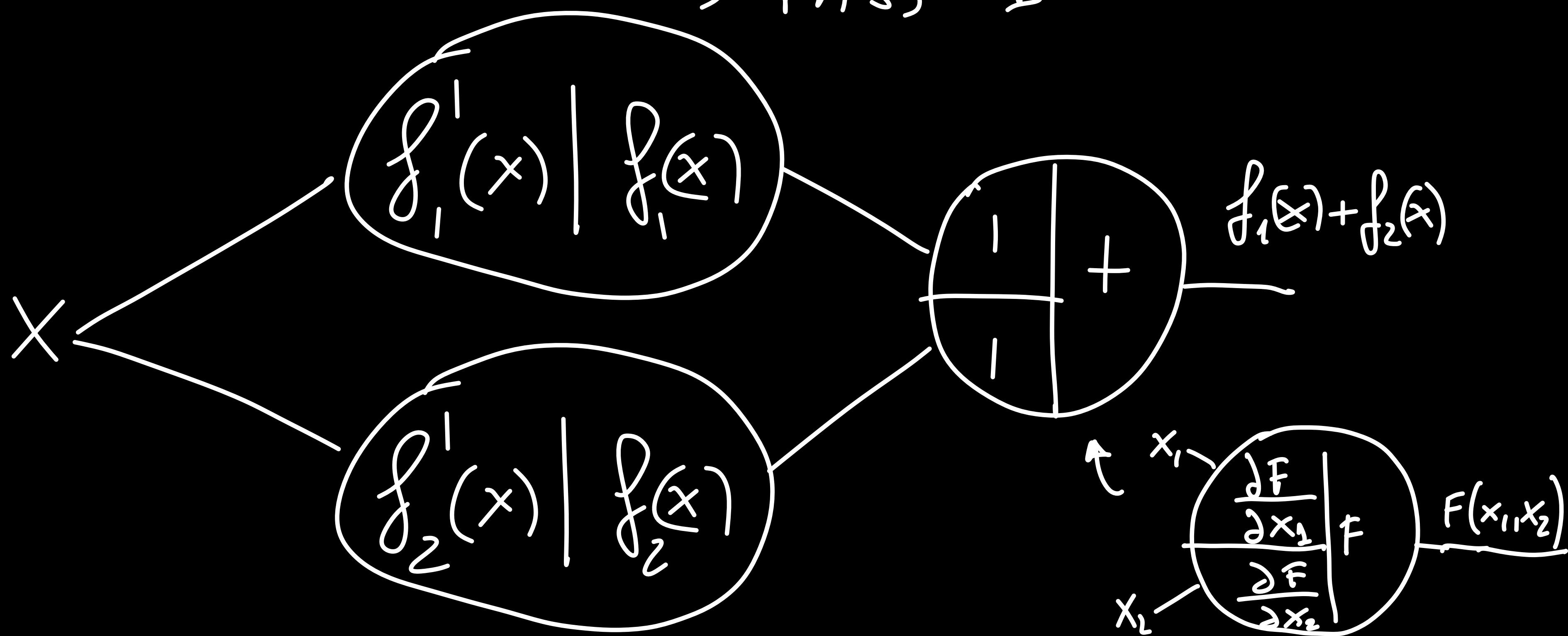
- BACKPROPAGATION PROVIDED A WAY TO IMPLEMENT THE CHAIN RULE PRINCIPLE
- INPUT FROM RIGHT SIDES IS ALWAYS 1 IN GENERAL:



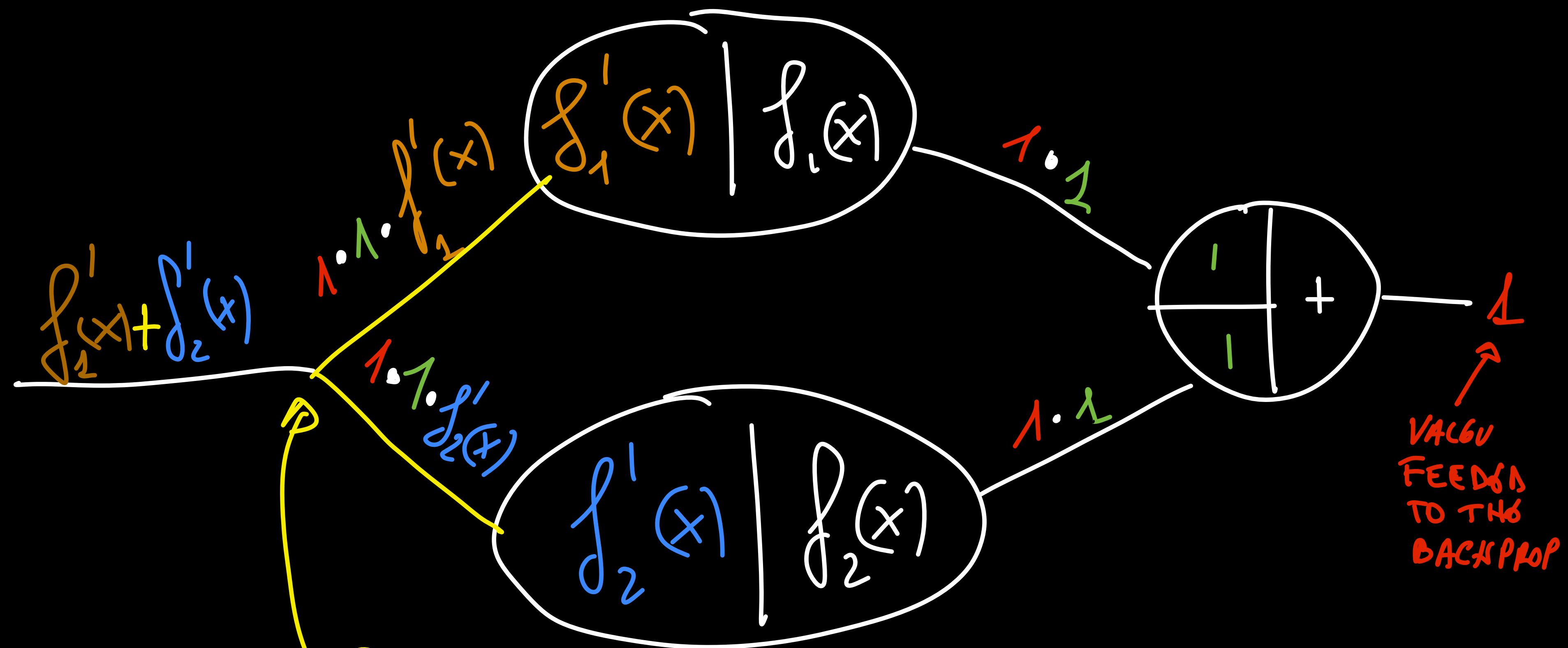
FUNCTIONS ADDITION (CASE 2)



FORWARD PASS →



BACKWARD PASS ←



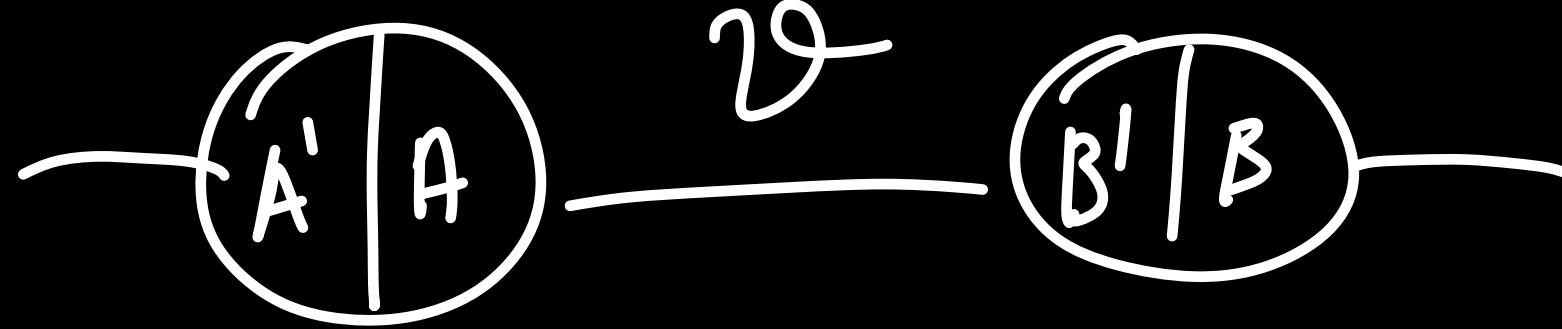
GOING BACKWARD THE VALUES ARE SUMMED WHEN COMING FROM THE SAME BRANCH (NODE)



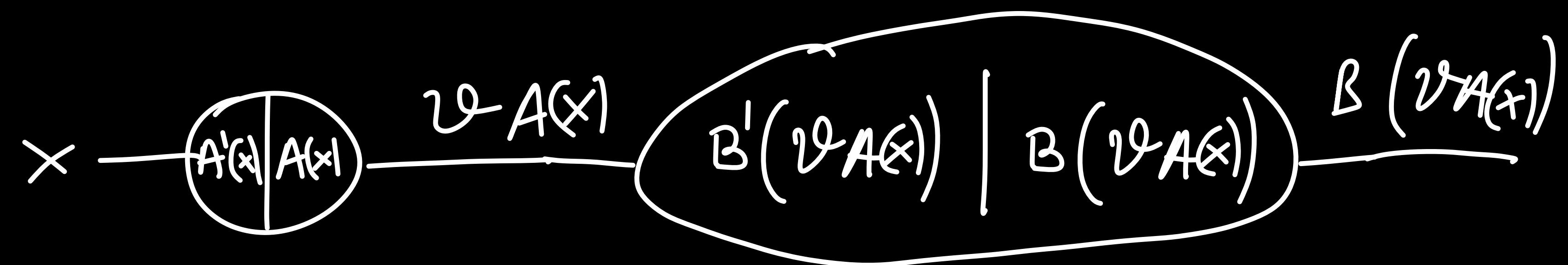
VALUE FEEDS TO THE BACKPROP

WEIGHTED DAGES (CASE 3)

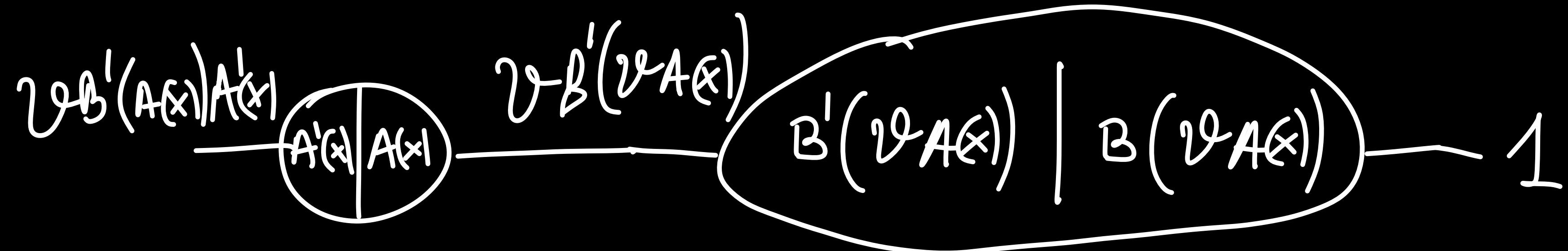
WEIGHTS
ARE NEGATIVELY
FORWARD AND BACKWARD



FORWARD \rightarrow



BACKWARD \leftarrow

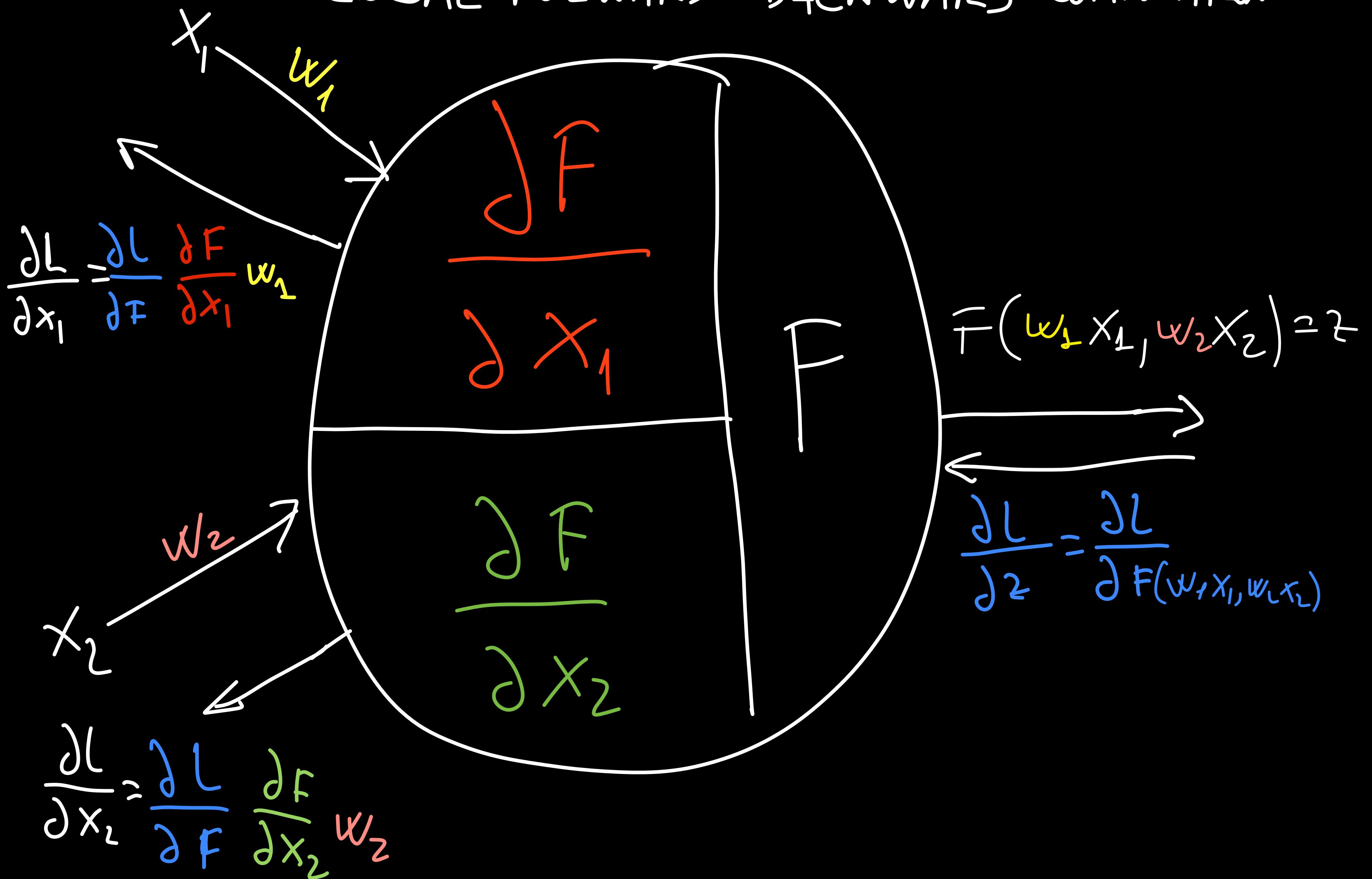


WEIGHTS MODULATES TRANSMITTED VALUES IN BOTH DIRECTIONS

Consider a network f with an input x ,
the derivative is computed as follows:

- 1) FORWARD PASS: The input is evaluated from left
to right considering the network. Values
related primitive functions and their
derivatives are computed and stored.
- 2) BACKWARD PASS: The constant 1 is imposed
as right input of the network. The
network is computed from right to left.
Information which arrived to a given node
are summed and the result multiplied with
the value saved on this left part of the net
and with the partial of the function.
The result is
the derivative of the function with respect to x .

LOCAL FORWARD - BACKWARD COMPUTATION

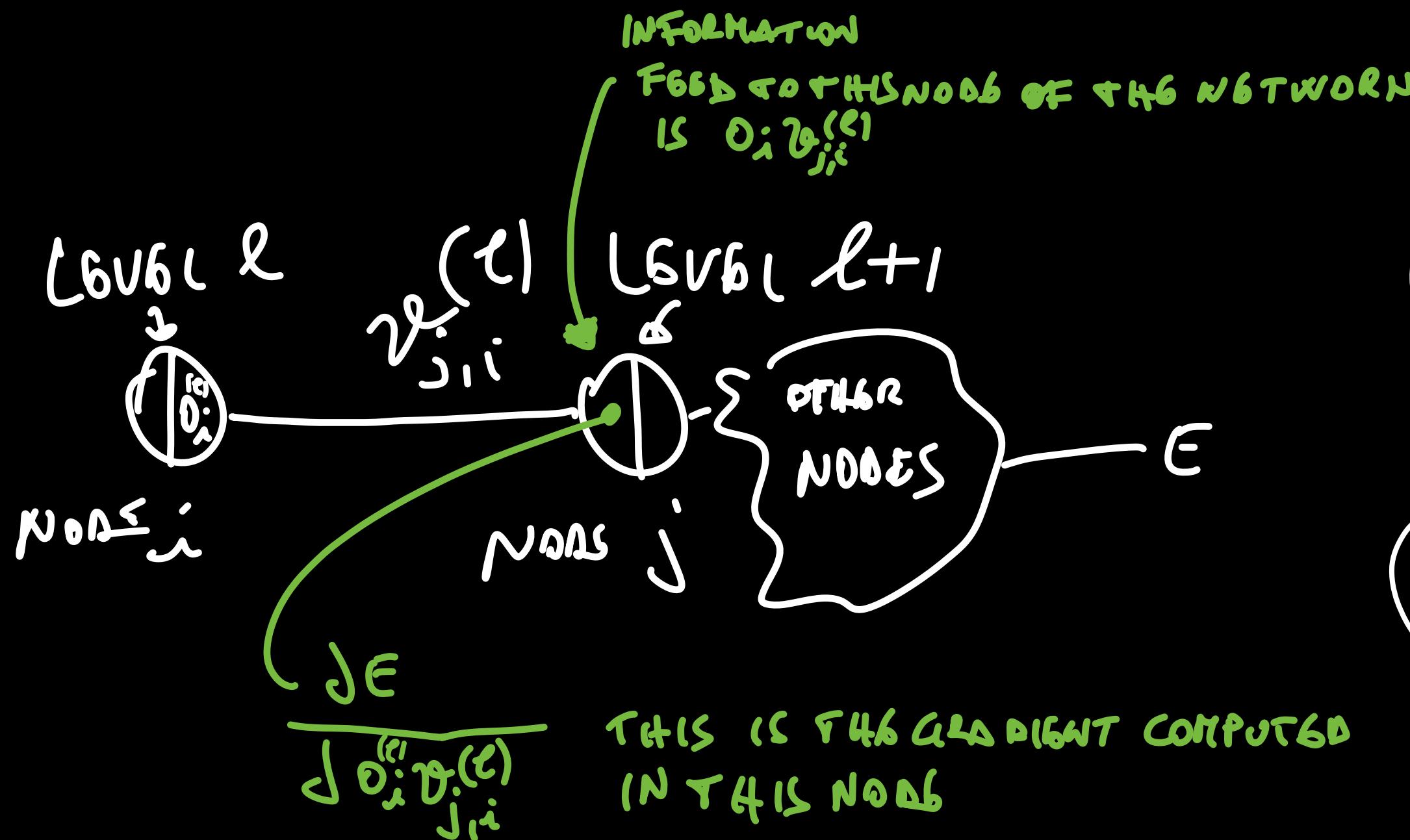


WITH REGARD TO THE EXTENDED NETWORK \mathcal{E}

WE WANT TO COMPUTE GRADIENT FOLLOWING FOR EACH PARAMETER $v_{ji}^{(l)}$

$$\frac{\partial \mathcal{E}}{\partial v_{ji}^{(l)}}$$

LEVEL l
FROM NODE i TO NODE j



$$\frac{\frac{\partial \mathcal{E}}{\partial O_i v_{ji}^{(l)}}}{\frac{\partial O_i}{\partial O_k} v_{jk}^{(l)}}$$

NOTE THAT $O_i^{(e)}$ IS ALREADY COMPUTED DURING BACKPROPAGATION.

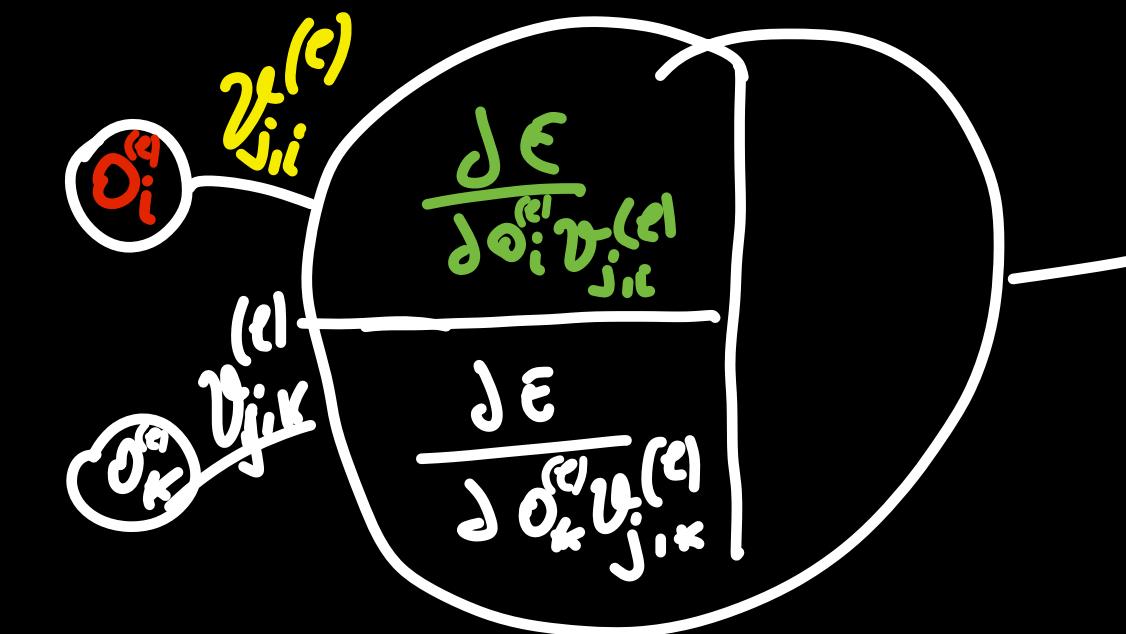
HENCE IT IS A CONSTANT DURING THE FORWARD PASS AND A CLASS TO WRITE THE FOLLOWING:

$$\frac{\partial E}{\partial v_{j,i}^{(e)}} = \frac{\partial E}{\partial O_i^{(e)} v_{j,i}^{(e)}}$$

$$\frac{\partial E}{\partial O_i^{(e)} v_{j,i}^{(e)}} =$$

$$= O_i^{(e)}$$

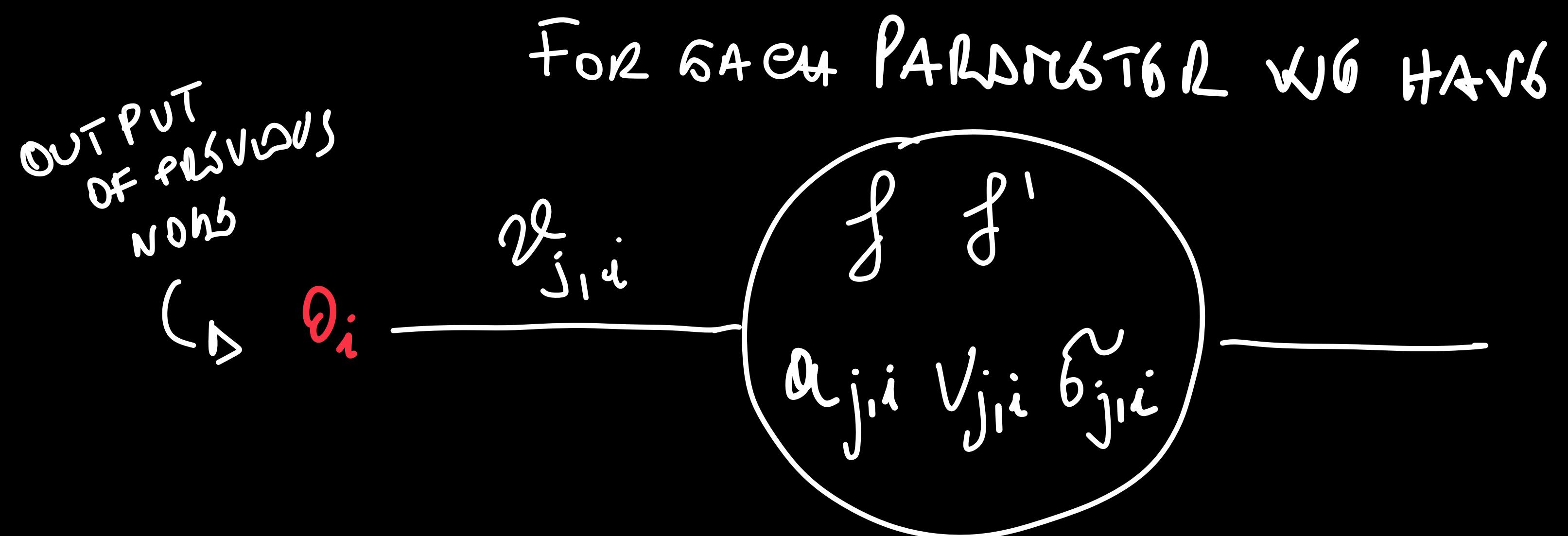
$\frac{\partial E}{\partial O_i^{(e)} v_{j,i}^{(e)}}$



OUR OBSERVATION

FOR THE
CATION
RUN

BACKPROPAGATION
ERROR AT NODE j
(USUALLY DENOTED
WITH $\delta_j^{(e)}$)



f : PRIMITIVE FUNCTION TO BE COMPUTED ON THE INPUT $\theta_i v_{j,i}$

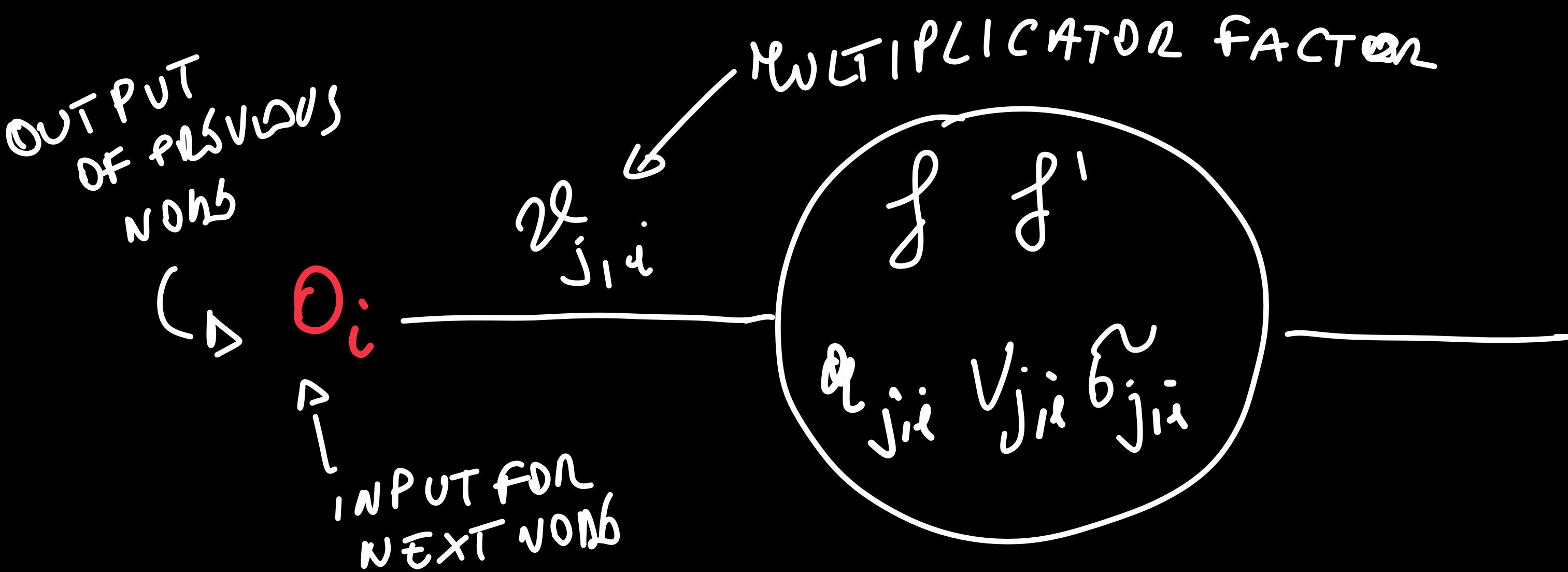
f' : DERIVATIVE OF THE PRIMITIVE FUNCTION

f' : TO BE COMPUTED ON THE INPUT $\theta_i v_{j,i}$

$a_{j,i}$: VALUE OF THE PRIMITIVE FUNCTION f COMPUTED
ON $\theta_i v_{j,i}$ IN THE FORWARD PASS

$v_{j,i}$: VALUE OF THE DERIVATIVE PRIMITIVE FUNCTION f' COMPUTED
ON THE VALUE $\theta_i v_{j,i}$ DURING FORWARD STEP

$\tilde{e}_{j,i}$: CUMULATIVE RESULT OF BACKWARD COMPUTATION
WHICH IS THE GRADIENT OF THE ERROR RESPECT TO $\theta_i v_{j,i}$



THIS MEANS THAT FOR EACH NODE WE CAN EASILY COMPUTE THE GRADIENT OF THE ERROR FUNCTION E WITH RESPECT TO THE PARAMETER ASSOCIATED TO THAT NODE

$$\frac{\partial E}{\partial v_{j,i}} = O_i b_{j,i} \Rightarrow \text{THIS CAN BE USED FOR GRADIENT DESCENT}$$

$$\Delta v_{j,i} = -\lambda O_i b_{j,i}$$

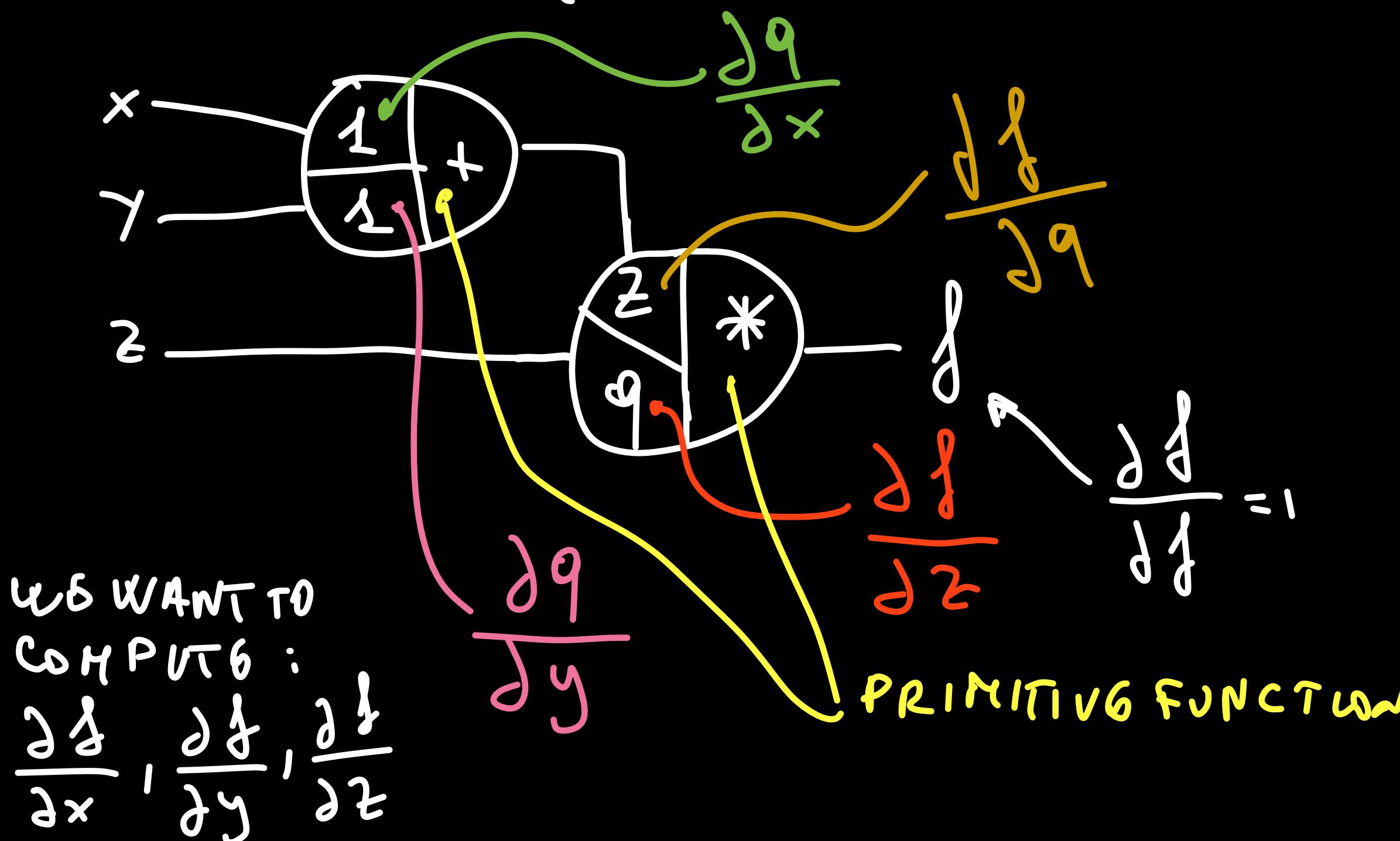
THIS CAN BE DONE FOR EACH PARAMETER OF THE NET

EXAMPLE OF FORWARD AND BACKWARD PROPAGATION

FUNCTION: $f(x, y, z) = (x+y)*z$

INPUT: $x = -2$ $y = 5$ $z = -4$

COMPUTATIONAL GRAPH



"SUBFUNCTIONS"

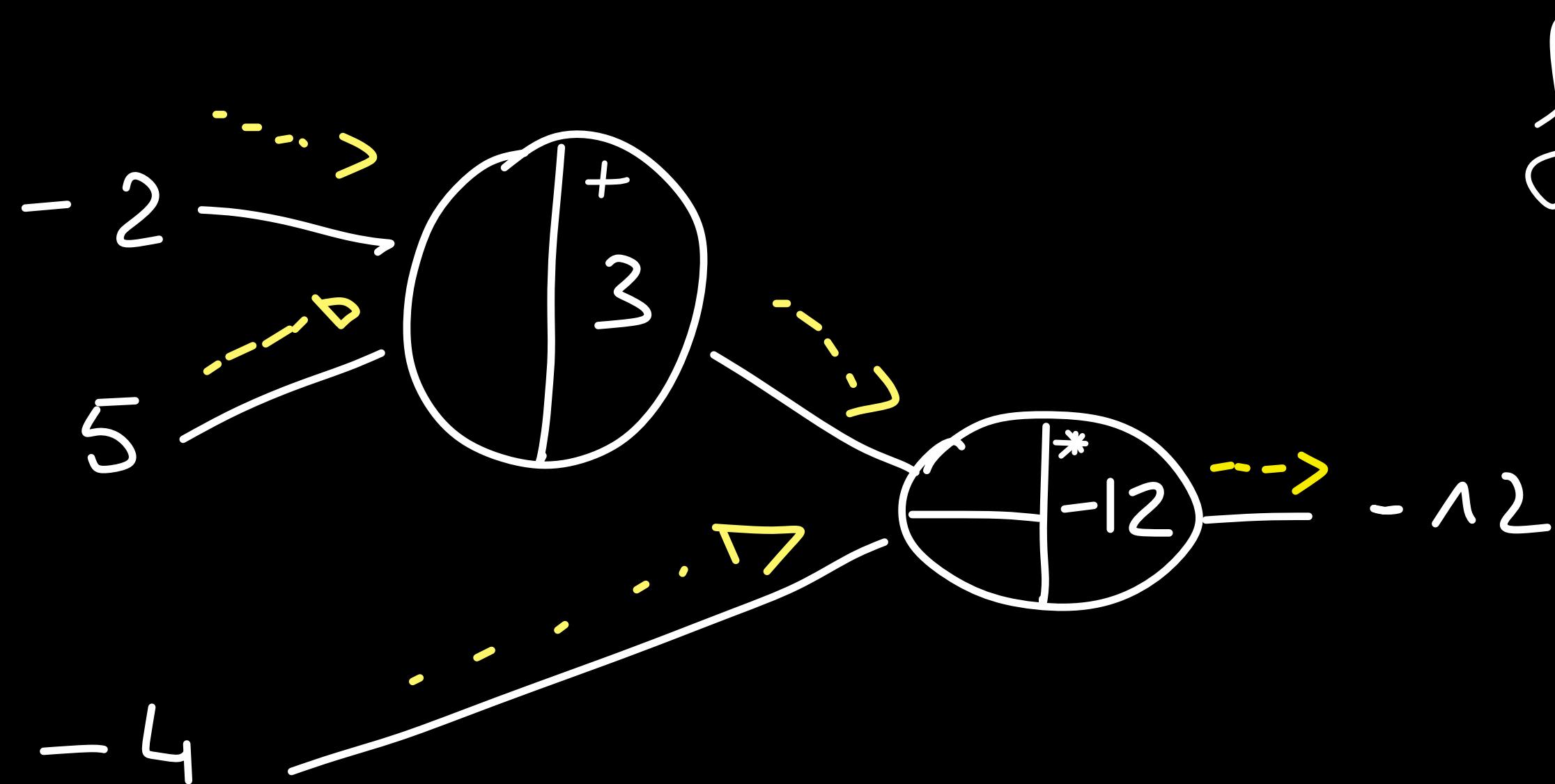
$$q = x + y \quad \begin{cases} \frac{\partial q}{\partial x} = 1 \\ \frac{\partial q}{\partial y} = 1 \end{cases}$$

$$f = q * z \quad \begin{cases} \frac{\partial f}{\partial q} = z \\ \frac{\partial f}{\partial z} = q \end{cases}$$

$$\frac{\partial f}{\partial q} = z$$

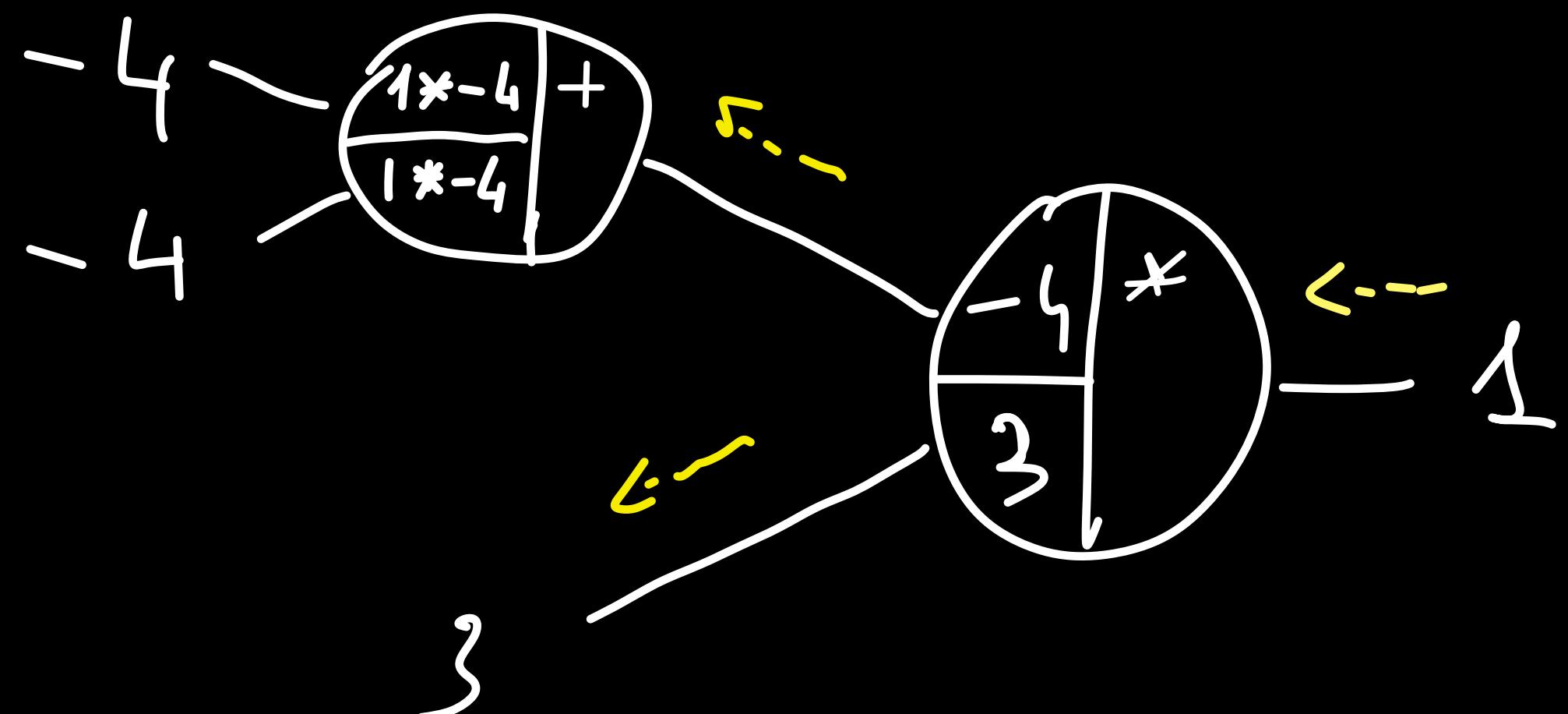
$$\frac{\partial f}{\partial z} = q$$

FORWARD PASS ON THE VALUES $x = -2$ $y = 5$ $z = 3$



$$f(x, y, z) = (x + y) * z$$

BACKWARD PASS CONSIDERING THE ABSOLUTE INPUT

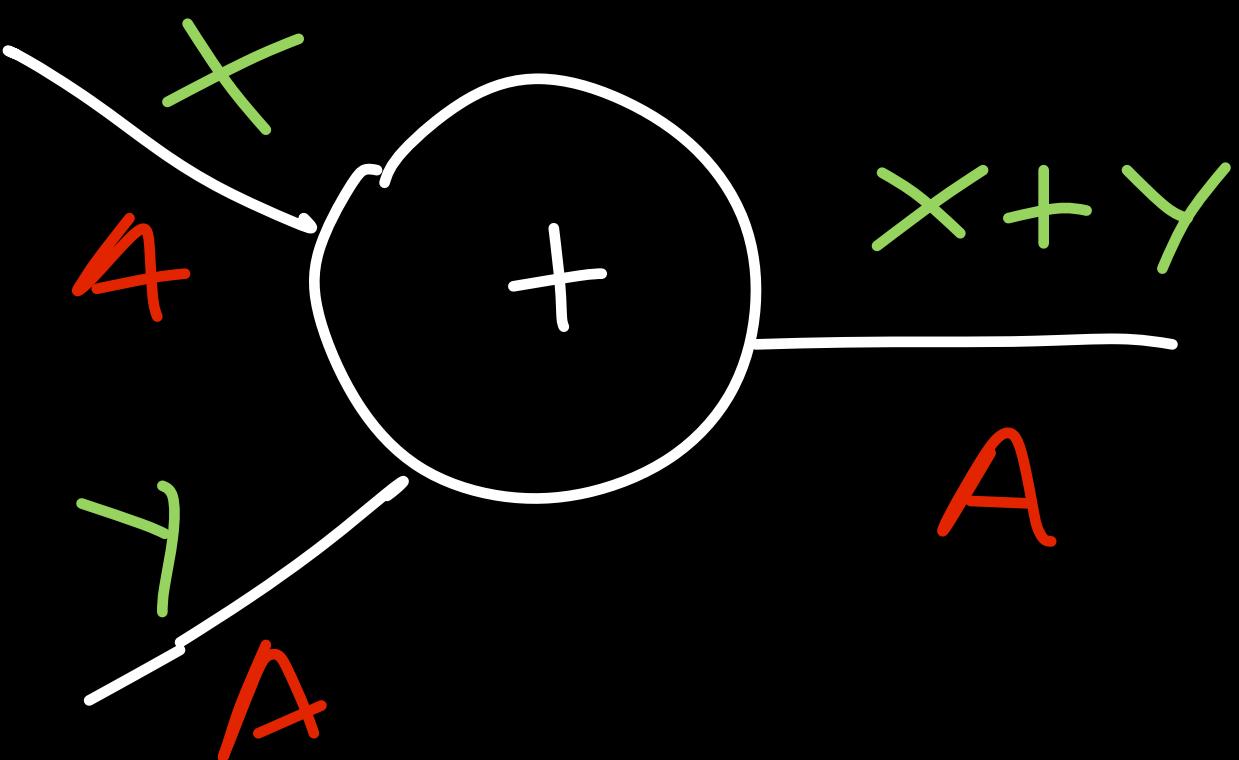


$$\frac{\partial f}{\partial x} = -4$$

$$\frac{\partial f}{\partial y} = -4$$

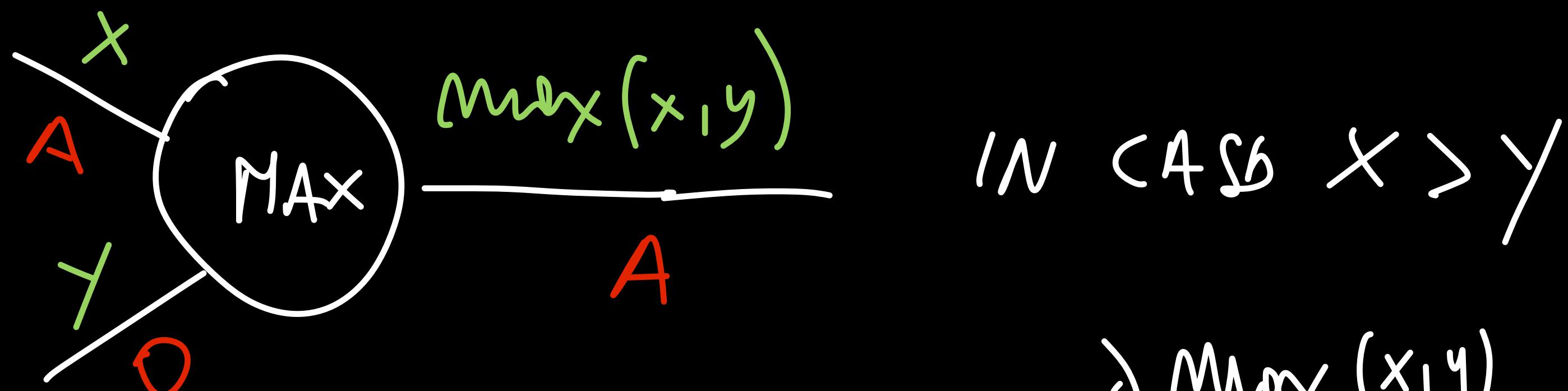
$$\frac{\partial f}{\partial z} = 3$$

ADD GATE : GRADIENT DISTRIBUTION



FORWARD
BACKWARD

MAX GATE : GRADIENT ROUTER

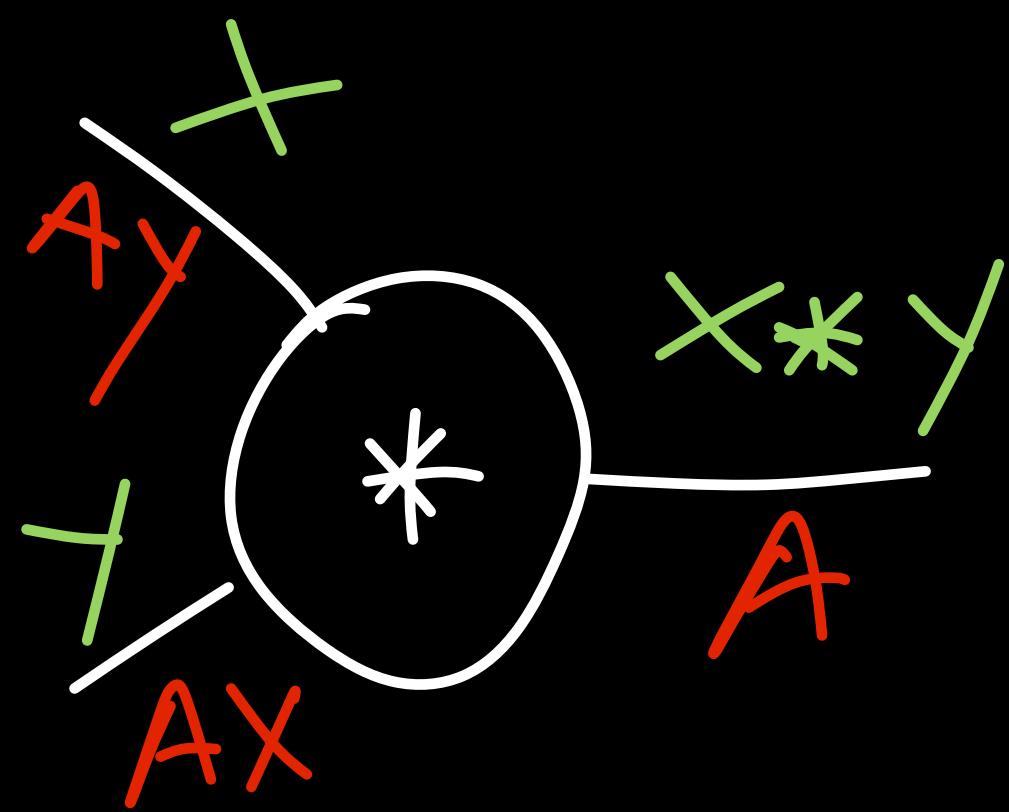


IN CASE $x > y$

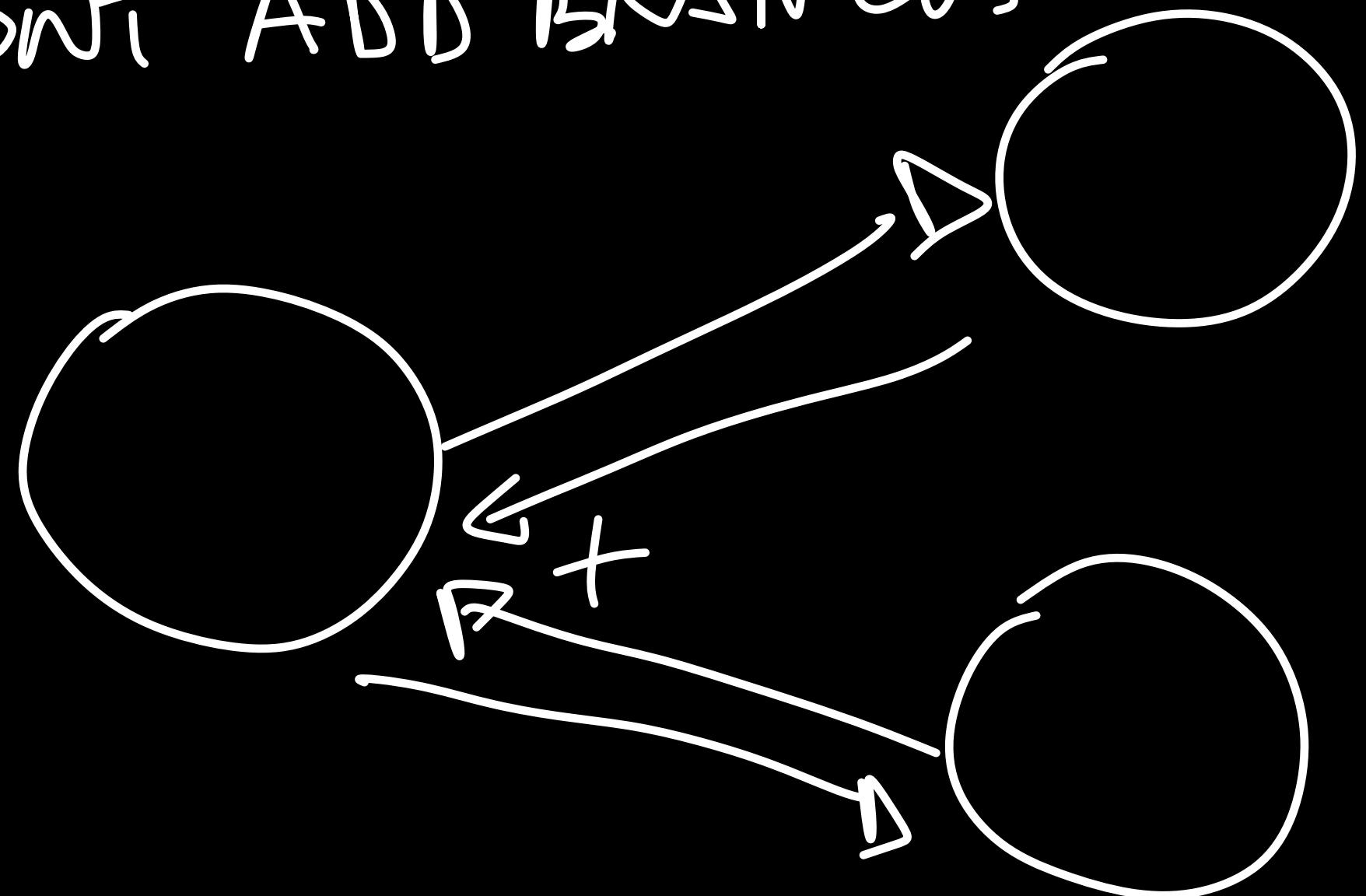
$$\frac{\partial \max(x, y)}{\partial x} = 1 \quad \{x > y\}$$

$$\max(x, y) \Rightarrow \frac{\partial}{\partial x} \quad \frac{\partial \max(x, y)}{\partial y} = 1 \quad \{y > x\}$$

ReLU GATE : GRADIENT SWITCH



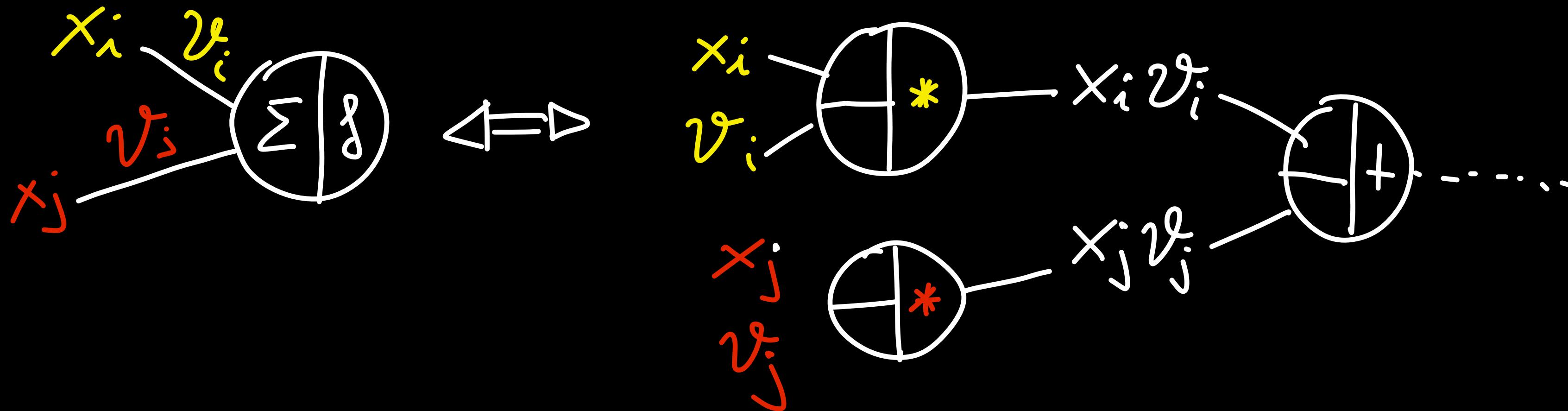
GRADIENT ADD BRANCHES



EXAMPLE WITH WEIGHTS - LOGISTIC REGRESSION

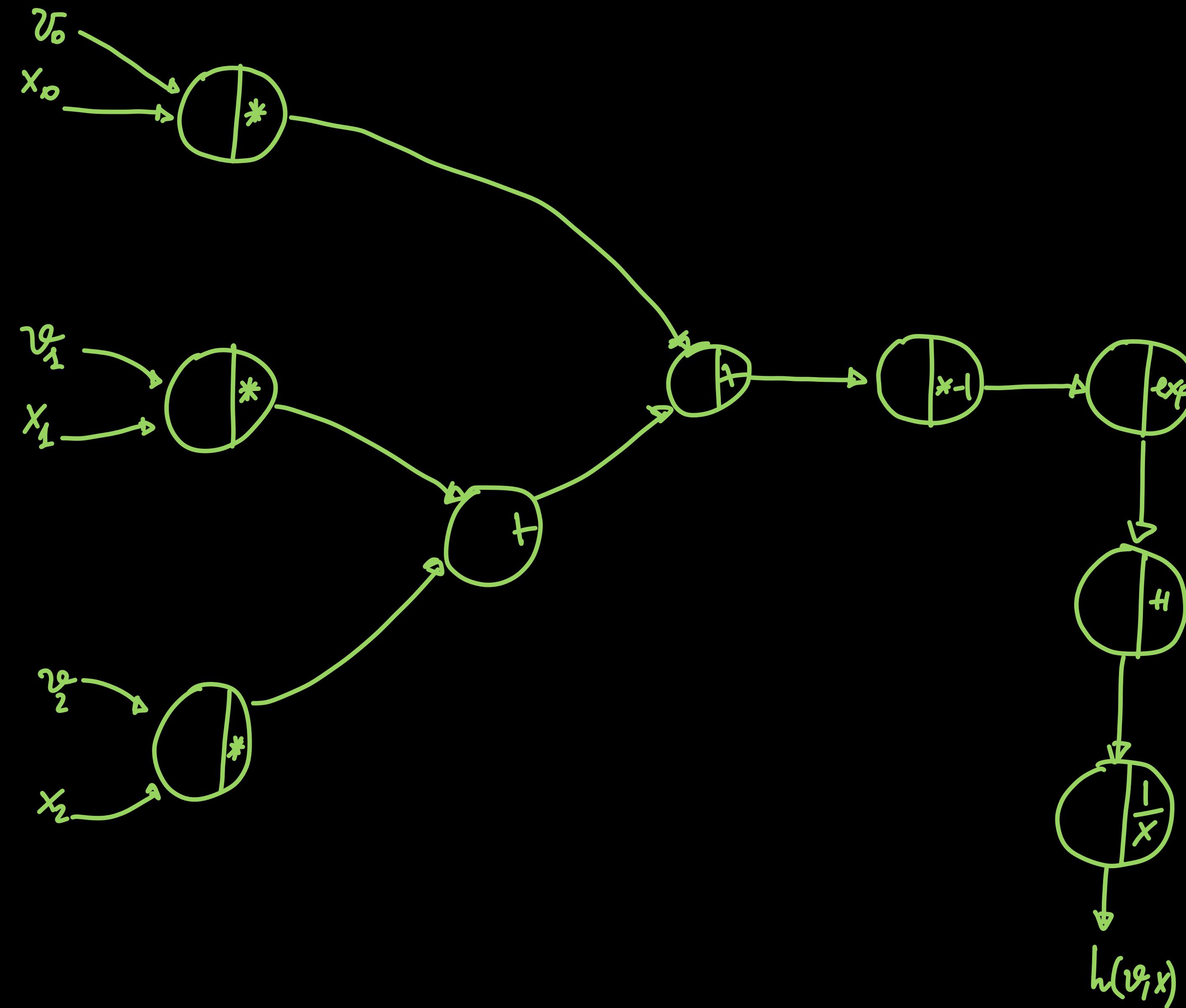
$$h(v_i, x) = \frac{1}{1 + e^{-(v_0 x_0 + v_1 x_1 + v_2 x_2)}}$$

$$v = \begin{bmatrix} -3 \\ -3 \\ 2 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ -2 \\ -1 \end{bmatrix}$$

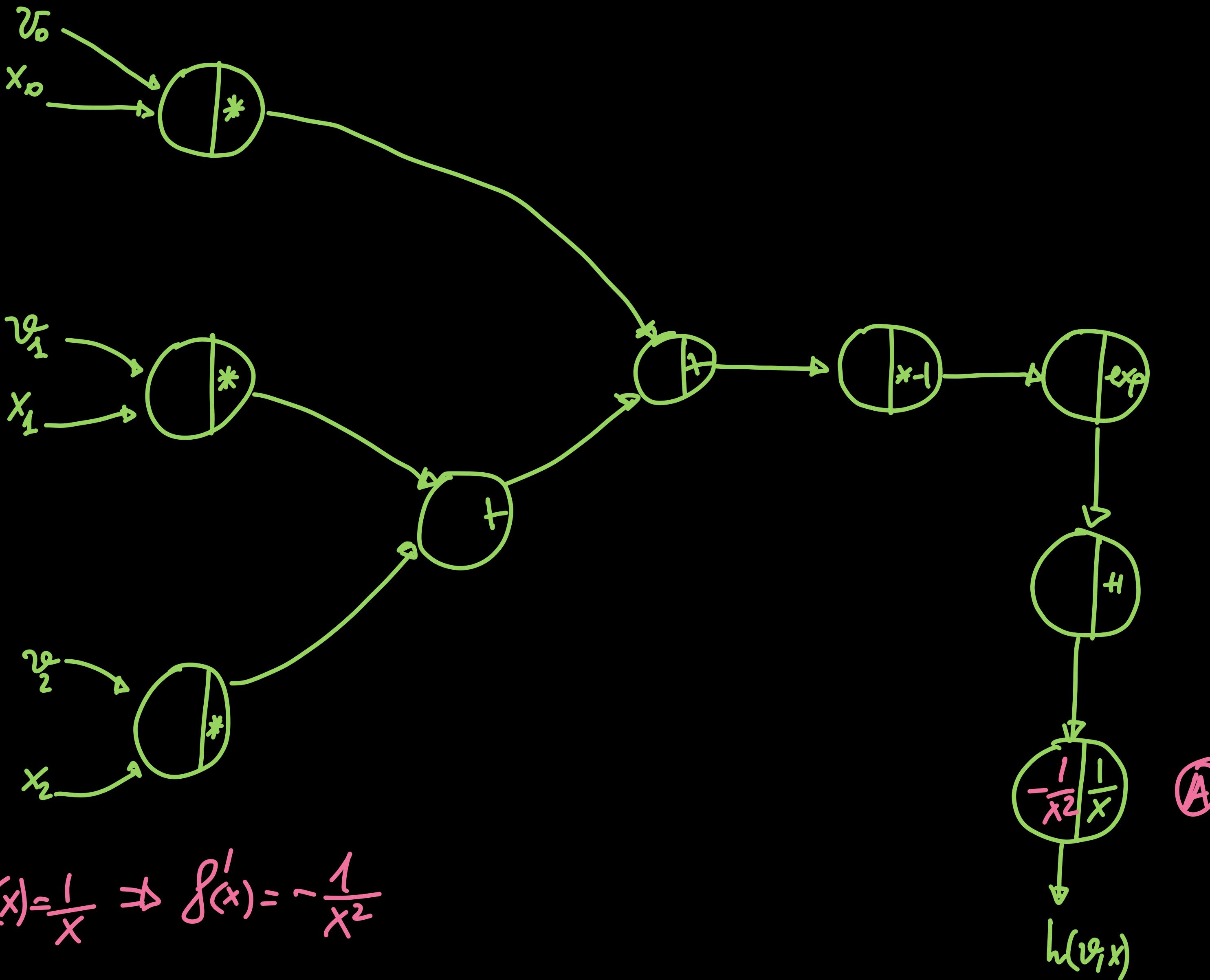


- ① DRAW THE COMPUTATIONAL GRAPH
- ② INSERT DERIVATIVES
- ③ ADD FORWARD VALUES
- ④ ADD BACKWARD VALUES
- ⑤ COMPUTE FINAL GRADIENTS WITH RESPECT TO v_j

① DRAW THE COMPUTATIONAL GRAPH



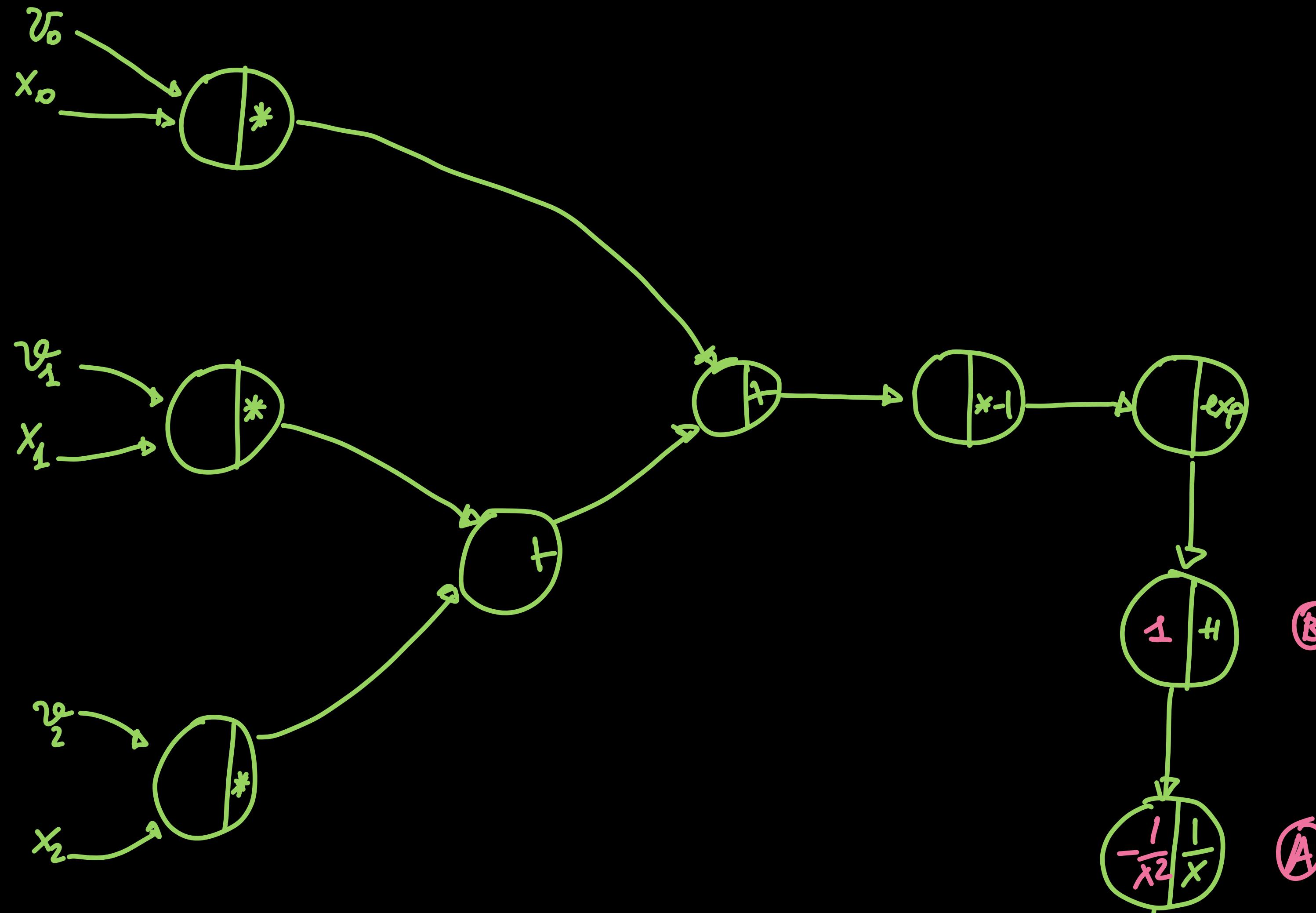
② INSERT DERIVATIVES



$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$h(v, x)$$

② INSERT DATA INTO

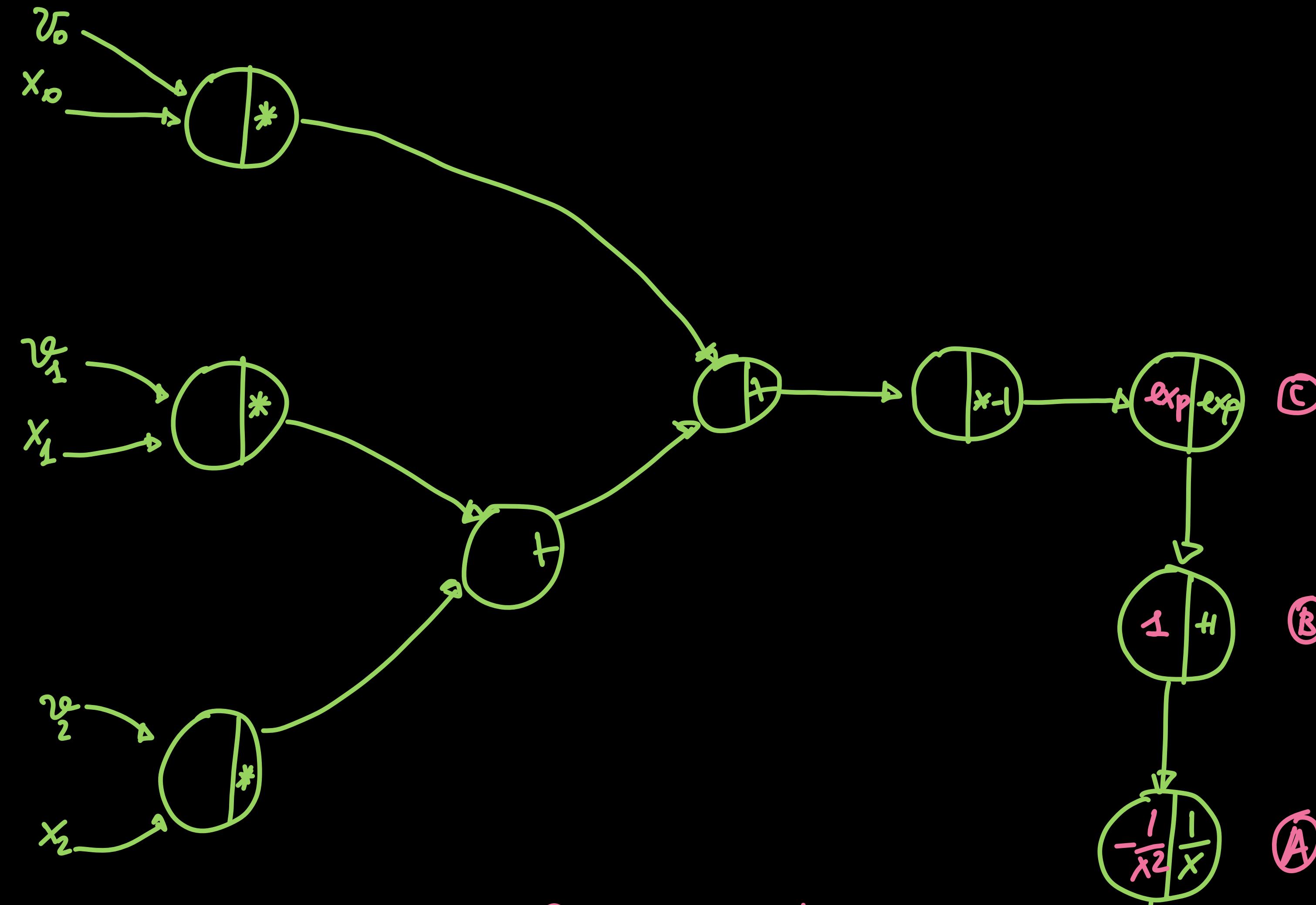


$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1$$

$$h(v_1, x)$$

② INSERT DEVIATIVES



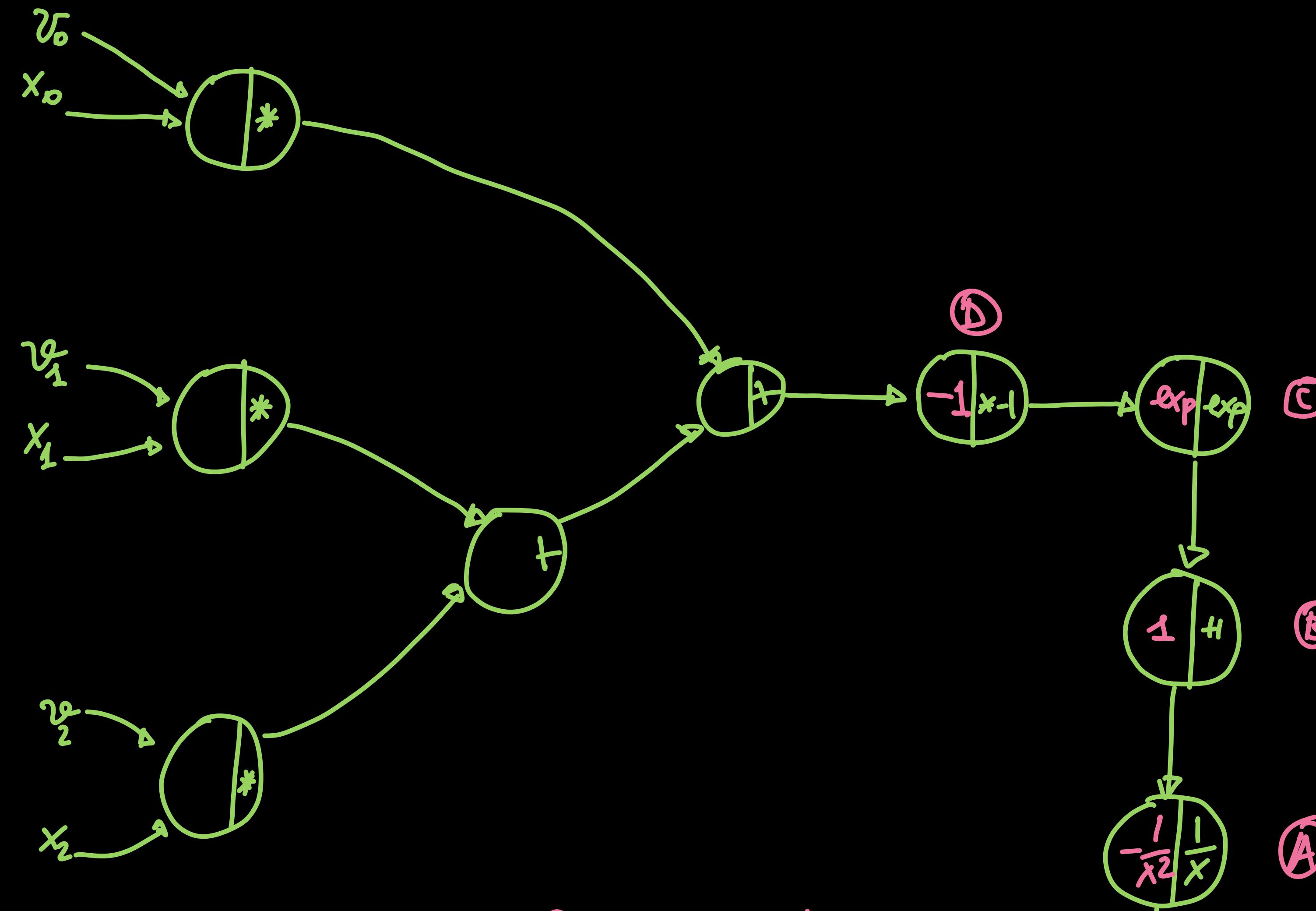
$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1$$

$$\textcircled{C} \quad f(x) = e^x \Rightarrow f'(x) = e^x$$

$$h(v_1, x)$$

② INSERT DEVIATIVES



$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

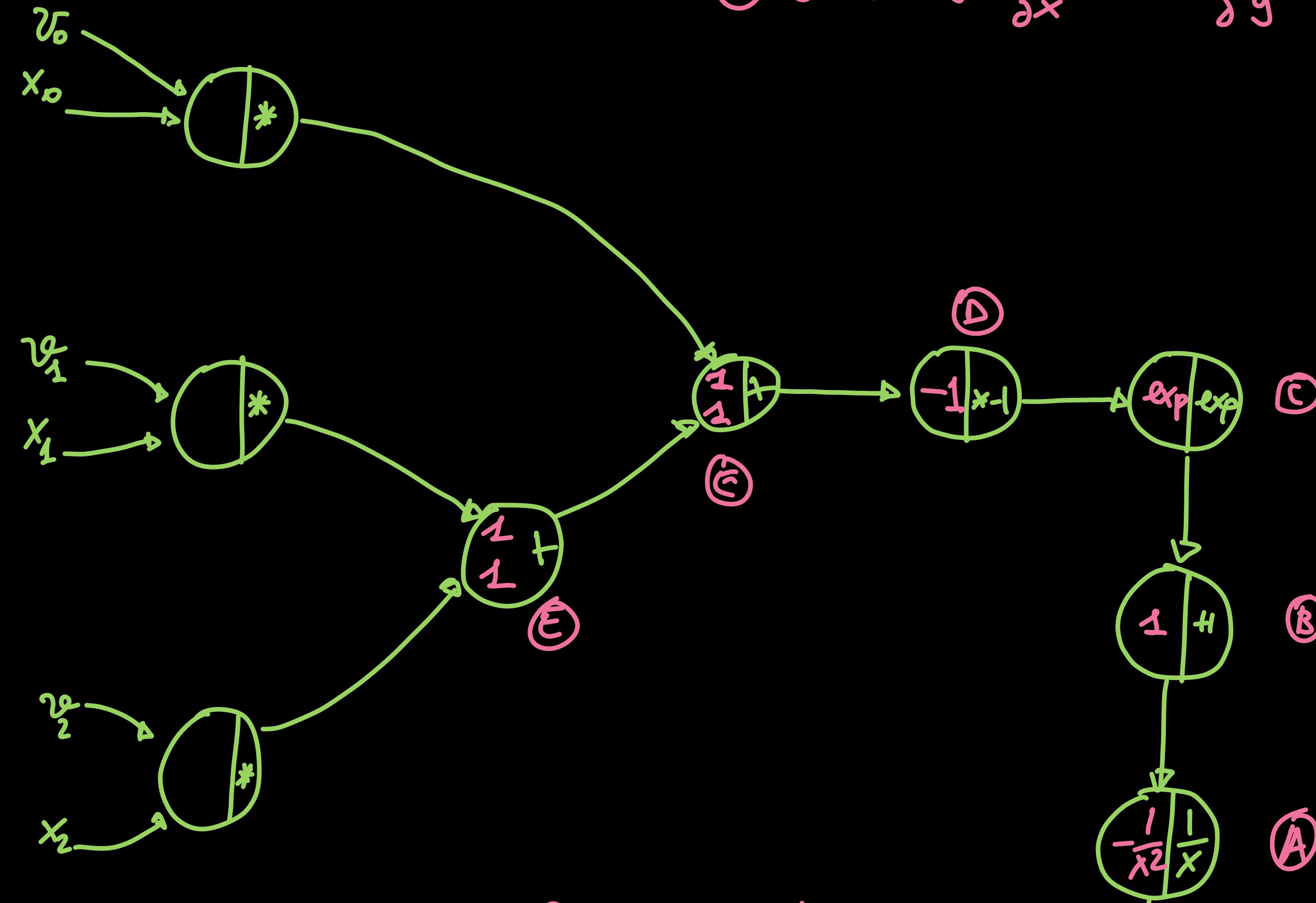
$$\textcircled{C} \quad f(x) = e^x \Rightarrow f'(x) = e^x$$

$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1 \quad \textcircled{D} \quad f_a(x) = ax \quad f'_a(x) = a$$

$$h(v_i, x)$$

② INSERT DATA INTO

$$⑥ f(x,y) = xy \quad \frac{\partial f}{\partial x} = 1 \quad \frac{\partial f}{\partial y} = 1$$



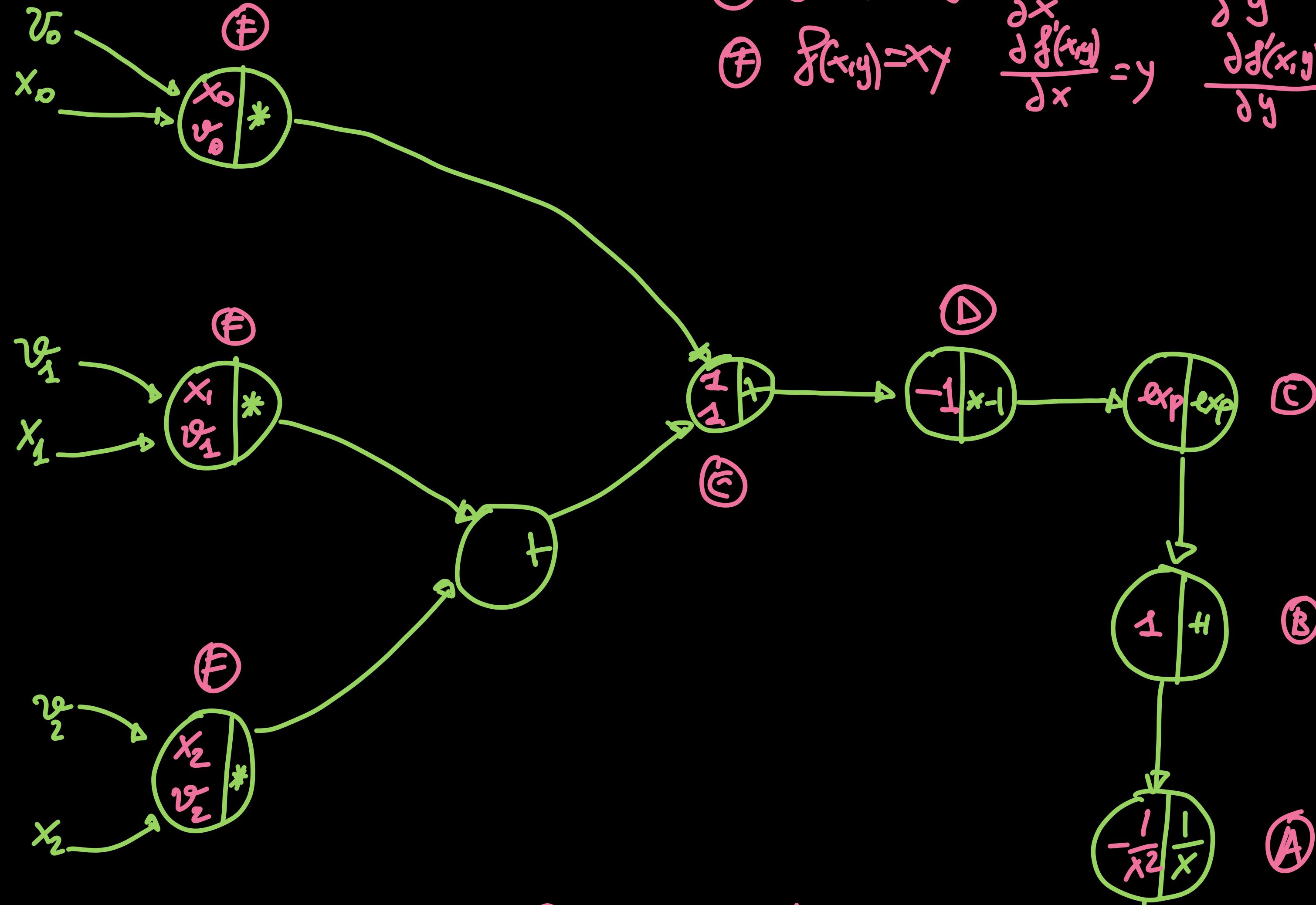
$$① f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$③ f(x) = c + x \Rightarrow f'(x) = 1 \quad ④ f_a(x) = ax \quad f'_a(x) = a$$

$$② f(x) = e^x \Rightarrow f'(x) = e^x$$

$$h(v_1, x)$$

② INSERT D621 VATIVS



$$\textcircled{G} \quad f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = 1 \quad \frac{\partial f'(x,y)}{\partial y} = 1$$

$$\textcircled{H} \quad f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = y \quad \frac{\partial f'(x,y)}{\partial y} = x$$

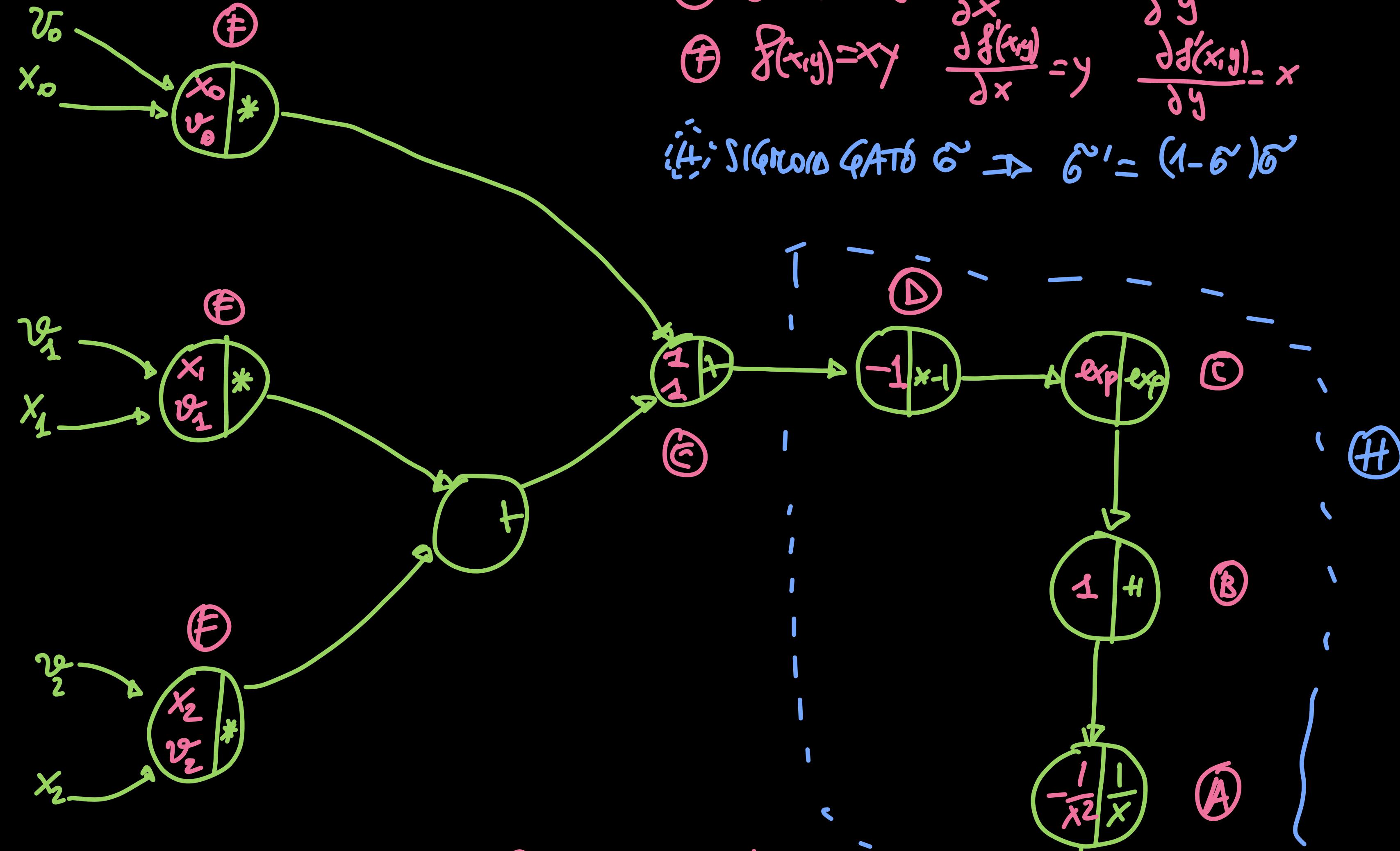
$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$\textcircled{C} \quad f(x) = e^x \Rightarrow f'(x) = e^x$$

$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1 \quad \textcircled{D} \quad f_a(x) = ax \quad f'_a(x) = a$$

$$h(v_i, x)$$

② INSERT D621 ATIVS



$$A) f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$B) f(x) = c + x \Rightarrow f'(x) = 1 \quad D) f_a(x) = ax \quad f'(x) = a$$

$$C) f(x) = e^x \Rightarrow f'(x) = e^x$$

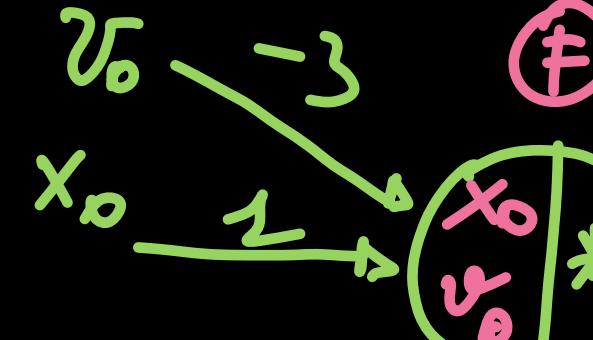
$$h(v_i x)$$

$$G) f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = 1 \quad \frac{\partial f(x,y)}{\partial y} = 1$$

$$H) f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = y \quad \frac{\partial f(x,y)}{\partial y} = x$$

$$I) \text{SIGMOID GATE } G \Rightarrow G' = (1-G)G'$$

③ ADD FORWARD VALUES



$$v = [-3, -3, -2]^T$$

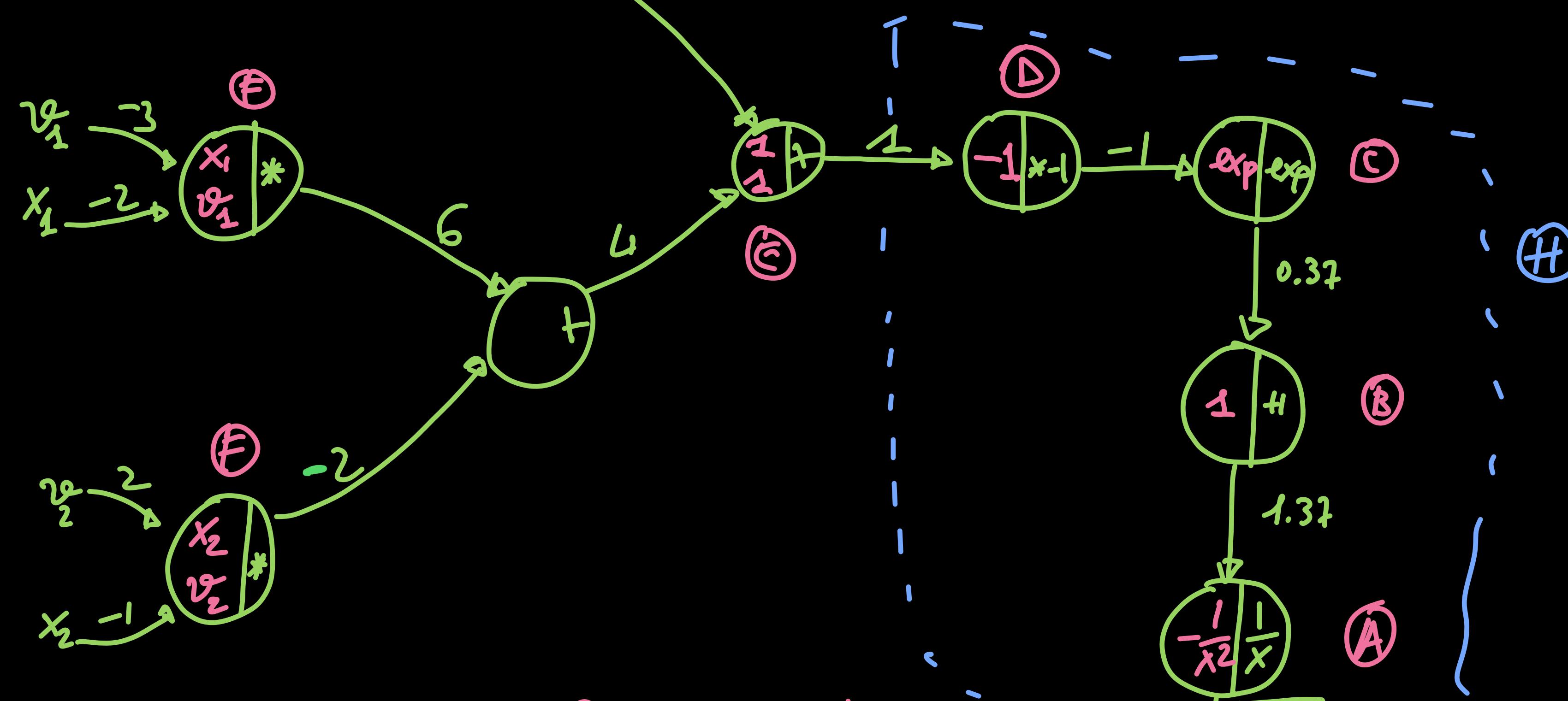
$$x = [1, -2, -1]^T$$

⑥ $f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = 1 \quad \frac{\partial f'(x,y)}{\partial y} = 1$

⑦ $f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = y \quad \frac{\partial f'(x,y)}{\partial y} = x$

⑧ SIGMOID GATE $G \Rightarrow G' = (1-G)G'$

FEED FORWARD AREA



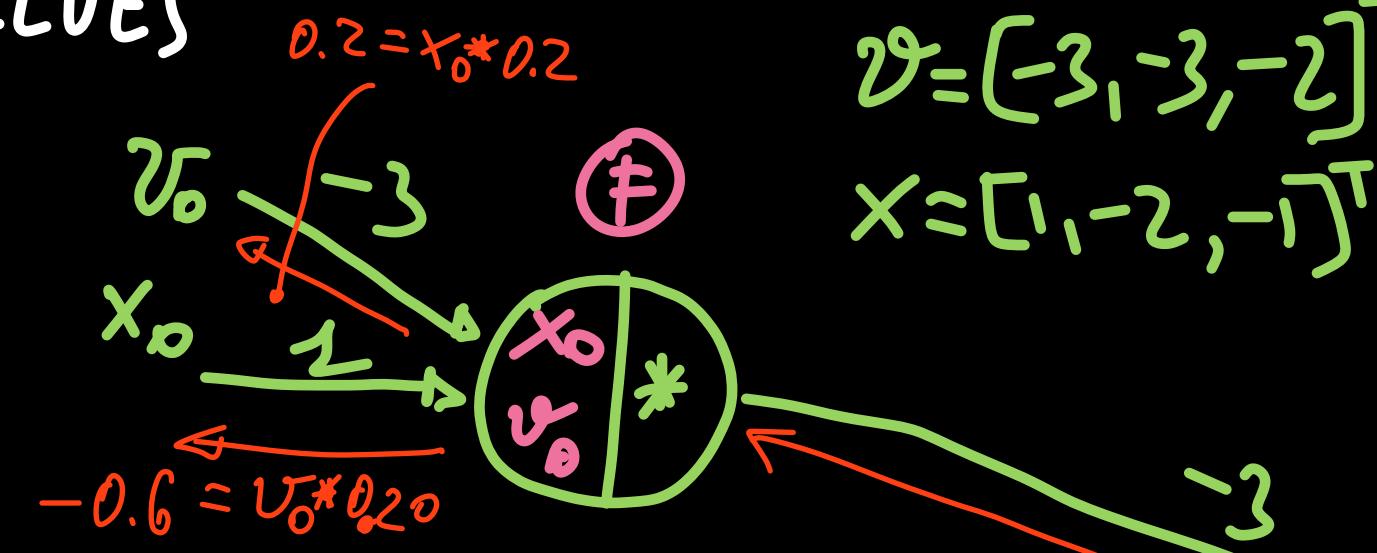
Ⓐ $f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$

Ⓑ $f(x) = c + x \Rightarrow f'(x) = 1$

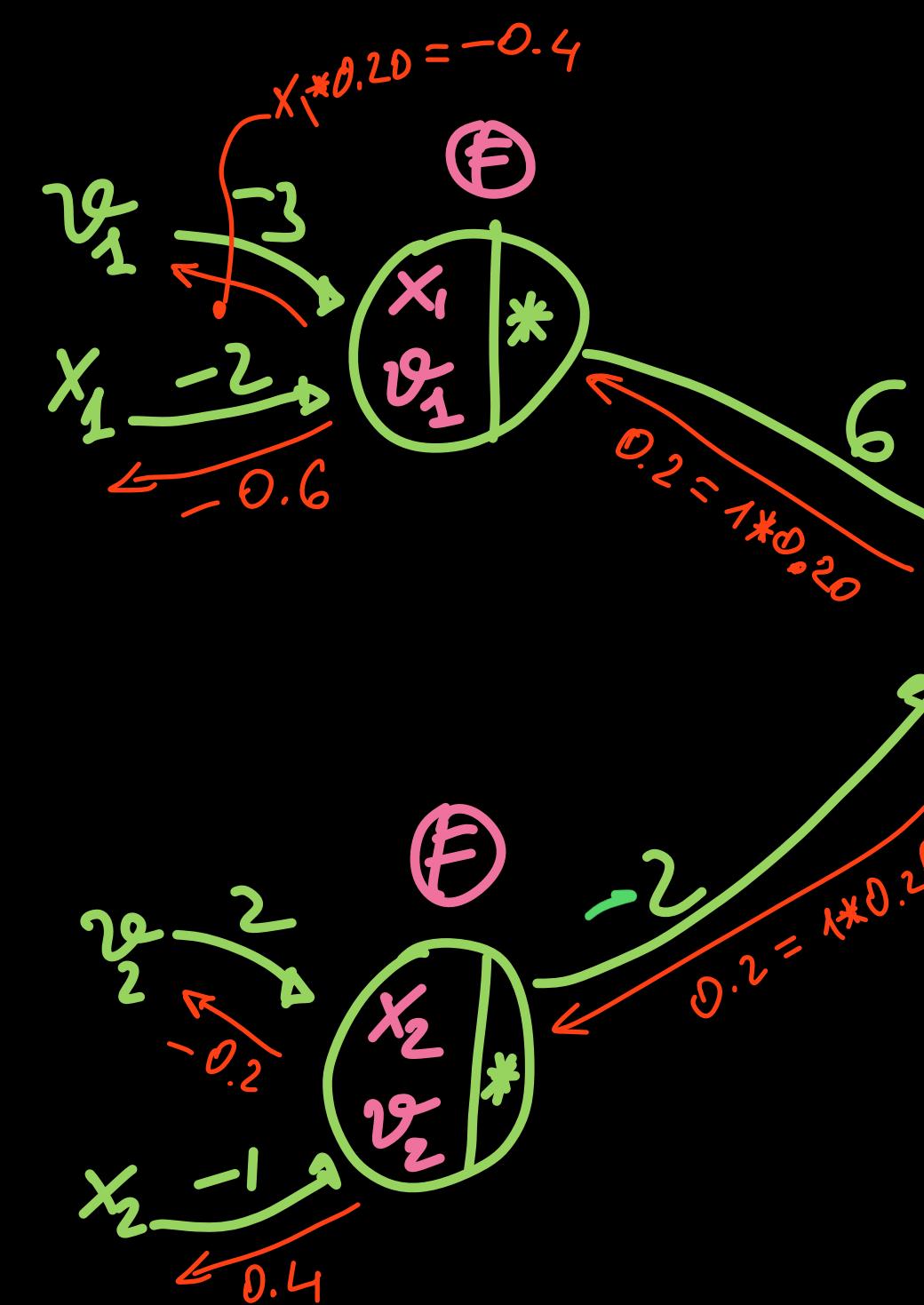
Ⓒ $f(x) = e^x \Rightarrow f'(x) = e^x$

Ⓓ $f_a(x) = ax \quad f'_a(x) = a$

④ ADD BACKWARD VALUES

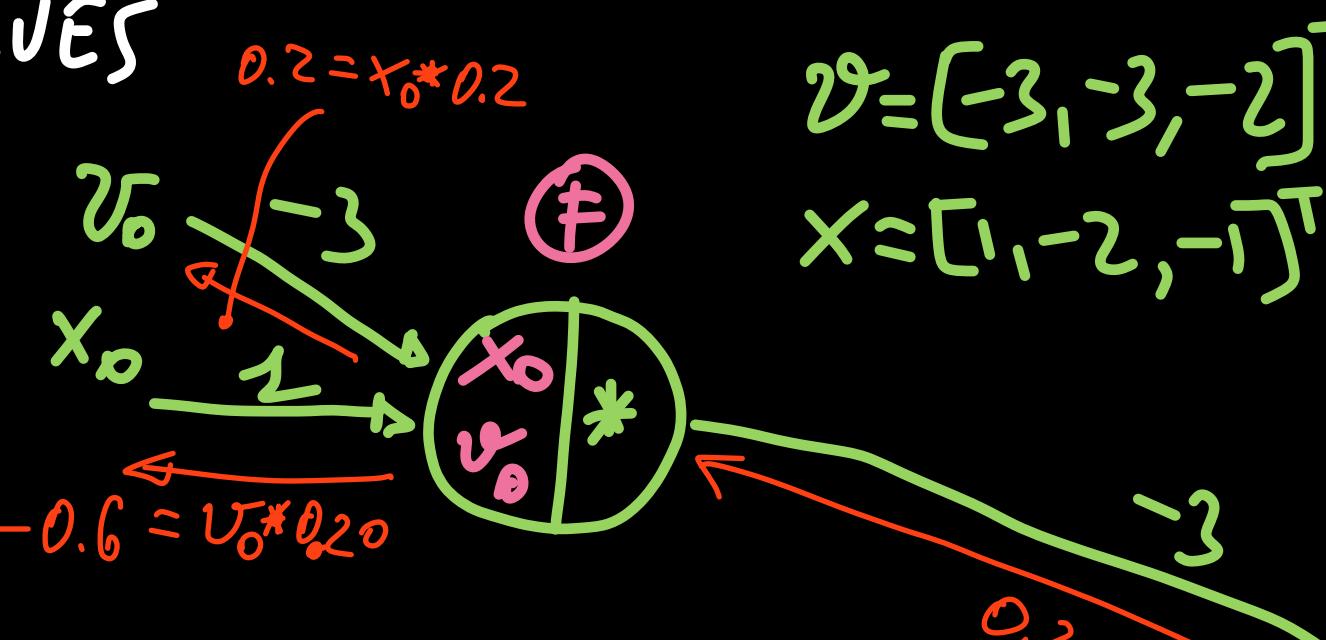


BACKWARD PASS



$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

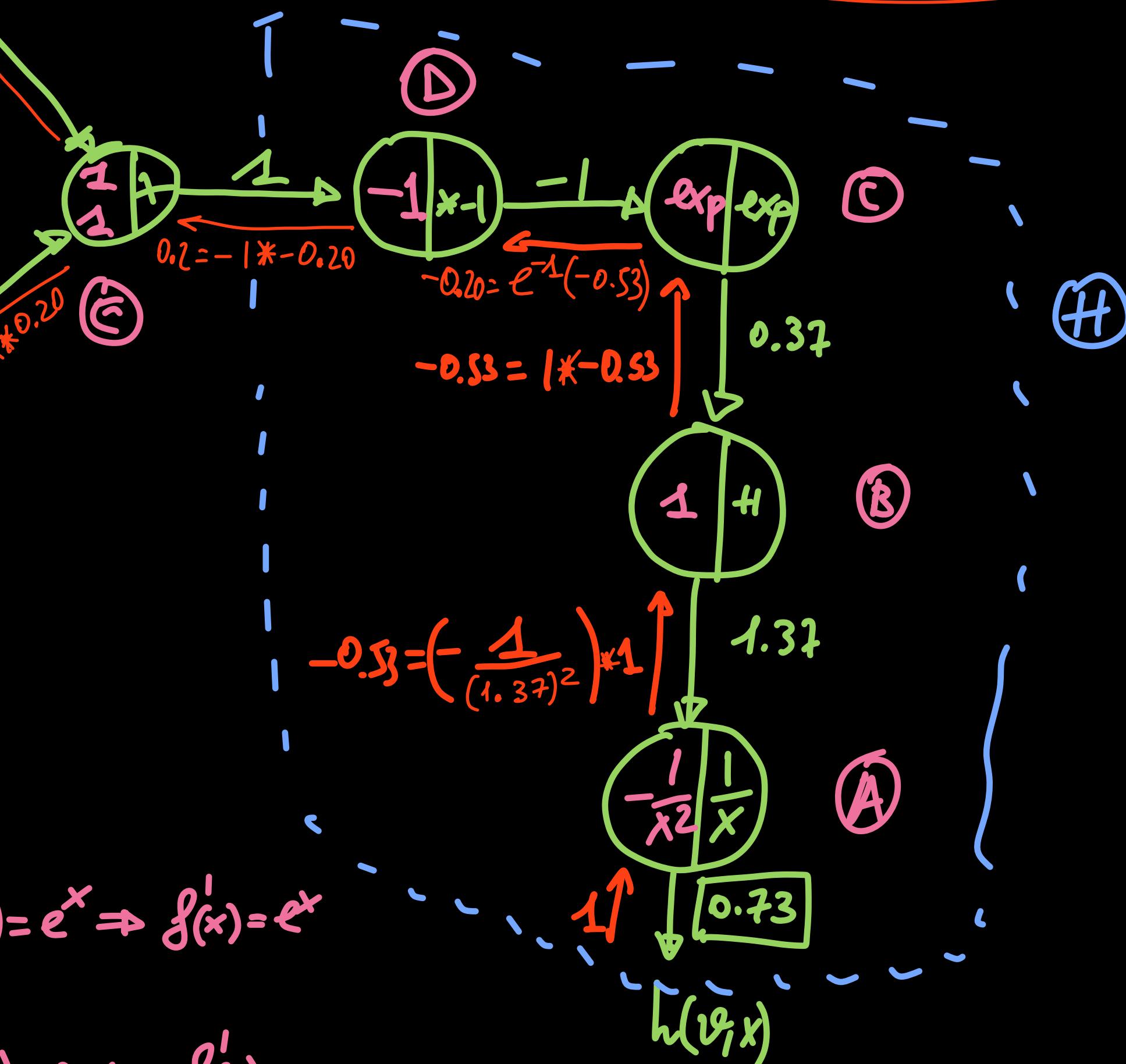
$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1 \quad \textcircled{D} \quad f_a(x) = ax \quad f'(x) = a$$



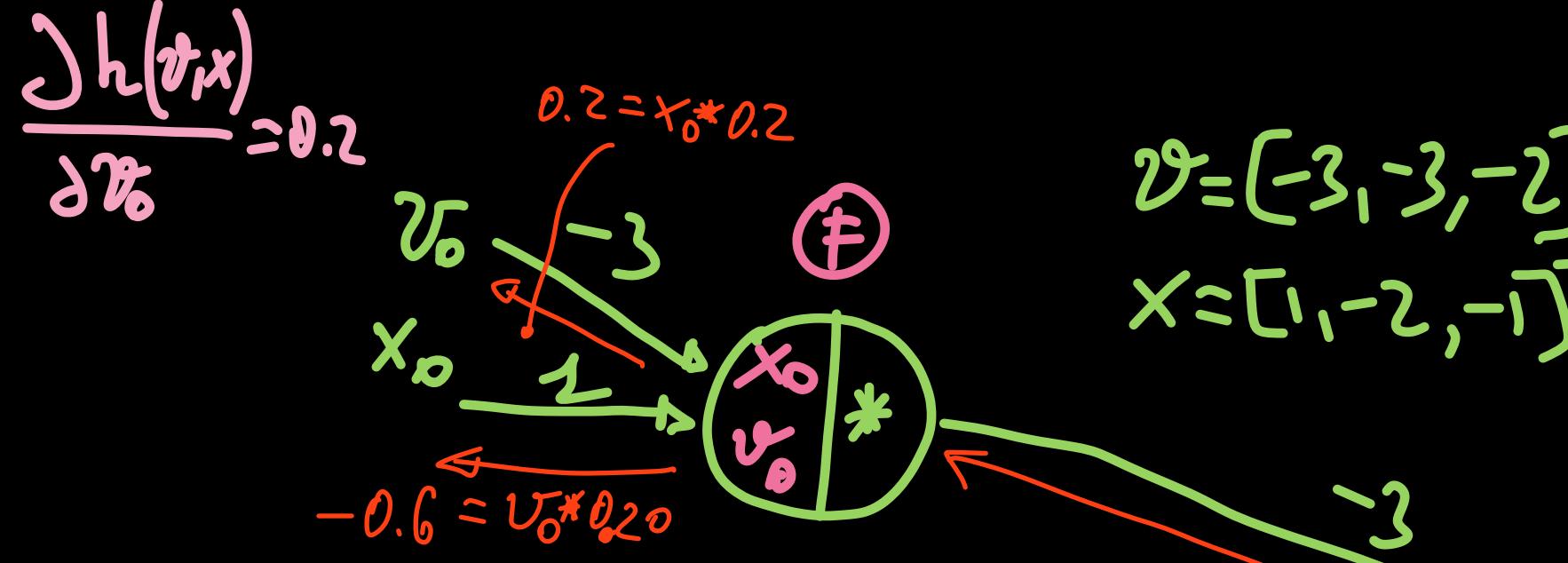
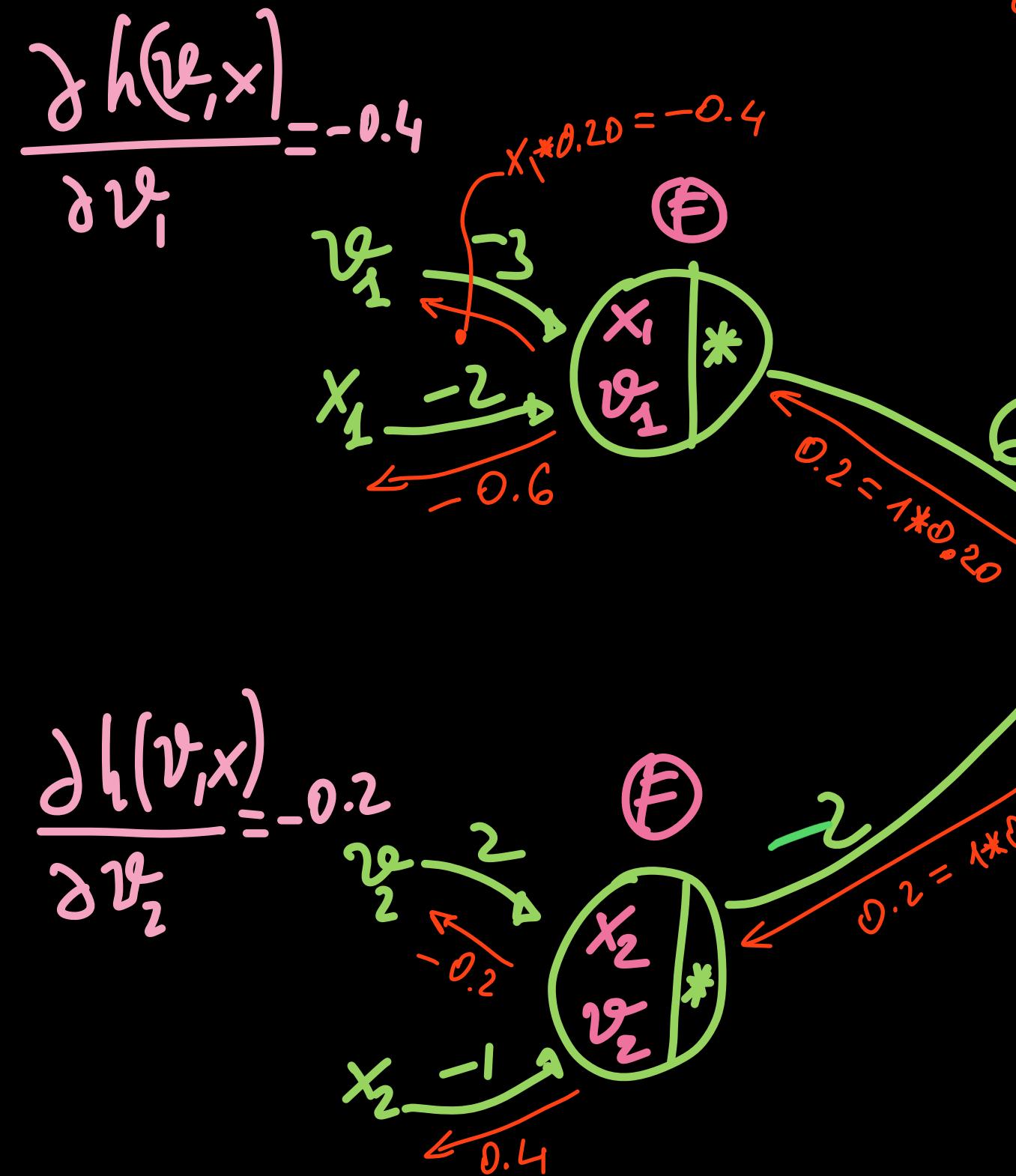
$$\textcircled{E} \quad f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = 1 \quad \frac{\partial f(x,y)}{\partial y} = 1$$

$$\textcircled{F} \quad f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = y \quad \frac{\partial f(x,y)}{\partial y} = x$$

$$\textcircled{H}: \text{SIGMOID GATE } 6' \Rightarrow 6' = (1 - 6')6' \quad \boxed{(1 - 0.73) \cdot 0.73 = 0.2}$$



⑤ Compute the final output of the network



$$v = [-3, -3, -2]^T$$

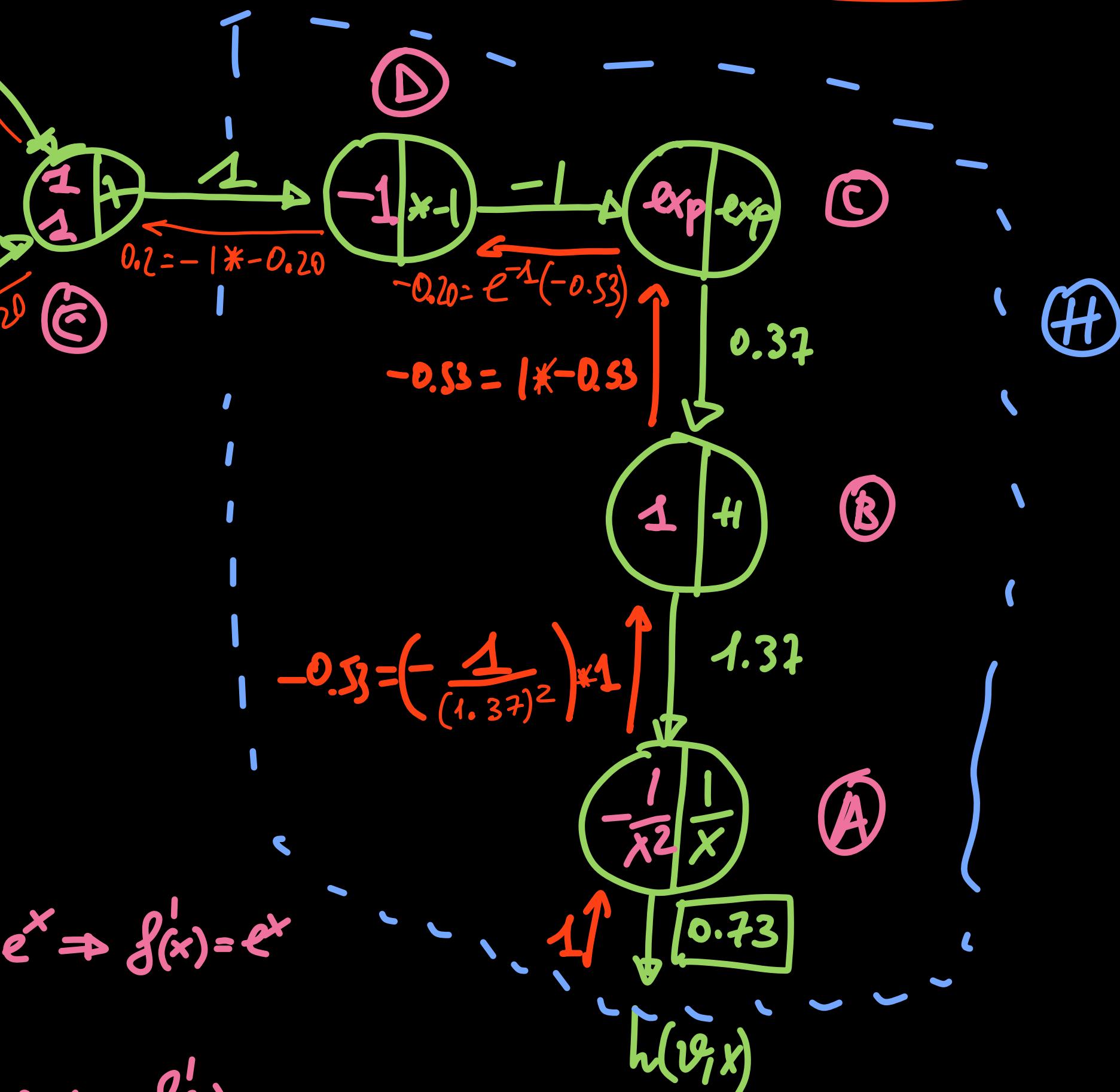
$$x = [1, -2, -1]^T$$

$$\textcircled{G} \quad f(x,y) = xy \quad \frac{\partial f(x,y)}{\partial x} = 1 \quad \frac{\partial f'(x,y)}{\partial y} = 1$$

$$\textcircled{F} \quad f(x,y) = x^y \quad \frac{\partial f(x,y)}{\partial x} = y \quad \frac{\partial f'(x,y)}{\partial y} = x$$

(4): sigmoid gate 6 $\Rightarrow \tilde{G}' = (1 - G')G'$

$$(1 - 0.73) \cdot 0.73 = 0.2$$



$$\textcircled{A} \quad f(x) = \frac{1}{x} \Rightarrow f'(x) = -\frac{1}{x^2}$$

$$\textcircled{B} \quad f(x) = c + x \Rightarrow f'(x) = 1$$

$$\textcircled{C} \quad f(x) = ax \Rightarrow f'(x) = a$$

HOMEWORK

Applicare l'algoritmo della backpropagation per derivare (utilizzando e disegnando grafo computazionale) la funzione riportata sotto rispetto ai parametri θ_i all'iterazione t -esima dell'algoritmo della discesa del gradiente considerando $x = [1,2,1]^T$ come input e $\theta = [1,1,1]^T$ al passo $t-1$:

$$f(x_0, x_1, x_2) = e^{-(\theta_0*x_0 + \theta_1*x_1^2 + \theta_2*x_2^3)}$$

Indicare quindi i valori di $\frac{\partial f(x_0, x_1, x_2)}{\partial \theta_i}$ per $i=0,1,2$. Calcolare quindi il nuovo vettore θ dei parametri una volta applicata la discesa del gradiente con learning rate $\alpha=2$.

Initializing Neural Networks

Further Readings

Go to the link for “**Tips and Tricks**” and further readings