

# Neural networks - 2

## Structure

### ① Recap

(15 mins)

- MLP Recap
- Softmax Recap
- Categorical Cross-Entropy
- Training Neural networks

## Problems Addressed

→ Multi-class Classification

↳ Stack multiple LR units

→ How to interpret as probability?

→ Softmax soln

↳ Better than normal division

→ Matrix View of Multiple LRU's

$$X = \begin{bmatrix} x_1^1 & \dots & x_n^1 \\ x_1^2 & \dots & x_n^2 \\ \vdots & & \vdots \\ x_1^m & \dots & x_n^m \end{bmatrix} \quad W = \begin{bmatrix} w_1^1 & w_1^2 & w_1^c \\ w_2^1 & w_2^2 & \dots \\ \vdots & \vdots & \vdots \\ w_n^1 & w_n^2 & w_n^c \end{bmatrix}$$

$m \times n$

$n \times c$

$b$

$m \times c$

$$b = \begin{bmatrix} b_1 & b_2 & b_3 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix} \quad \text{Broadcasting}$$

→ Entropy

→ Measure of Surprise

- → Cross Entropy Loss

$$H(P, Q) \nrightarrow \sum P(x) \log Q(x)$$

$H(P, Q) \nrightarrow$  measures how well

one prob distribution  $Q$  (predicted)  
approximates another distribution

P (true distribution)

"Expected no. of bits needed

to encode events from true

distribution  $P$ , using optimal

code for predicted distribution"

→ If  $Q$  matches  $P$ , then cross  
entropy reduces to entropy  
of  $P$ .

→ If  $Q$  diverges from  $P$ , then  
 $H(P, Q) \uparrow$ , because more bits  
are needed to encode events  
from  $P$  using code of range by  $Q$

## Backpropagation (60 mins)

Problem?

→ We want to minimize  
our loss,

= So we have a function

How do we optimize a  
function?

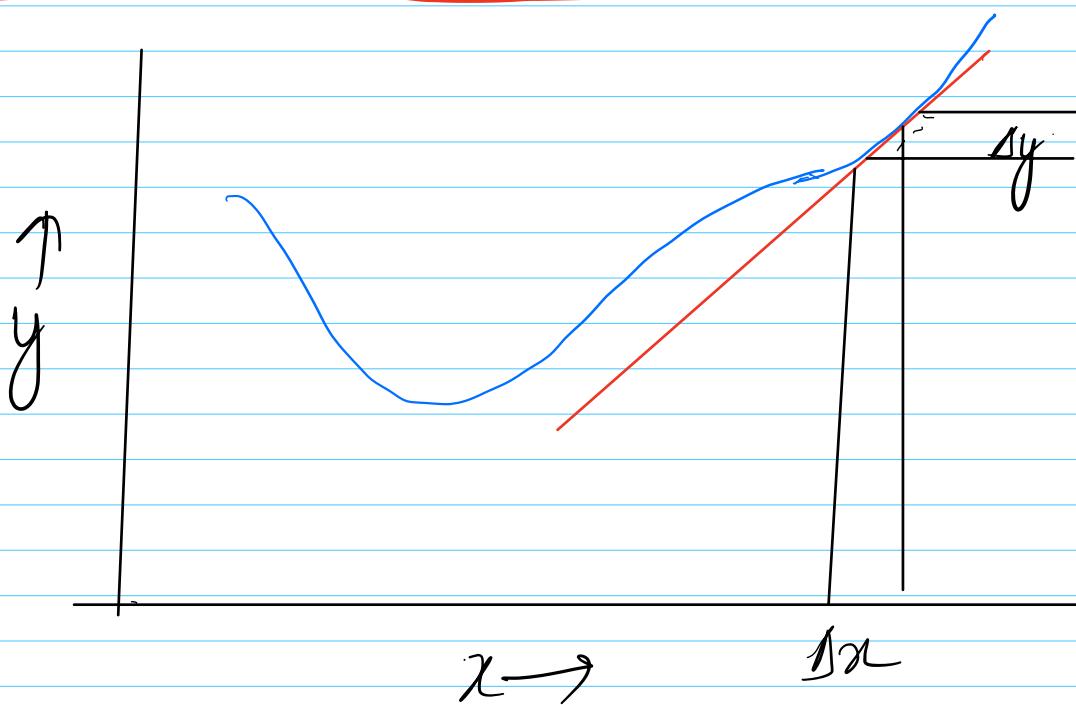
$$\text{Loss}(w) = \frac{1}{n} \sum \text{loss}(x_i, w, y_i)$$

T

$$\hat{w}_{\text{opt}} = \underset{w}{\operatorname{argmin}} \text{Loss}(w)$$

→ minimize error over training samples and hope to work ☺  
Maybe to next come back

## Quick Intro to Derivatives



→ wts derivative

$$\Delta y \approx \alpha \Delta x$$

(At fine enough resolution,  
Any smooth func<sup>n</sup> is locally  
(linear))

Cases

①  $x$  &  $y$  Scalar

$$y = f(x)$$

$$\Delta y \approx \alpha \Delta x$$

②  $y$  scalar,  $x$  vector

$y = \text{Scalar value}$

$x = \text{vector}$

$$\Delta x = \Delta$$

$$\Delta x = [\Delta x_1, \Delta x_2, \dots, \Delta x_B]$$

$$\Delta y = \alpha \Delta x$$

Partials

where

$$\alpha = \nabla_x y = \left[ \frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_n} \right]$$

→ This is derivative

→ Gradient is just transpose  
of above

→ Here derivative  $\nabla_x f(x)$  is

a scalar function  $f(x)$

of multi-variate input  $X$

$$\boxed{\nabla_x f(x)^T = \text{Gradient}}$$

$$\nabla_x f(x) = \left[ \frac{\partial f(x)}{\partial x_1}, \frac{\partial f(x)}{\partial x_2}, \dots, \frac{\partial f(x)}{\partial x_n} \right]$$

$$\boxed{(dfx) = \nabla_x f(x) \cdot dx}$$

We want this to be maximized  $\rightarrow$  Vector inner product

When is this inner product maximized?

→ When both vectors align

i.e  $d\mathbf{x}$  and  $\nabla_{\mathbf{x}} f(\mathbf{x})^T$  are in same direction

i.e  $d\mathbf{x}$  is exactly in dir<sup>n</sup> of gradient

Hence we take a step in direction of gradient

→ Why we need this ??

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$$

and solve for  $\mathbf{x}$

→ Becomes complex, to solve

$$f(x_1, x_2, x_3)$$

$$= x_1^2 + x_1(1-x_2) + x_2^2 - x_2x_3 \\ + x_3^2 + x_3$$

Grad

$$\nabla_{\mathbf{x}} f^T = \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 + 1 \end{bmatrix}$$

$$\nabla_{\mathbf{x}} f = 0 \Rightarrow \begin{bmatrix} 2x_1 + 1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 + 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

3 eqns, 3 unknowns

$$A\mathbf{x} = \mathbf{B} \Rightarrow A^{-1}A\mathbf{x} = A^{-1}\mathbf{B}$$

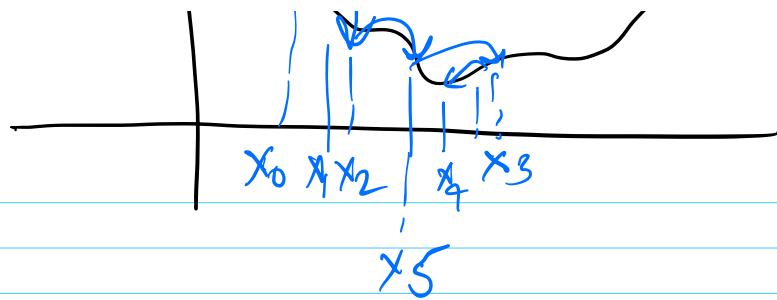
NN's have millions of params, doing matrix inversion is computationally expensive at such dimension, so we need to find other way

What can we use then?

→ Iterative Solutions

① Start from guess



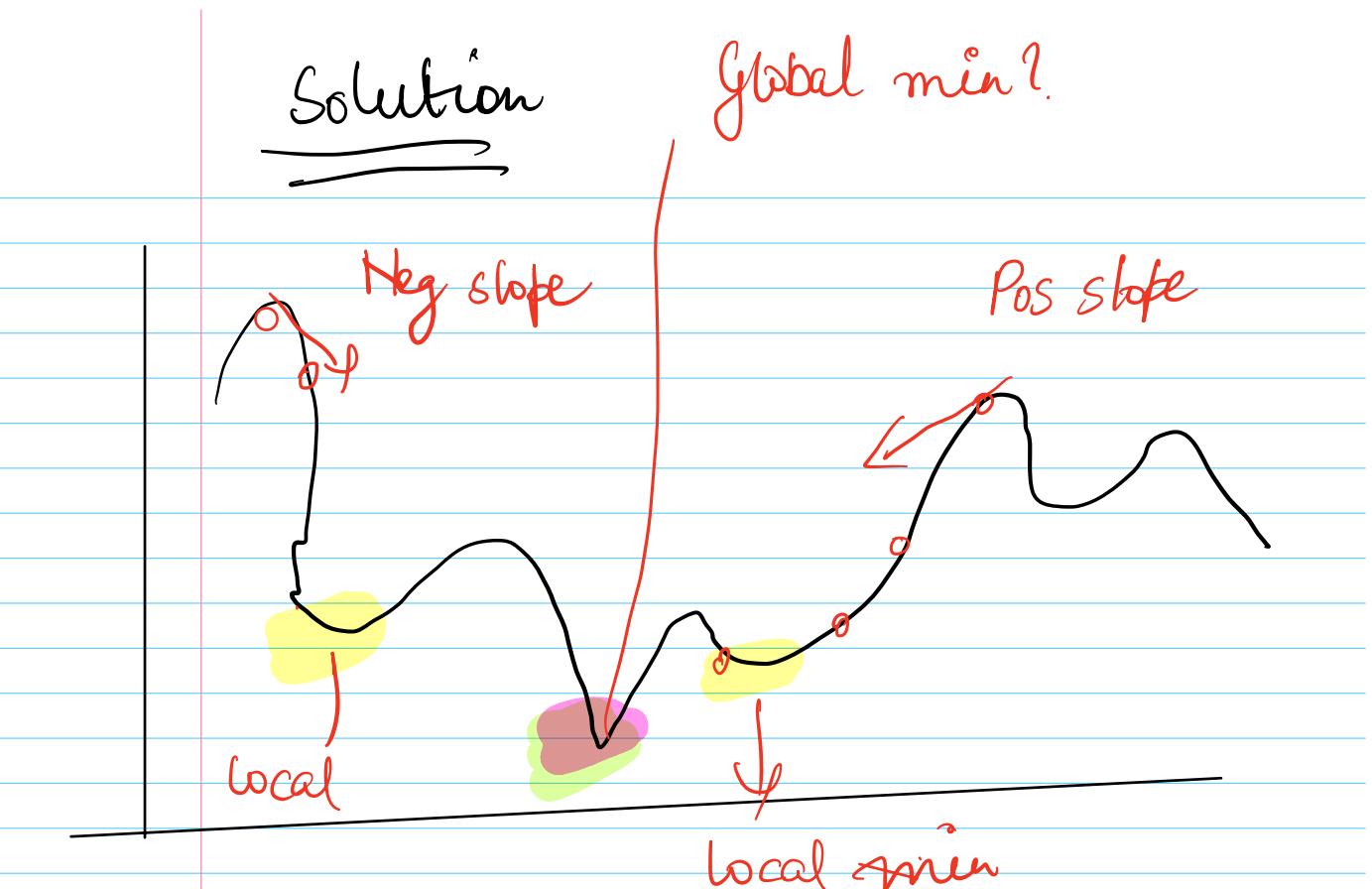


- ② Take a step towards hopefully better value of  $X$ .
- ③ Stop when happy or  $f(X)$  no longer  $\uparrow$

## ② Hm, & k Problems

→ How to decide which dir<sup>n</sup> to go to? in a sea of million params

→ How much to hop, step size?



① Where to move

-ve derivative  $\rightarrow$  move right

+ve  $\rightarrow$  move left

So, a rough "good" soln

① Init X

② while  $f'(x^k) \neq 0$   
if sign ( $f'(x^k)$ ) +ve  
 $x^{k+1} = x^k - \text{step}$

else

$$x^{k+1} = x^k + \text{step}$$

② How much to step?

Q.) Select all true about derivatives

✓ At any loc  $X$ , there may many dir $^n$  in which we can step such that  $f(x) \uparrow$

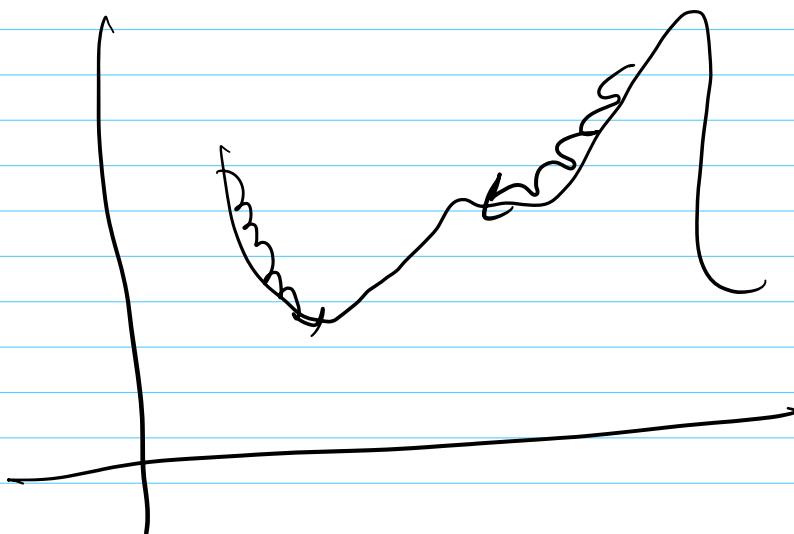
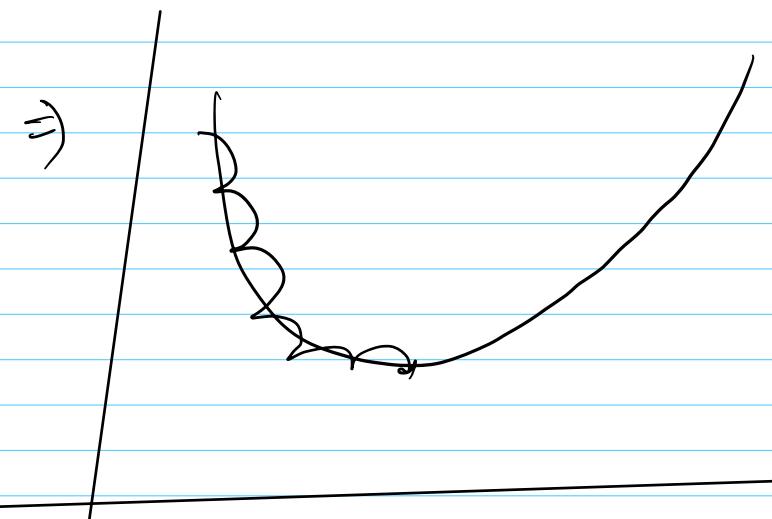
✗ Dir $^n$  of gradient is

dir<sup>n</sup> of steepest / fastest ↑

→ Gradient & derivative

→ Dir<sup>n</sup> of gradient, f ↑

\* More opposite, f ↓



→ for convex, appropriate step size will always find min

→ for non-convex it will find local minima/ saddle/inflection points

Hmm.. We should use this to find min of our neural network.

Algo

① Init  $w^0$ ,  $k \in \text{Step}$

② do  $w^{k+1} = w^k - \eta^k \nabla \text{Loss}(w^k)^T$

$$k = k + 1$$

while

$$|\text{loss}(\omega^k) - \text{Loss}(\omega^{k+1})| > \epsilon$$

Coming back to Our NN

⇒ Our function

a NN

⇒ We do init, as discussed  
above

⇒ We create function which  
is NN

⇒ We attach loss function

to our function , and  
finally compose a  
differentiable func<sup>n</sup> ??

Why? So that we can  
use iterative solution  
to find out best params

$$\text{loss} = f(\text{Inp}, \omega, \text{outputs})$$

a

$$\text{loss} = \frac{1}{T} \sum_{N=1}^N \text{loss}(f(x, \omega), \text{label})$$

$\Rightarrow$  So we do forward pass  
 $\rightarrow$  to get this loss

→ Then get gradient

→ Then take a step

⇒ So why we prefer one loss function over another?

⇒ Better loss surface,

the gradient should push us to faster convergence

when prediction deviates from label

↳ which depends on loss function

Hmn, we are getting closer to solve over biexcess !!!

## Now problem

① We must figure out  
how to compute  
derivative of loss w.r.t  
each weight

### A Detour

We saw

$$\Delta y \approx \frac{dy}{dx} \Delta x$$

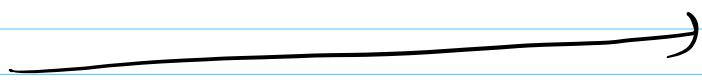
Influence diagram :

Edge = derivat

$\frac{dy}{dx}$

Node & edge

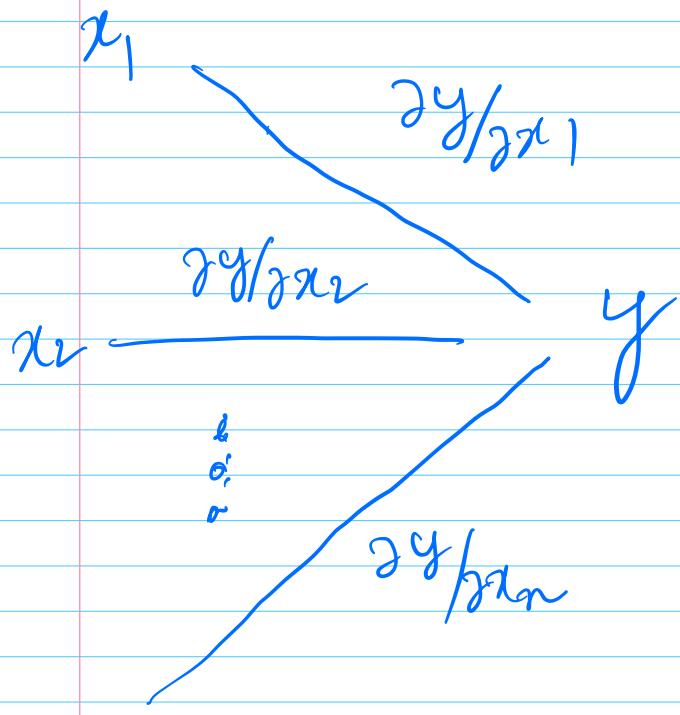
multiply



$$\Delta y = \frac{dy}{dx} \Delta x$$

9) What about scalar out,  
vector in

$$y = f(x_1, x_2, x_3 \dots x_m)$$



$$\Delta y \approx \frac{\partial y}{\partial x_1} \Delta x_1 +$$

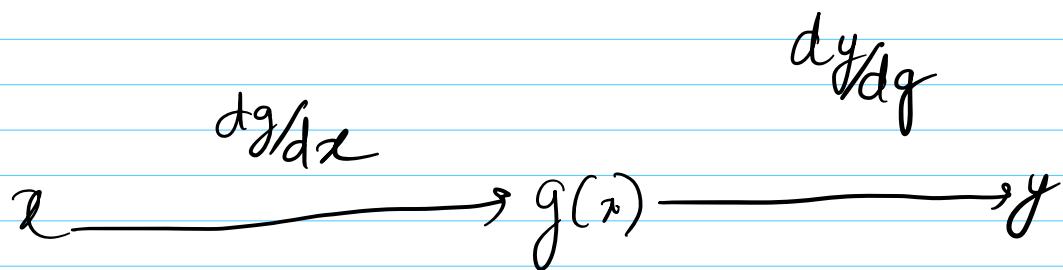
$$\frac{\partial y}{\partial x_2} \Delta x_2 +$$

⋮

$$\frac{\partial y}{\partial x_m} \Delta x_m$$

What about

$$y = f(g(x))$$

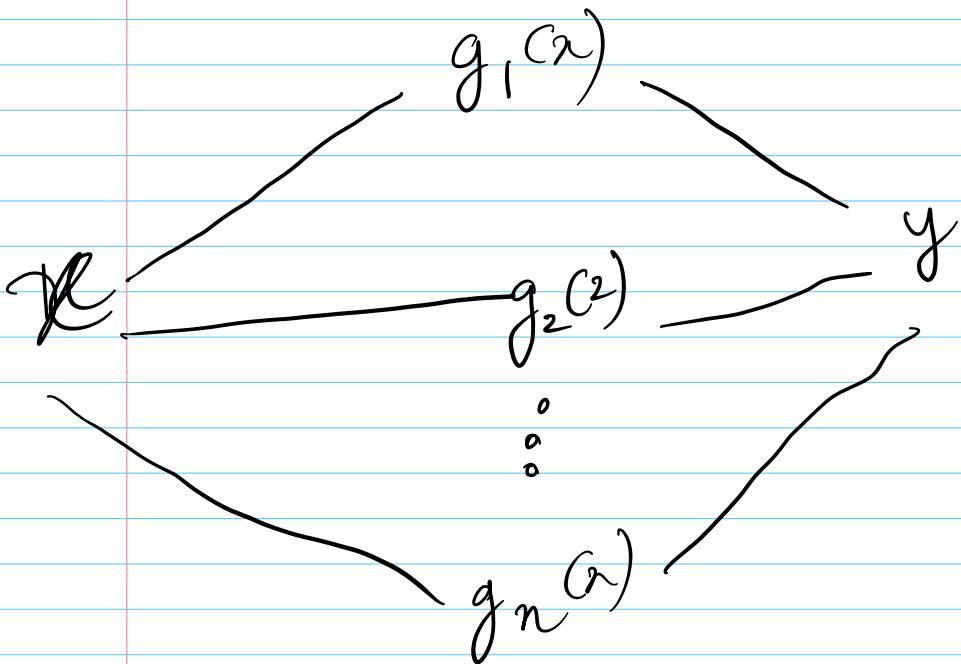


$$\Delta y = \frac{dy}{dg} \Delta g; \Delta g = \frac{dg}{dx} \Delta x$$

$$\Delta y = \frac{dy}{dg} \frac{dg}{dx} \Delta x$$

What about

$$y = f(g_1(x), g_2(x), \dots)$$



$$\Delta y = \frac{dy}{dg_1} \frac{\Delta g_1}{dx} + \frac{dy}{dg_2} \frac{\Delta g_2}{dx} + \dots$$

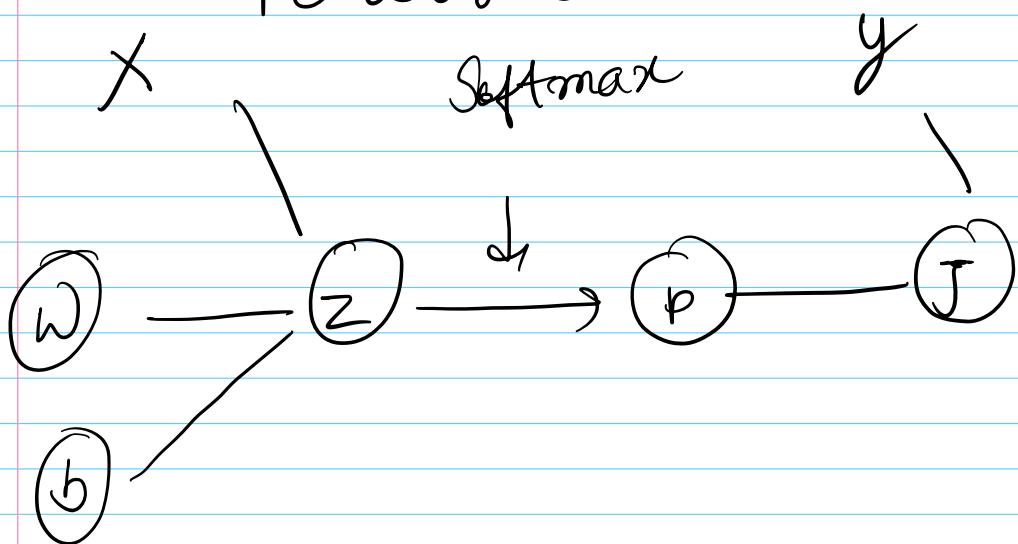


chain rule

Coming back

lets take our Neural

Network



$$\text{Loss} = \sum_i y_i \log p$$

⇒ forward pass is required  
to compute derivative,

$$\frac{d \text{Loss}(Y, p)}{1 \dots k}$$

$$d w_{i,j}$$

which requires intermediate  
and final output

⇒ So using chain rule,

$$\frac{\partial J}{\partial w} = \underbrace{\frac{\partial J}{\partial p}}_{[1 \times n_c]} \times \underbrace{\frac{\partial p}{\partial z}}_{n_c \times n_c} \times \underbrace{\frac{\partial z}{\partial w}}$$

for now, without going into  
math, let me tell you

that

$$\frac{\partial z}{\partial w} = \frac{\partial (w^T x)}{\partial w} = x$$

$$\frac{\partial J}{\partial p} \times \frac{\partial p}{\partial z} = (p_k - y_i) d_2$$

$$\text{Eg: } y_i = 2 \quad [0, 1, 0] \\ p = [0.2, 0.3, 0.5)$$

$$[0.2-0, 0.3-1, 0.5-0]$$

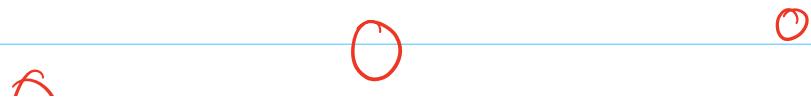
$$= [0.2, -0.7, 0.5]$$

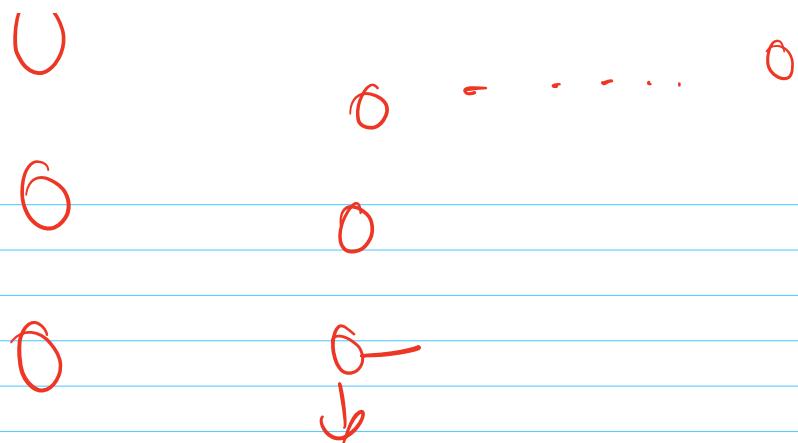
Refer to Notebook

$$\frac{\partial \text{Div}}{\partial w_{ij}} = \sum \frac{\partial \text{Div}}{\partial z_j} \frac{\partial z}{\partial y_i}$$

$$\frac{\partial \text{Div}}{\partial w_{ij}^{(N)}} = y_i \frac{\partial \text{div}}{\partial z_j^{(N)}}$$

for any inside layer

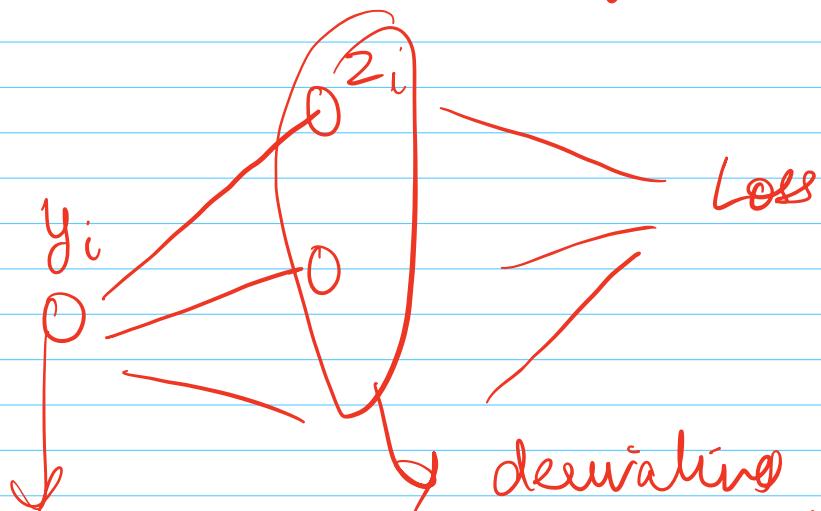




If we do back  
then derivative

$\frac{\partial \text{Loss}}{\partial z_i}^{(N+1)}$  has already  
been  
computed

which is simply



derivative of these  
already  
computed

So