

BACK PROPAGATION

- Computation Graphy
- why
 - good data structure for computing gradient
- Node → Variable (could be scalar vector, matrix or operations etc)
 - Variables can feed values to operations
- Edge → Funcⁿ arguments
data dependency
- Variables: Parameters of operation
Eg: $Wx + b$ x = Input
 w = Variable

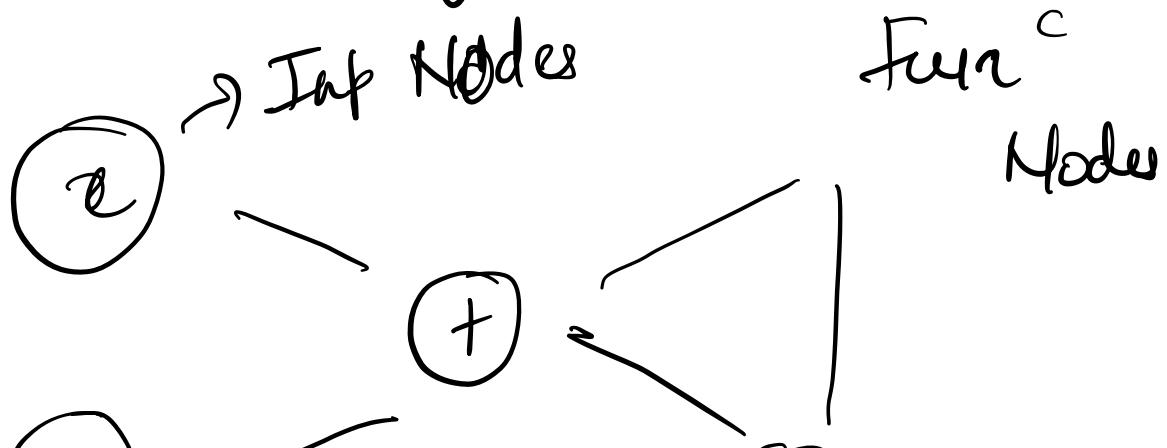
→ operations - Compute = Output given input

list of inputs = which can be nodes variable or other ops

list of consuming =

Eg: $f(x) = \overline{(x+y)^*}^2$

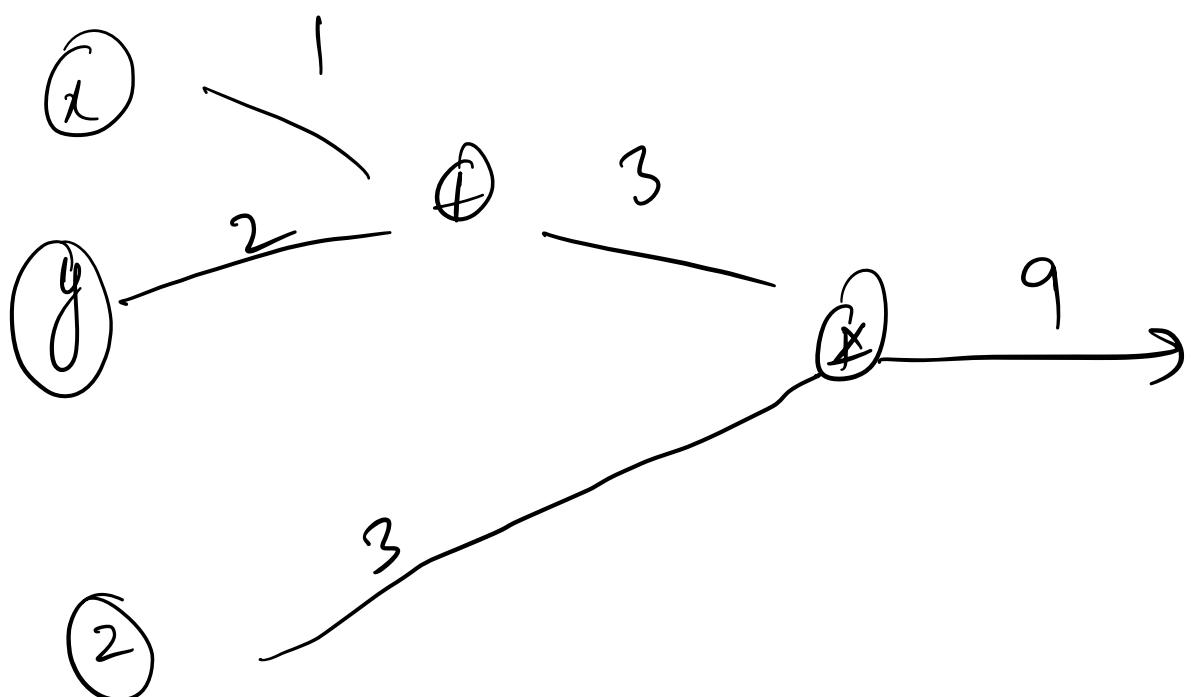
Eg: $(x+y)^*^2$



y

x

⇒ Values flow from edges,

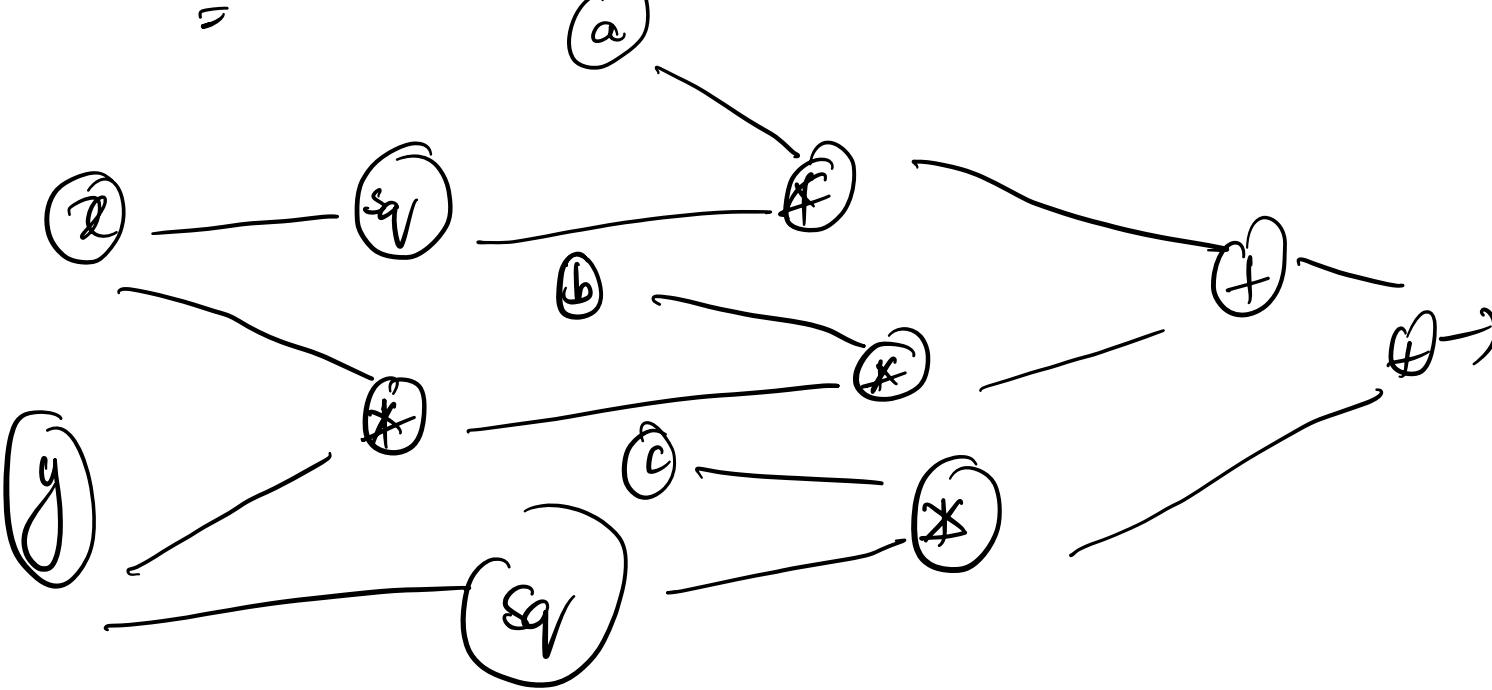


Neural Nets as Computation
graph

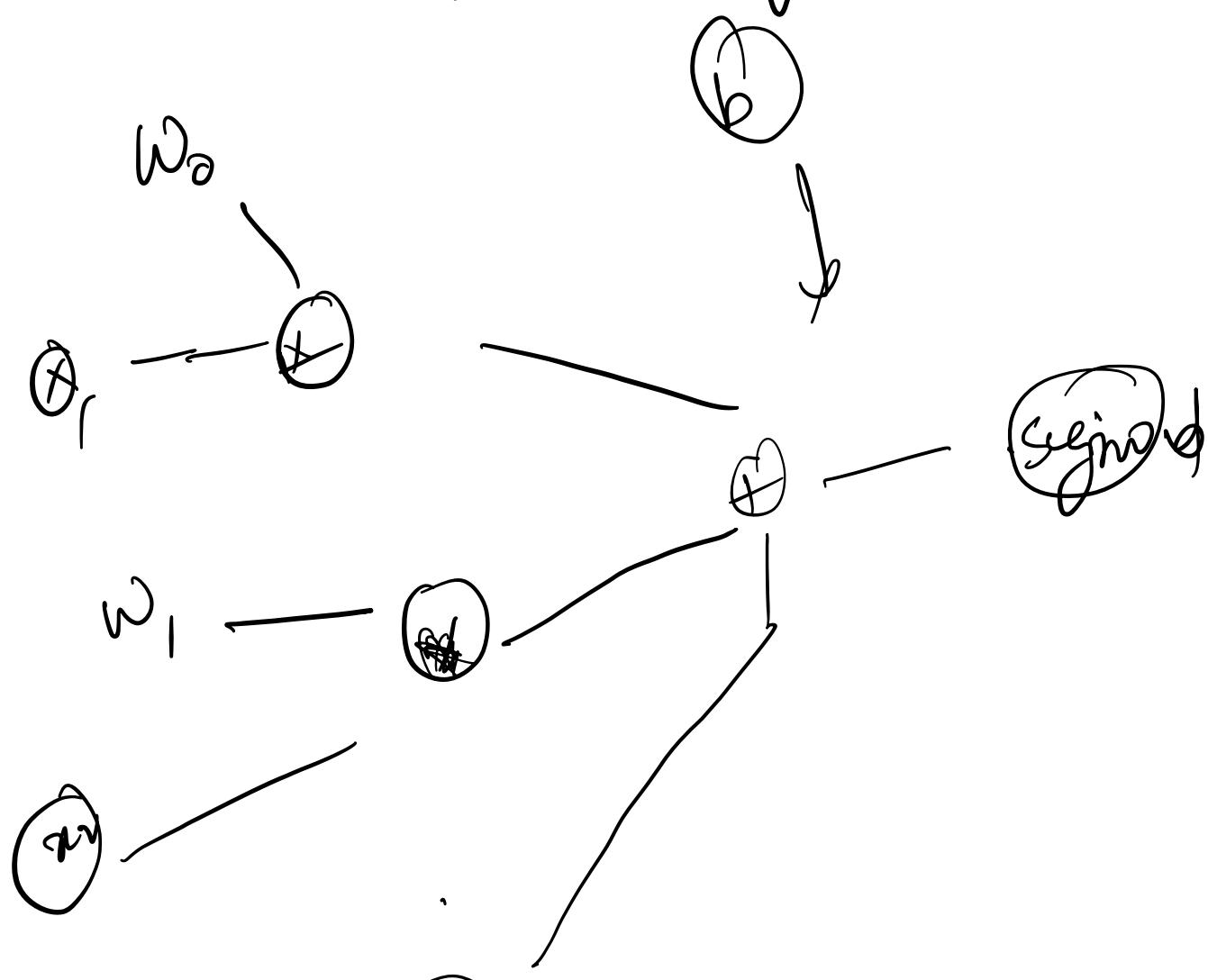
$$g: ax^2 + bxy + cx^2$$

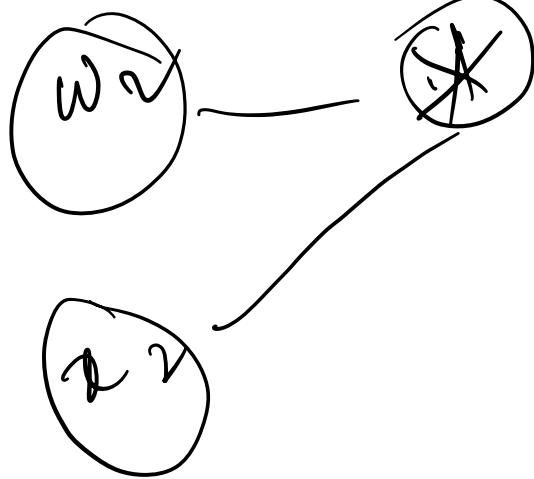
a, b, c = weights

x, y = inputs



Perceptron Graph





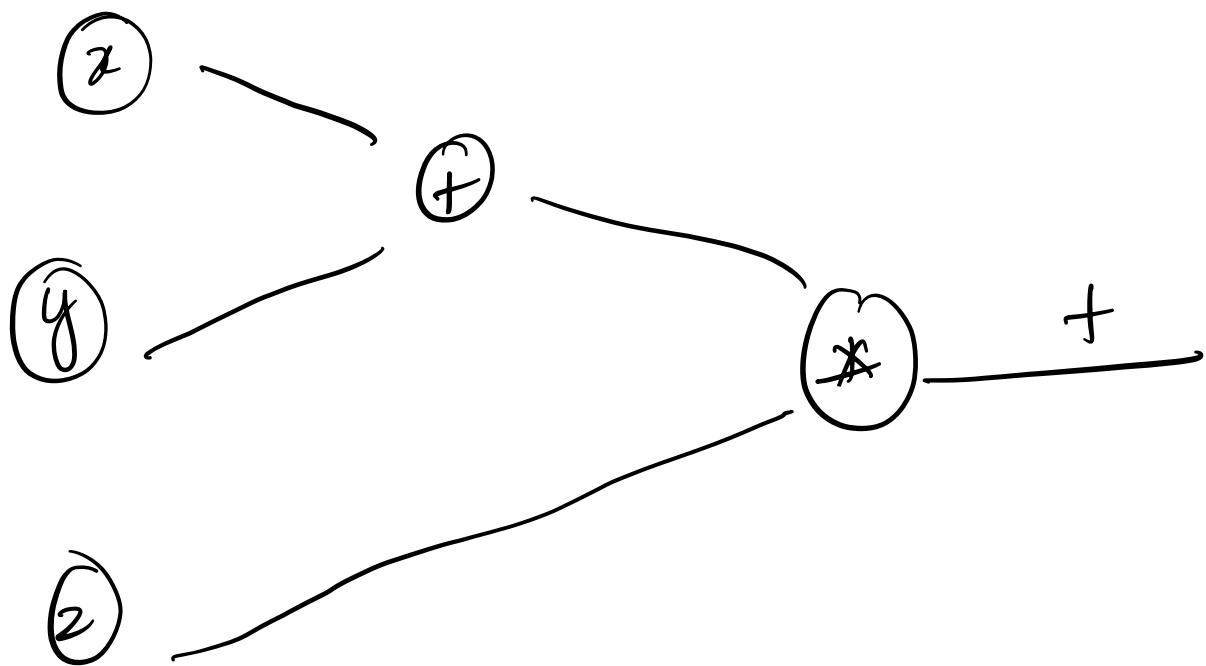
→ Why we need Con Graph?

- ↑ Deriving grads for deep net works everytime is complex one, relationship is easy with graphs.
- Data Structure aligns well with backprop and provides modularity.
- ↑ Graph traversal algorithms can help us in calculating grads easily

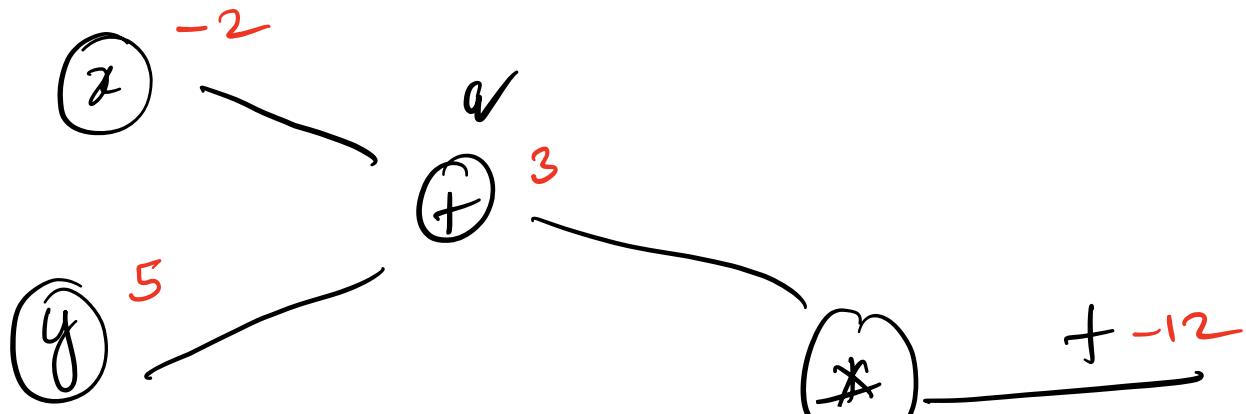
Backprop Simple Example

② Let's take

$$f(x, y, z) = (x + y)^2$$



① Forward pass



② 4

- ⇒ left to right
- ⇒ Compute output values.

② Backward pass

→ Want to compute grads wrt
each of input or output

i.e

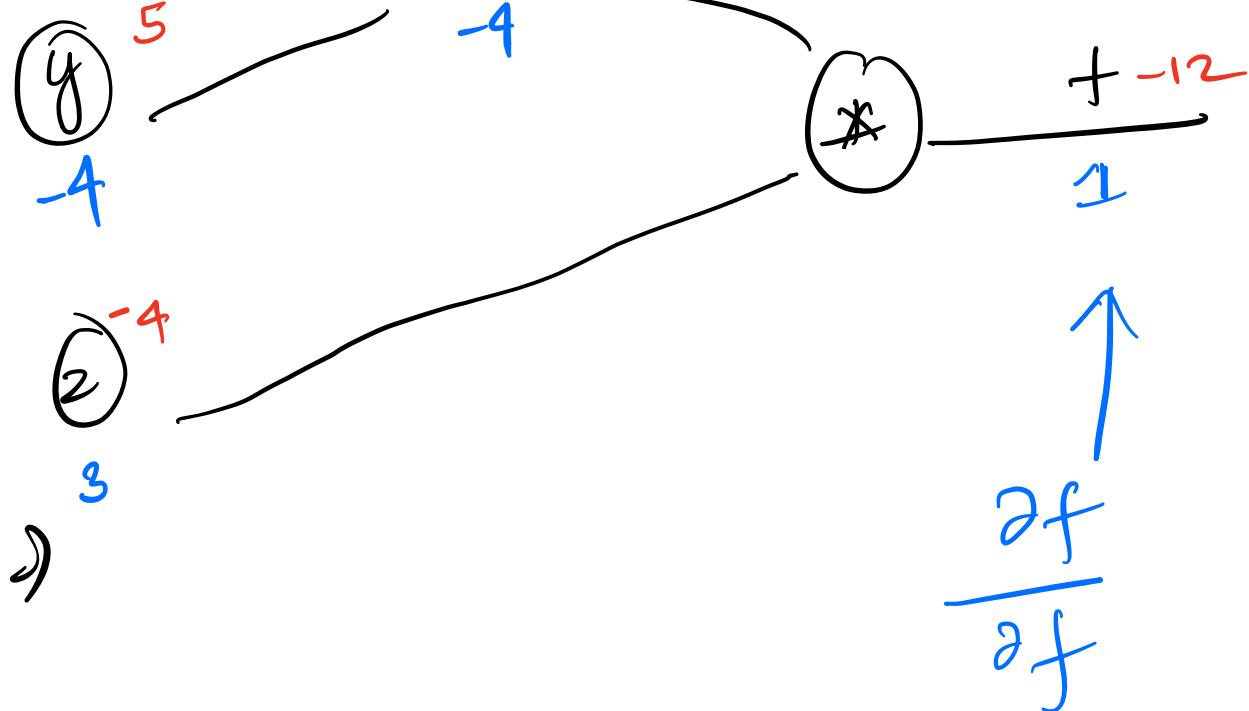
$$\frac{\partial f}{\partial x} \cdot \frac{\partial f}{\partial y} \frac{\partial f}{\partial z}$$

→ Right to left

Up → Values
Down → grads

② -2
4

a
+ 3



$$q = x + y ; \quad t = q^2$$

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

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

Now

$$\frac{\partial f}{\partial y} = \frac{\partial q}{\partial y} \times \frac{\partial f}{\partial q}$$

Since y is not directly connected

Downstream
gradient

↓
this is value
of
derivative
to
be
computed

Local
gradient

local
effect of
 y
affecting
intermediate
value

Upstream
gradient

flow much
output of
this piece
of
graph, affects
final output

$$q = x + y$$

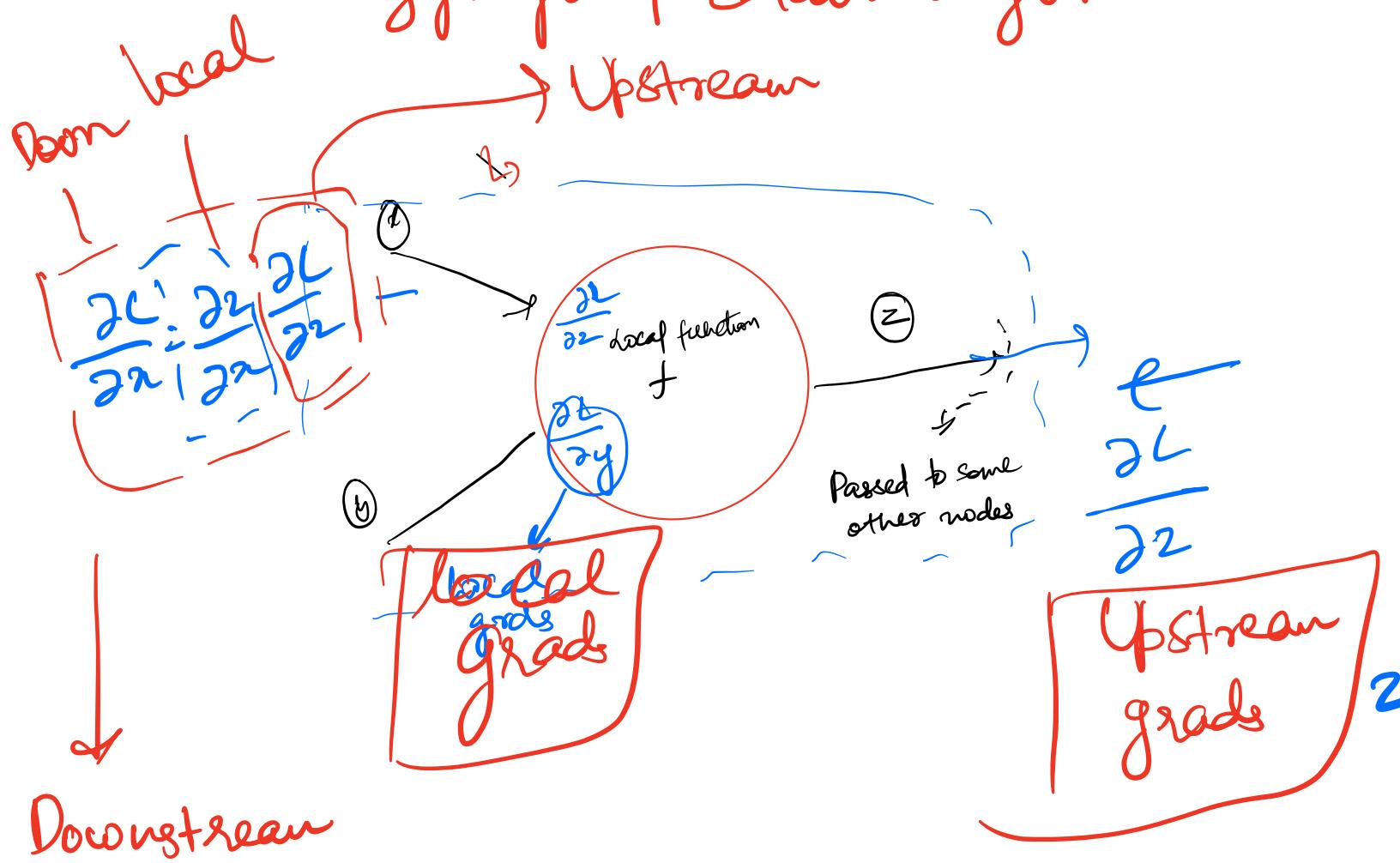
$$\frac{\partial f}{\partial q} \approx 1 \quad = \text{local derivative}$$

$$\frac{\partial f}{\partial y} \approx 1x - 4 = -4$$

- See so easy
- Just do forward
and
then do backward
- This way of thinking about
computation graph is very
modular

ZOOM IN ON NODE

- Each node doesn't need to know or care about rest of the graph.
- Just do local computation within each node, and by aggregating calculate global der.



grads

- Output passed to some other nodes
- Eventually loss is computed w.r.t this
- Eventually backprop passes some grad w.r.t this output to this node
- local grads: How much each output gets affected by each input of node
- Again downstream we passed along to other nodes in graph
- At the end we would have traversed graph entirely b

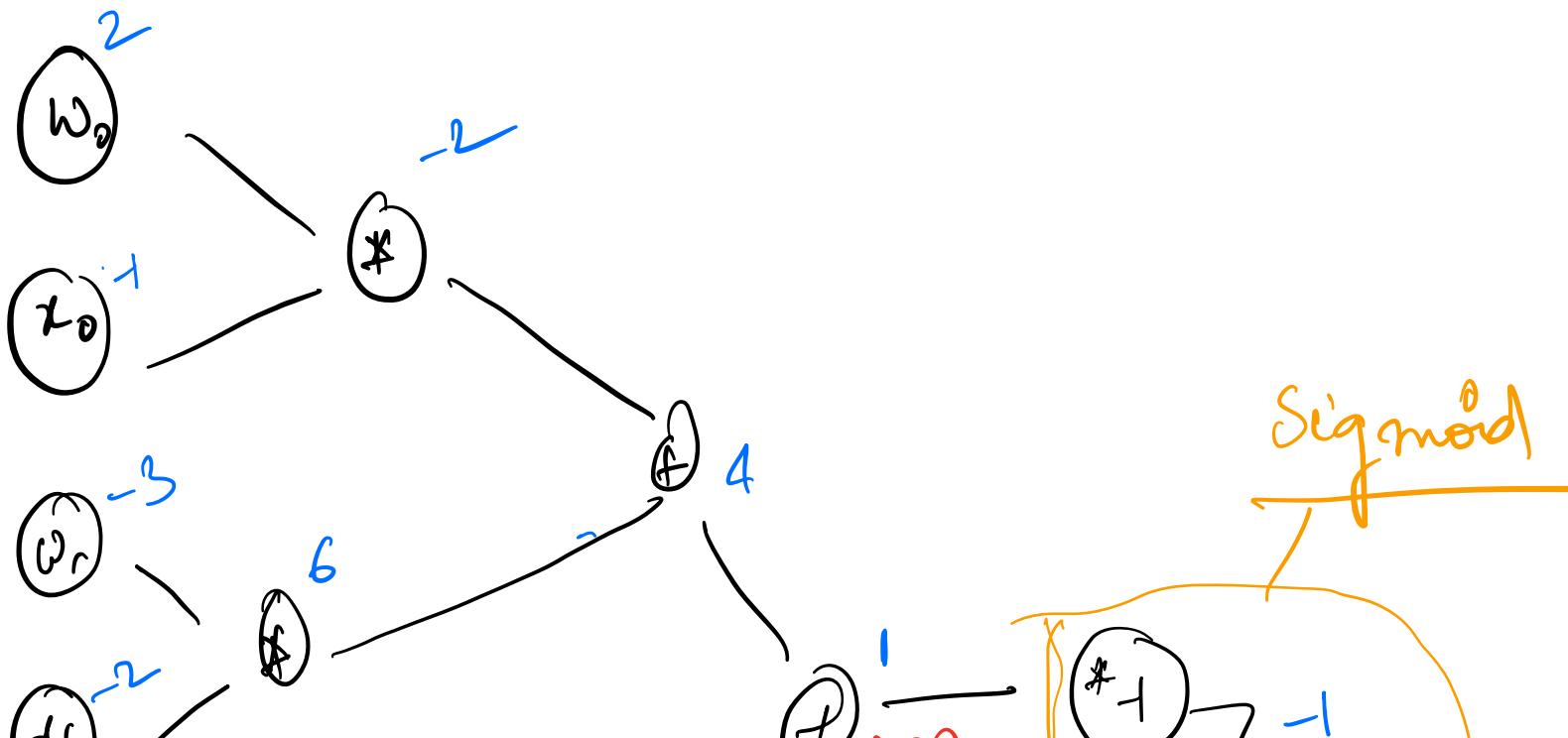
Cheerleaders competing the grads
entirely.

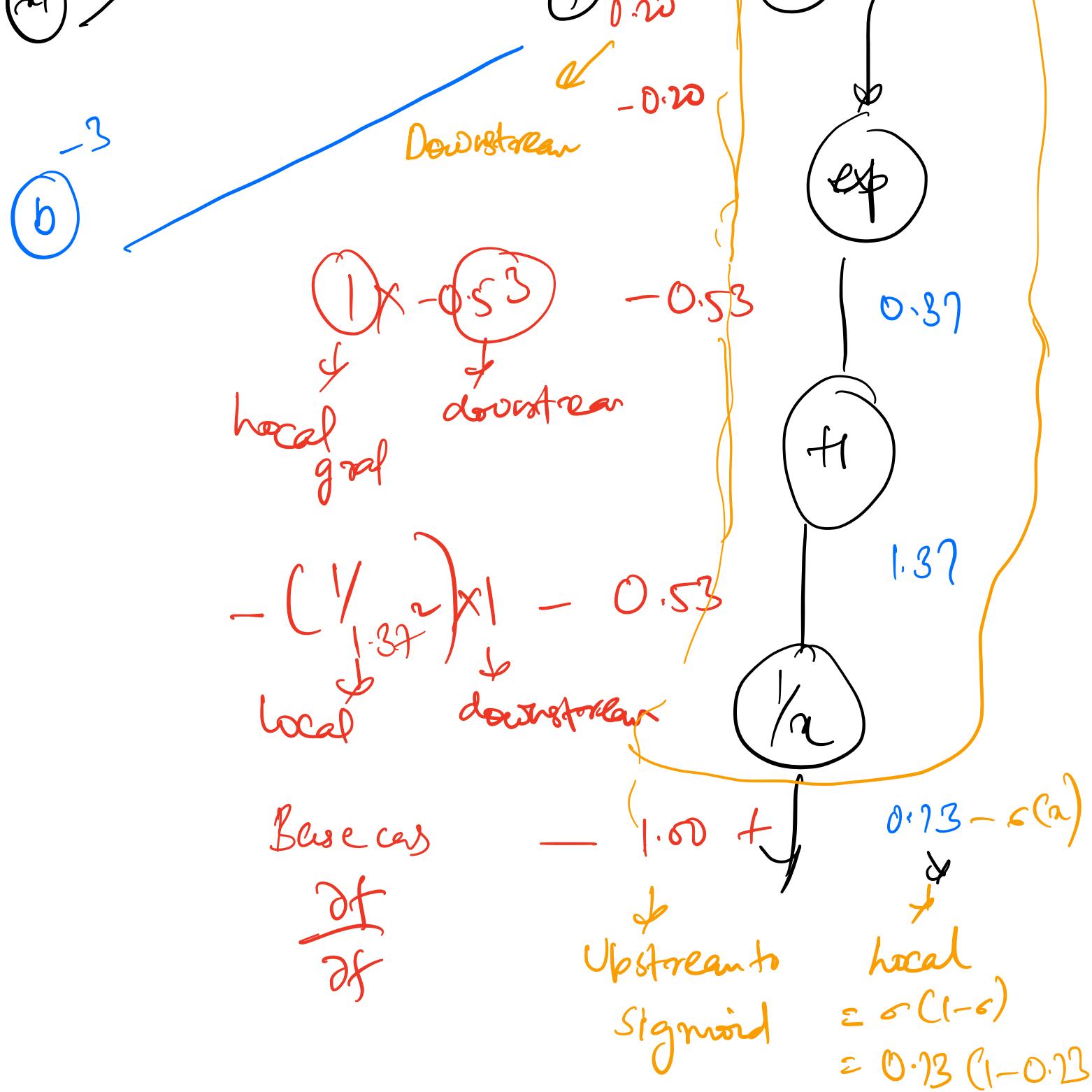
→ So GRAPH is good data structure

ANOTHER EXAMPLE

$$f(x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1)}}$$

↑
looks familiar ??.





\rightarrow There are multiple ways
 \rightarrow pick the pieces of conf
 use graph.

→ Type of primitives can be defined
on our own to be more
complicated

→ Eg: Sigmoid shows up most of
the time in ANN

→ So picking this as primitive
would be easy to have
as building block

local grad of σ

$$\frac{\partial f(x)}{\partial x} = (1 - \sigma(x)) (\sigma(x))$$

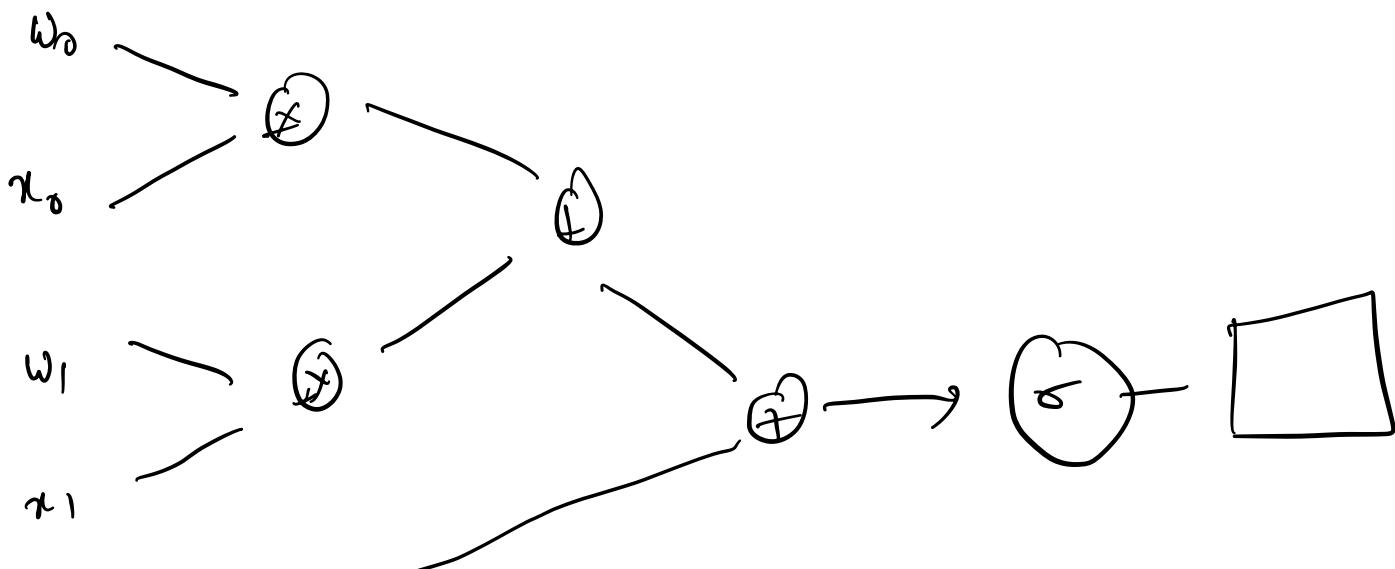
→

$$S_0 \\ 0.2 = (1 - 0.73) \times 0.73 \times 1 = 0.2$$

→ Defining more complex pointers
gives more efficiency and semantic
meaning

⇒ Implementation "flat" code

~~④~~



b

def $f(w_0, x_0, w_1, x_1, w_2)$

$$s_0 = w_0 * x_0$$

$$s_1 = w_1 * x_1$$

$$s_2 = s_0 + s_1$$

$$s_3 = s_2 + w_2$$

$$L = \text{sigmoid}(s_3)$$

Back

$$\text{grad}_L = 1.0$$

$$\text{grad}_s^3 = \text{grad}_L * (1-L) * L$$

$$\text{grad}_w^2 = \text{grad}_s^3$$

$$\text{grad}_c^2 = \text{grad}_s^3$$

$$\text{grad}_s^0 \approx \text{grad}_s^2$$

$$\text{grad}_s^1 = \text{grad}_s^2$$

2nd Way

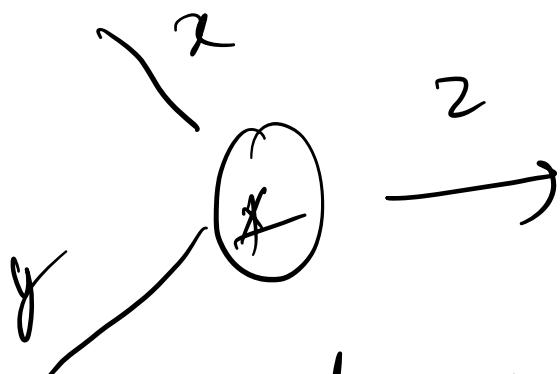
Modular API

class ComputeGraph

def forward () :

def backward () :

Eg:



class Multiply

def forward (x, y)

ctx.save(x, y) → Need the
in backprop
and

$$z = x \times y$$

return z

Upstream
grad
↓

def backward(ctx, grad_z)

x, y = ctx.get(x, y)

$$\text{grad } x = y \times \text{grad } z$$

$$\text{grad } y = x \times \text{grad } z$$

return grad_x, grad_y

Local
grads

→ What about Vectors, Sir??

→ Vector valued functions

→ Derivative is Jacobian

$$\frac{\partial y}{\partial z} \in \mathbb{R}^{N \times M}$$

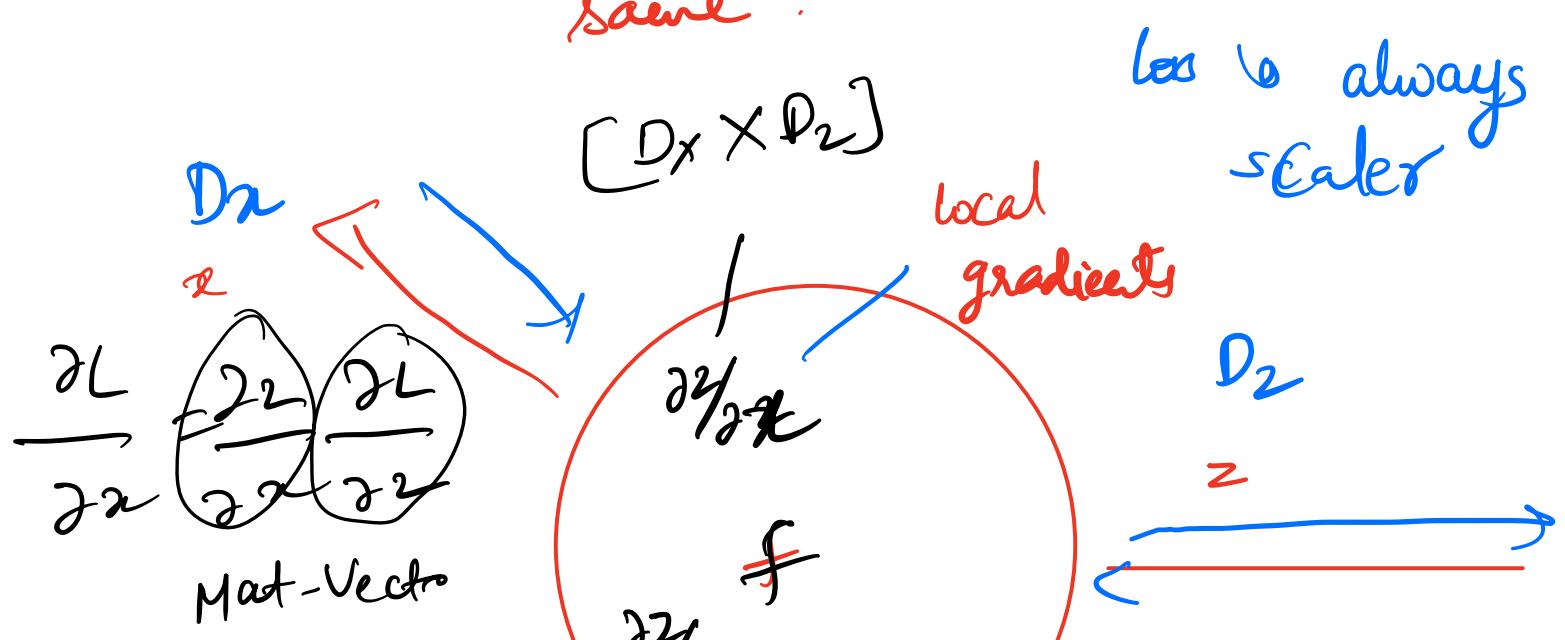
$$\left(\frac{\partial y}{\partial z} \right)_{n,m} = \frac{\partial y_m}{\partial z_n}$$

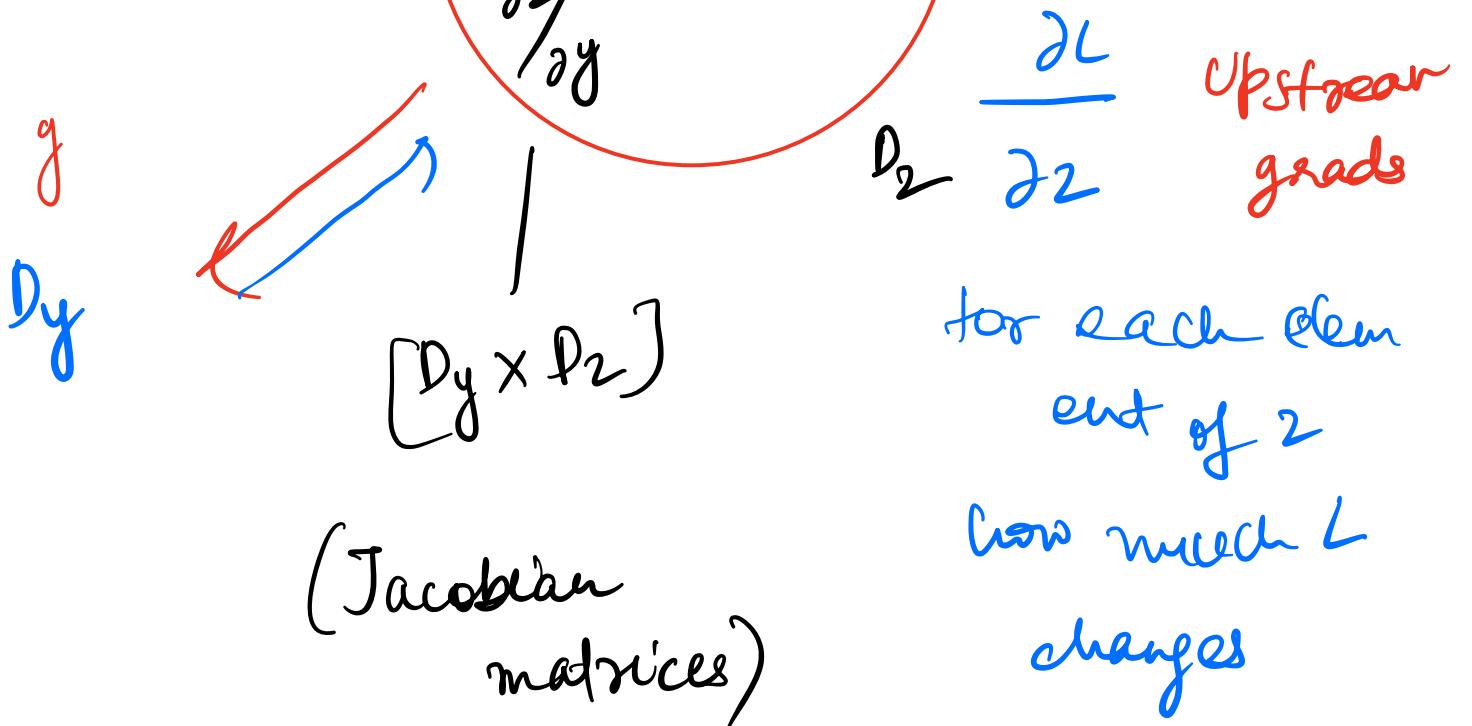
$$x \in \mathbb{R}^N$$
$$y \in \mathbb{R}^M$$

Backprop with Vectors ↗

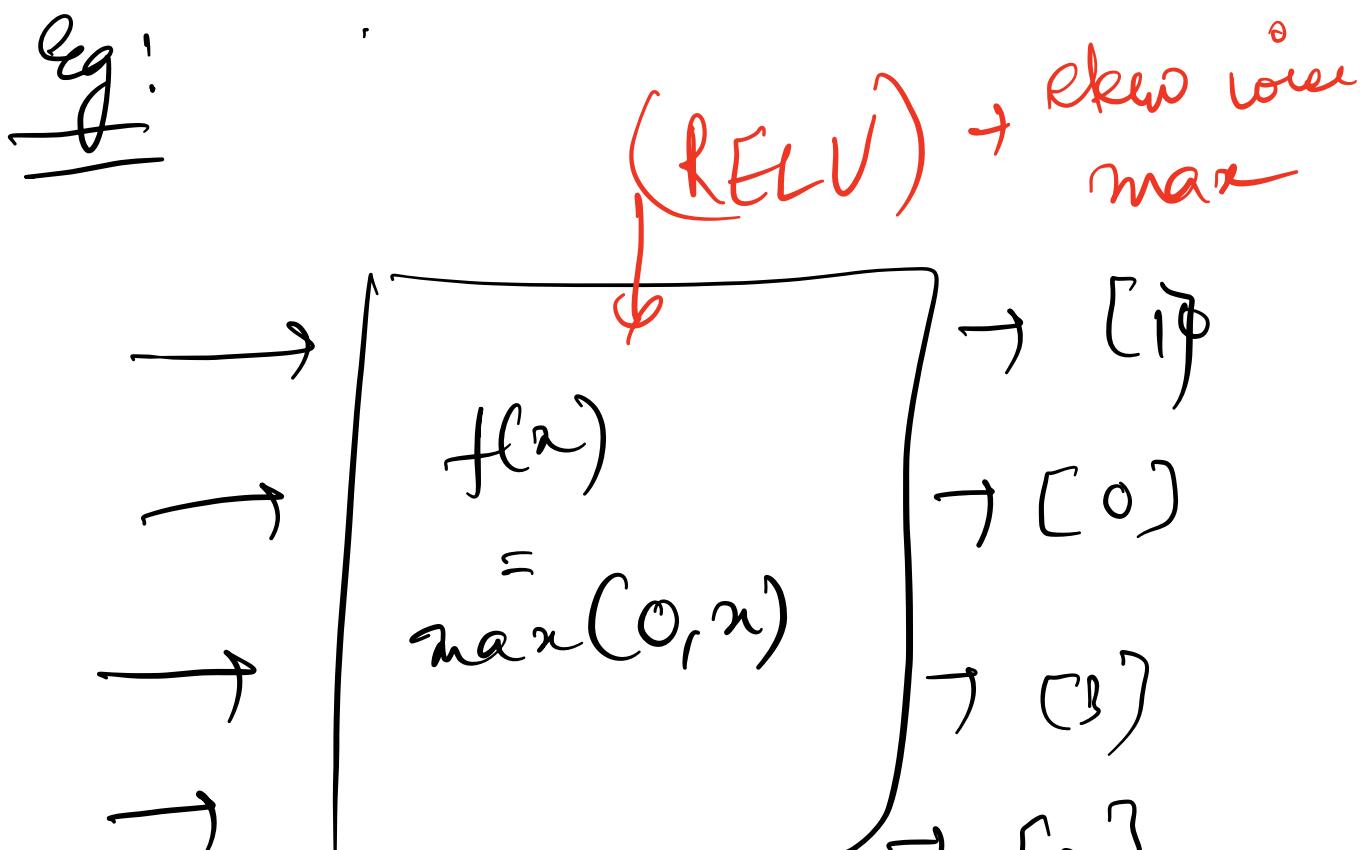
→ We know, we can just focus on one node

→ And then everything else is same.





for each output
elem, how much
change for each
input elem



$$\begin{array}{l}
 \text{Jacobian} \\
 \frac{\partial L}{\partial x} \left[\frac{\partial L}{\partial y} \right] \leftarrow [4]
 \\
 \begin{matrix}
 (4) & \leftarrow \\
 (0) & \leftarrow \\
 (5) & \leftarrow \\
 [0] & \leftarrow
 \end{matrix}
 \quad
 \begin{bmatrix}
 1 & 0 & 0 & 0 \\
 0 & 1 & 0 & 0 \\
 0 & 0 & 1 & 0 \\
 0 & 0 & 0 & 1
 \end{bmatrix}
 \begin{bmatrix}
 4 \\
 -1 \\
 1 \\
 9
 \end{bmatrix}
 \leftarrow [-1] \\
 \leftarrow [5] \\
 \leftarrow [9]
 \end{array}$$

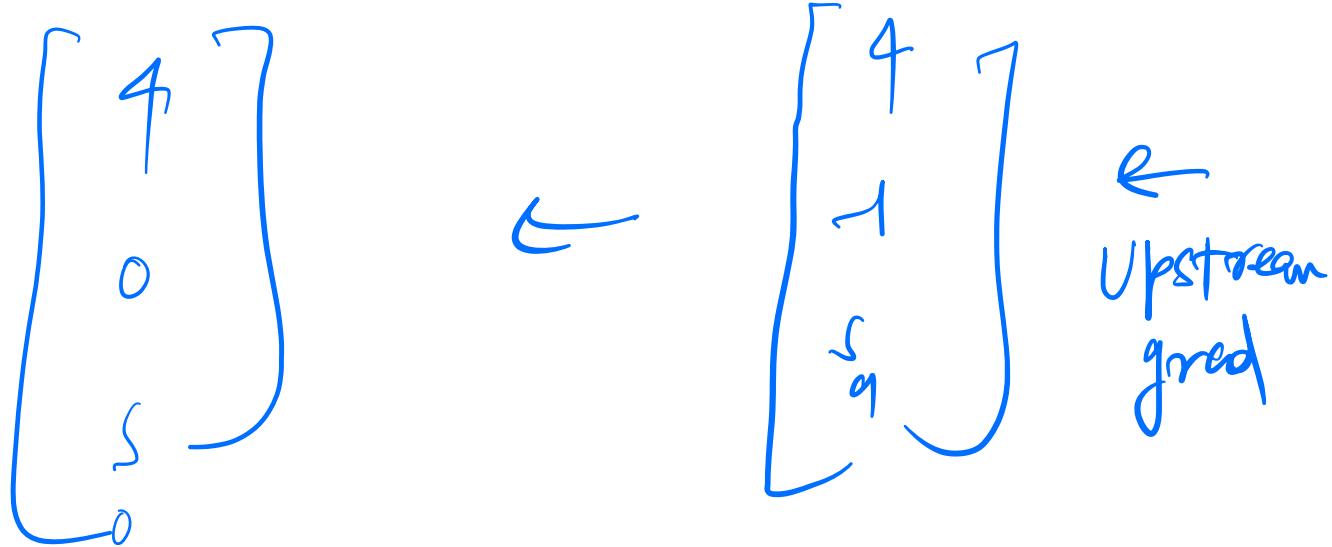
Upstream

\rightarrow very sparse
 usually
 Jacobian for
 most functions

\rightarrow Hence \rightarrow sparse
 trick here is
 to identify
 efficient way \rightarrow
 get
 same output

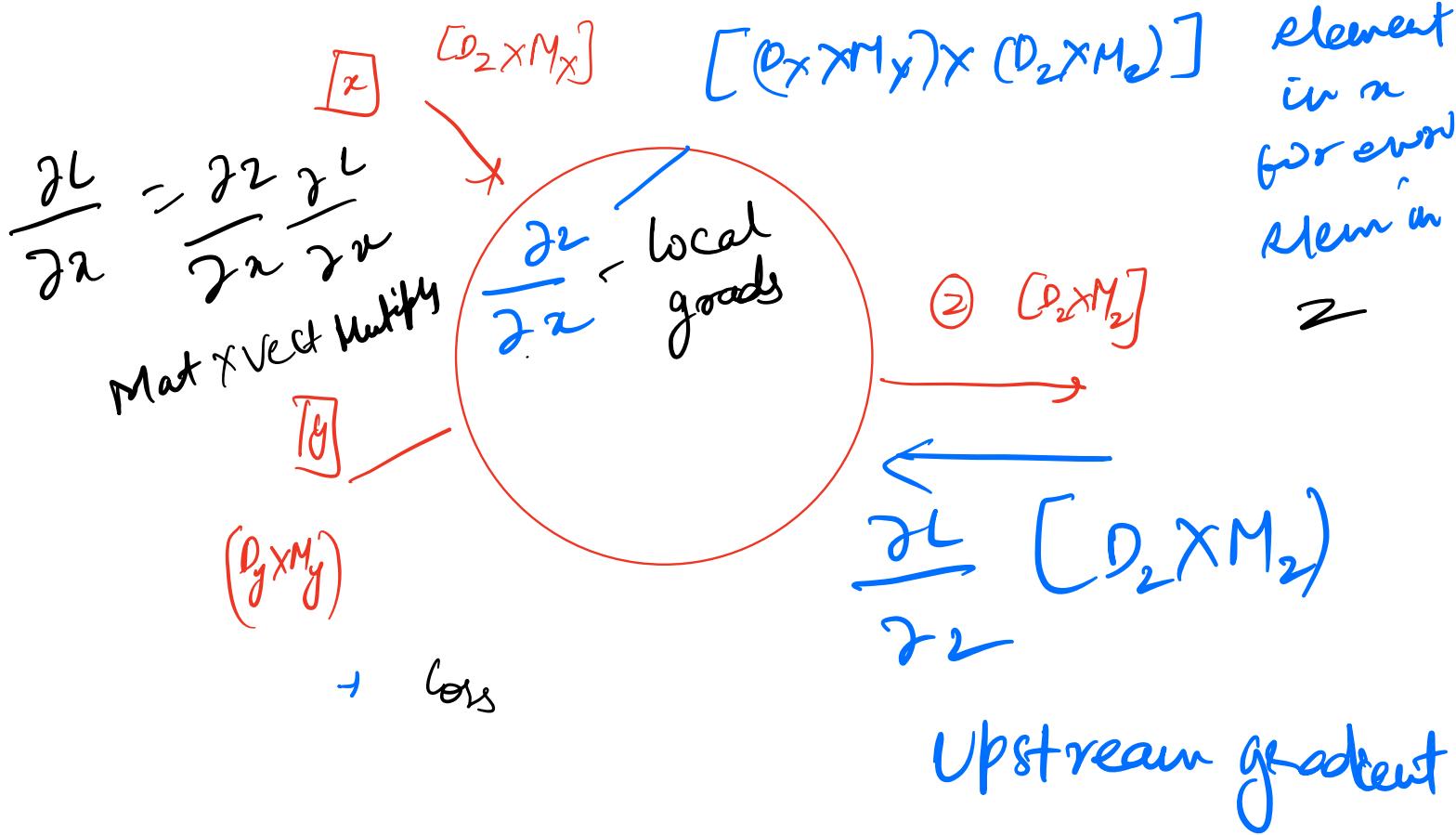
(6/1)

$$\left(\frac{\partial L}{\partial x} \right)_i = \begin{cases} \left(\frac{\partial f}{\partial x} \right)_i & x_i > 0 \\ 0 & \text{otherwise} \end{cases}$$



① backprop with Matrices

shape ie for each



still a scalar

$\rightarrow \frac{\partial L}{\partial x}$ will always have
same shape as x

\rightarrow easy to understand:

① They of in element
and element derivatives

lets take w example

$$x \\ [N \times D]$$

$$\begin{bmatrix} 2 & 1 & -3 \end{bmatrix}$$

$$\begin{bmatrix} -3 & 4 & 2 \end{bmatrix}$$

$$w \\ [D \times M]$$

$$\begin{bmatrix} 3 & 2 & 1 & -1 \end{bmatrix}$$

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

$$\begin{bmatrix} 3 & 2 & 1 & -2 \end{bmatrix}$$

Matrix mult

$$y_{n,m} = \sum_d x_{nd} w_{dm}$$

$$y \quad \begin{bmatrix} -1 & -1 & 2 & 6 \end{bmatrix}$$

$$\begin{bmatrix} 2 & 3 & -3 & 9 \end{bmatrix}$$

$$\begin{bmatrix} 5 & 2 & 11 & 7 \end{bmatrix} \quad \begin{bmatrix} -8 & 1 & 16 \end{bmatrix}$$

$\frac{\partial L}{\partial y}$

Upstream
grad

Jacobian Shape

$$\frac{\partial y}{\partial x} = [(N \times D) \times N \times M]$$

$$\frac{\partial y}{\partial w} = [(D \times M) \times (N \times M)]$$

Eg: $N = 64 \quad D = 4$,
 $M = 4096$

≈ 256 GB of mem

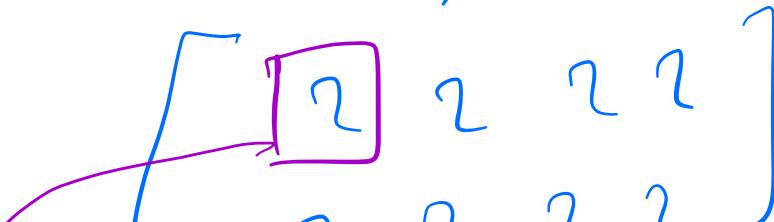
(Trick is to find a way to do this implicitly)

→ Let's look at just a slice of local grad

$$\frac{\partial L}{\partial x_{1,1}}$$

$$= \left(\frac{\partial y}{\partial x_{1,1}} \right) \left(\frac{\partial L}{\partial y} \right)$$

↓ local grad slice



$$\left(\begin{array}{c} y_1 \\ y_2 \\ y_3 \end{array} \right) = \left(\begin{array}{ccc} 1 & 2 & 3 \\ 2 & 1 & 0 \\ 3 & 0 & 1 \end{array} \right) \left(\begin{array}{c} x_1 \\ x_2 \\ x_3 \end{array} \right)$$

$$\rightarrow \frac{dy_{1,1}}{dx_{1,1}}$$

$$y_{1,1} = x_{1,1} w_{1,1} + x_{1,2} w_{2,1} + x_{1,3} w_{3,1}$$

$$\frac{dy_{1,1}}{dx_{1,1}} = \boxed{w_{1,1}}$$

So, we know what this is

\rightarrow So you see a pattern

$$\boxed{\frac{dy_{1,2}}{dx_{1,1}}}$$

Just copying

$$= [3 \ 2 \ 1 \ -1]$$

$$[0 \ 0 \ 0 \ 0]$$

91

local gradient slice



So now

$$\frac{\partial L}{\partial x_{1,1}}$$

$$= \frac{dy}{dx_{1,1}} \cdot \frac{dL}{dy}$$

$$= (w_1, \dots) \times (dL/dy)$$

$$= 3^{*2} + 2^{*3} + 1^{*}-3 + (-1)^{*9} = 0$$



→ Now keep on doing this

→ So think of it as one element at a time

$$\frac{\partial L}{\partial x} = \left[\begin{array}{c} \gamma_{xy} \end{array} \right] w^\top$$

$N \times D$ $[N \times M]$ $[M \times D]$

Another View



Chain rule $\frac{\partial L}{\partial w_0} = \left(\frac{\partial x_1}{\partial w_0} \right) \left(\frac{\partial x_2}{\partial x_1} \right) \left(\frac{\partial x_3}{\partial x_2} \right) \left(\frac{\partial L}{\partial x_3} \right)$

$$D_0 \times R_3 \quad D_0 \times D_1 \quad D_1 \times D_2 \quad D_2 \times D_3 \quad D_3$$

→ We do right to left,
so we always do vector
matrix multiply

→ This is called reverse mode
auto diff

→ Using same understanding we
can use to calculate high
order derivative as well