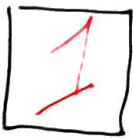




# Multi-layer perceptron from scratch:

Input  $X$ :  $28 \times 28$    $28 \times 28$   ...   $28 \times 28$   
 $m$  images

Output  $Y$ : 1 2 ... 5

1- Reshaping the input:

$$X[m, 28, 28] \Rightarrow X[m, 784] \Rightarrow X[N_x, m]$$

$N_x$ : # features  
 $m$ : # samples

$$Y[m, ] \Rightarrow Y[1, m]$$

$X$ :

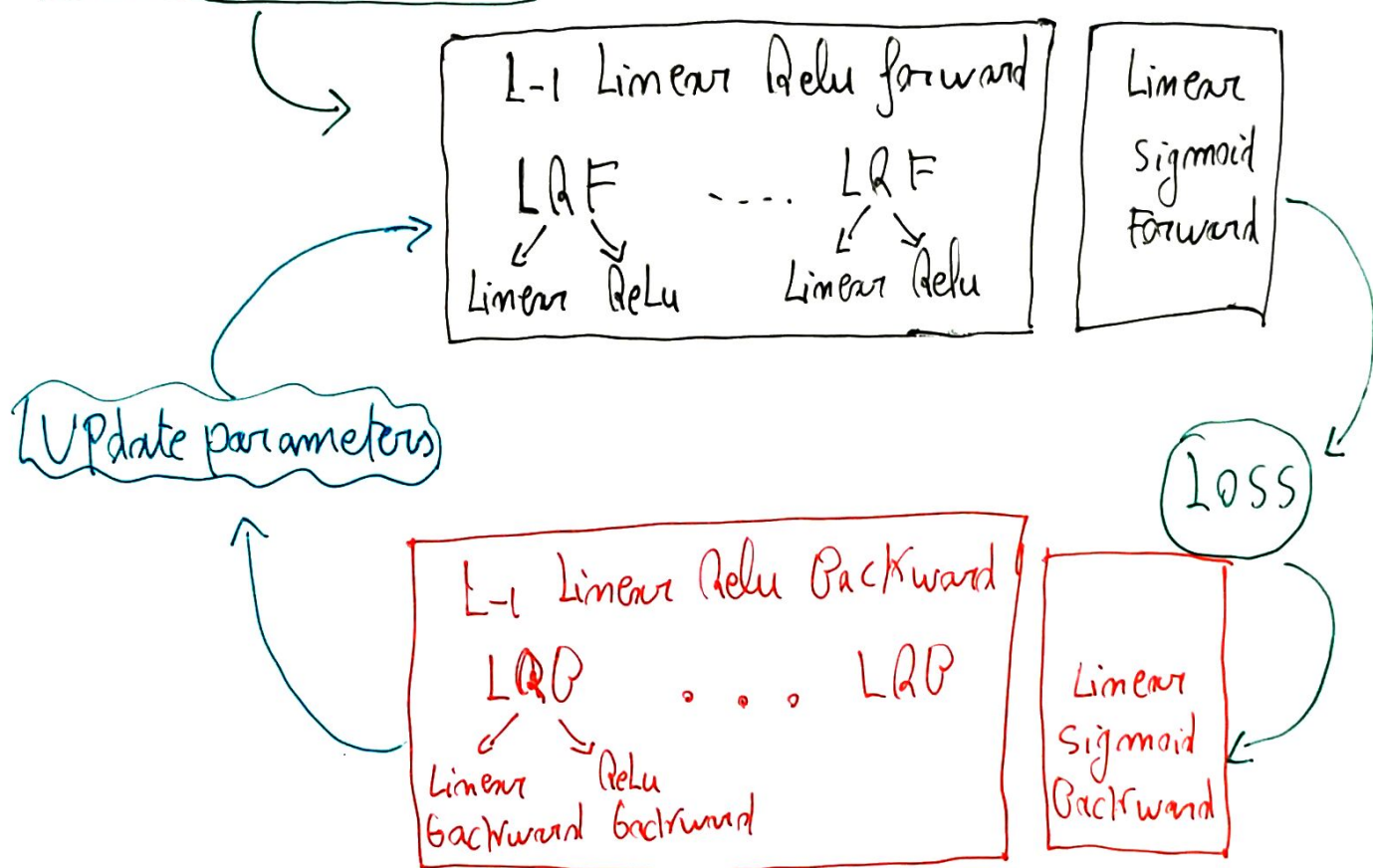
	$P_1$	$P_2$	$P_3$	...	$P_{784}$
$i_1$	0	0	0	-	-
$i_2$	0	0	0	-	-
$\vdots$	0	255	0	-	-
$i_m$	0	0	56	192	0

$$Y: [1 \ 1 \ \dots \ 5]$$

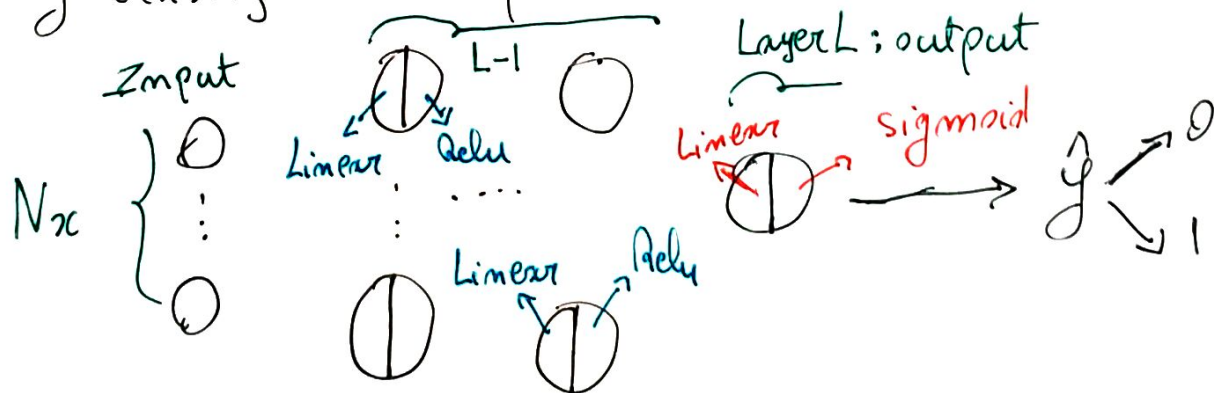
Objective: Building an L-layer neural network

- Strategy:
- initialize the parameters for  $L$ -layer NN
  - Forward propagation
  - Compute loss
  - Backpropagation
  - Update parameters.

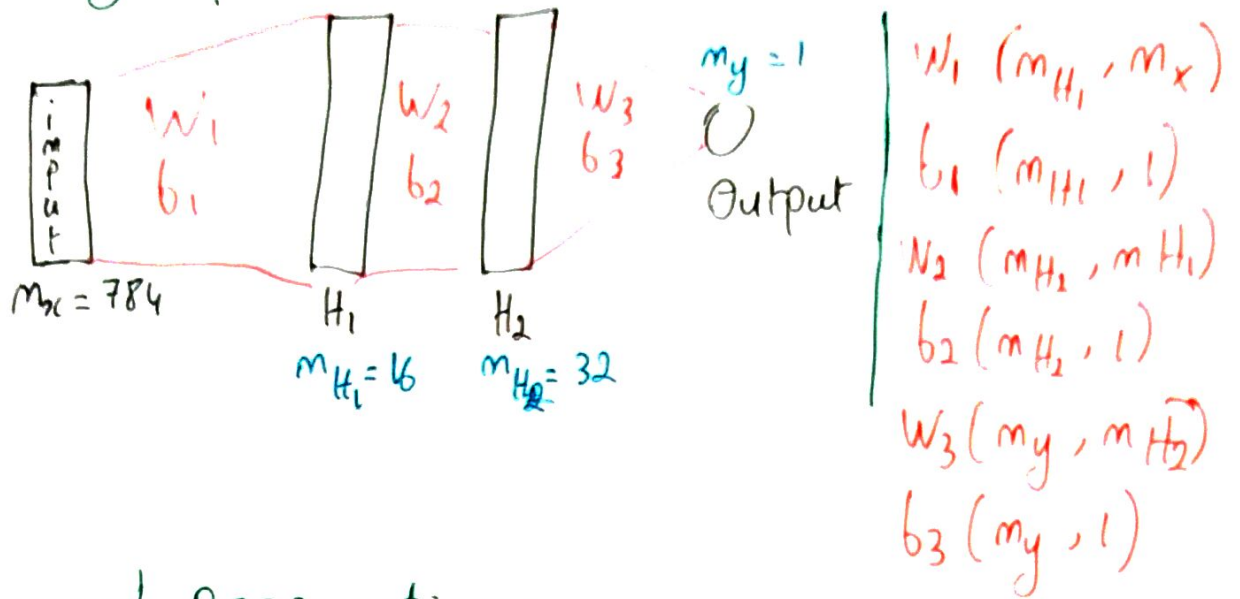
Initialize all params  
 $W_1, b_1 \dots W_L, b_L$



Binary classification problem:



## 2- Initialize parameters



## 3- Forward propagation:

Input  $X(m_x, m)$

$\Downarrow \begin{matrix} w_1 \\ b_1 \end{matrix}$

Linear  $Z_1 = w_1 X + b_1 (m_{H_1}, m)$

Activation  $A_1 = G(Z_1) = \text{Relu}(Z_1) = \max(0, Z_1) (m_{H_1}, m)$

$\Downarrow \begin{matrix} w_2 \\ b_2 \end{matrix}$

Linear  $Z_2 = w_2 A_1 + b_2 (m_{H_2}, m)$

Activation  $A_2 = \max(0, Z_2) (m_{H_2}, m)$

$\Downarrow \begin{matrix} w_3 \\ b_3 \end{matrix}$

Linear  $Z_3 = w_3 A_2 + b_3 (\overset{\pm}{m_y}, m)$

Activation  $A_3 = G(Z_3) = \text{Sigmoid}(Z_3) = \frac{1}{1 + e^{-Z_3}} (m_y, m)$

4- Compute cost

output  $A_1$   $(1, m)$

True labels  $Y$   $(1, m)$

binary classification

$\Downarrow$   
binary crossentropy

$x$ : images

label: number = 3 , number  $\neq 3$

what is the probability of an image being = 3 ideal = 1

Positive class : image = 3

Negative class : image  $\neq 3$

Our model will predict a probability of an image being = 3. How good or bad are the predicted probabilities?

Loss function  $\rightarrow$  Low values for good predictions  
 $\rightarrow$  High values for bad predictions

$$J(w_1, b_1, \dots, w_L, b_L) = -\frac{1}{m} \sum_{i=1}^m y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i))$$

$p(y_i)$  is the predicted probability of an image being = 3



## Why it works?

Suppose:  $y=1$   
 $p(y)=0.9$



$$y \log(p(y)) + (1-y) \log(1-p(y))$$

○

$$1 \log 0.9 = \boxed{-0.105}$$

$y=1$   
 $p(y)=0.2$



$$y \log(p(y)) + (1-y) \log(1-p(y))$$

○

$$1 \log 0.2 = \boxed{-1.609}$$

$y=0$   
 $p(y)=0.1$



$$y \log(p(y)) + (1-y) \log(1-p(y))$$

○

$$1 \log(1-0.1)$$
$$1 \log(0.9) = \boxed{-0.105}$$

$y=0$   
 $p(y)=0.8$



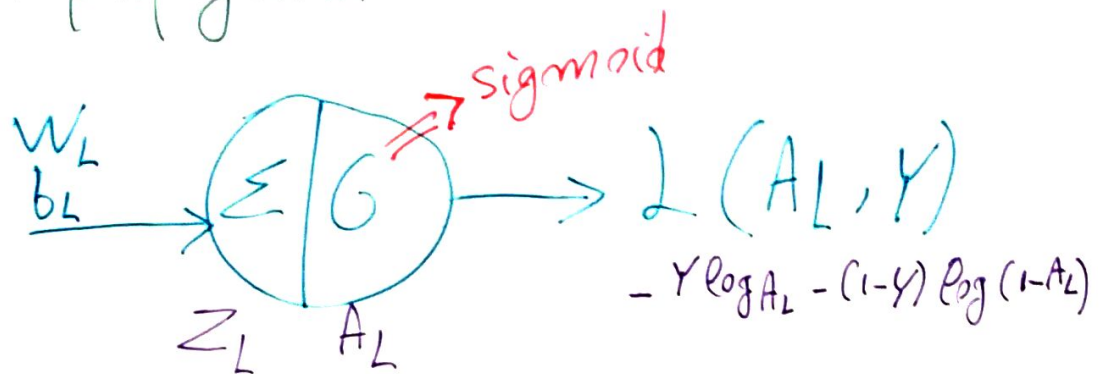
$$y \log(p(y)) + (1-y) \log(1-p(y))$$

○

$$1 \log(1-0.8)$$
$$1 \log 0.2 = \boxed{-1.609}$$

we add the minus sign since the objective is to minimize error.

## 5- Backward propagation



Chain rule:  $\frac{\partial L}{\partial w_L} = \frac{\partial L}{\partial A_L} \cdot \frac{\partial A_L}{\partial z_L} \cdot \frac{\partial z_L}{\partial w_L}$

$$\frac{\partial L}{\partial A_L} = -\frac{Y}{A_L} + \frac{1-Y}{1-A_L}$$

$\hookrightarrow \frac{dz_L}{dz_L} A_{L-1}^T$

$$\frac{\partial L}{\partial z_L} = \frac{\partial L}{\partial A_L} \cdot \frac{d}{dz_L} g(z_L) = dA_L \cdot g'(z_L)$$

Pseudo code:

$$dA[L] = -\left(\frac{Y}{A[L]}\right) + \frac{1-Y}{1-A[L]}$$

For each layer  $l$ :

input:  $dA[l]$

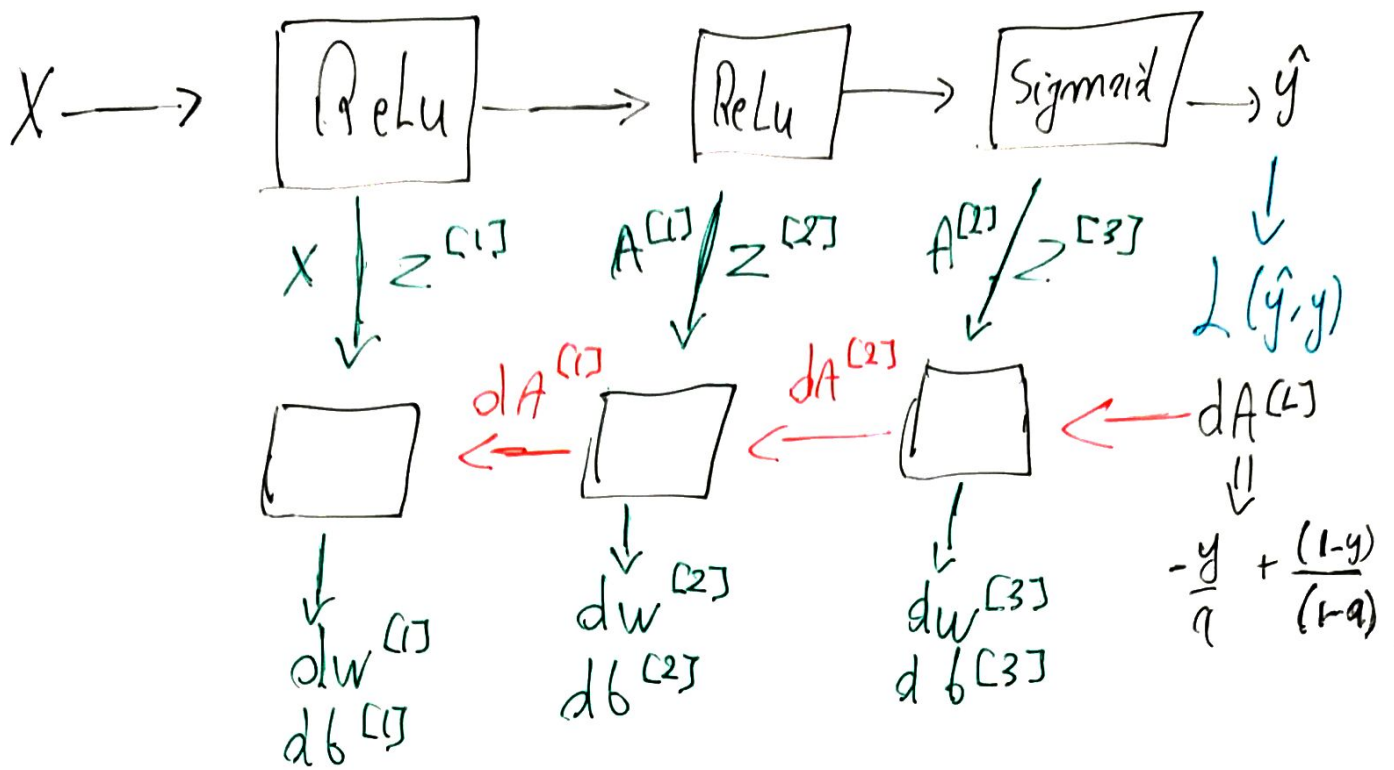
$$dz[l] = dA[l] \cdot g'(l)(z[l])$$

$$dw[l] = \frac{1}{m} dz[l] A[l-1]^T$$

$$db[l] = \frac{1}{m} \sum dz[l]$$

$$dA[l-1] = w[l]^T dz[l]$$

output  $dA[l-1]$



6- update parameters

for each layer  $l$ :

$$w_l = w_l - \alpha dw_l$$

$$b_l = b_l - \alpha db_l$$

$\alpha$  is the learning rate.

$\hat{z}$  - predict

