

Model Architecture

→ On MNIST handwritten digits dataset $\xrightarrow{\text{image}}$ shape = 28×28 $\xrightarrow{\text{flatten}}$ shape = 1×784

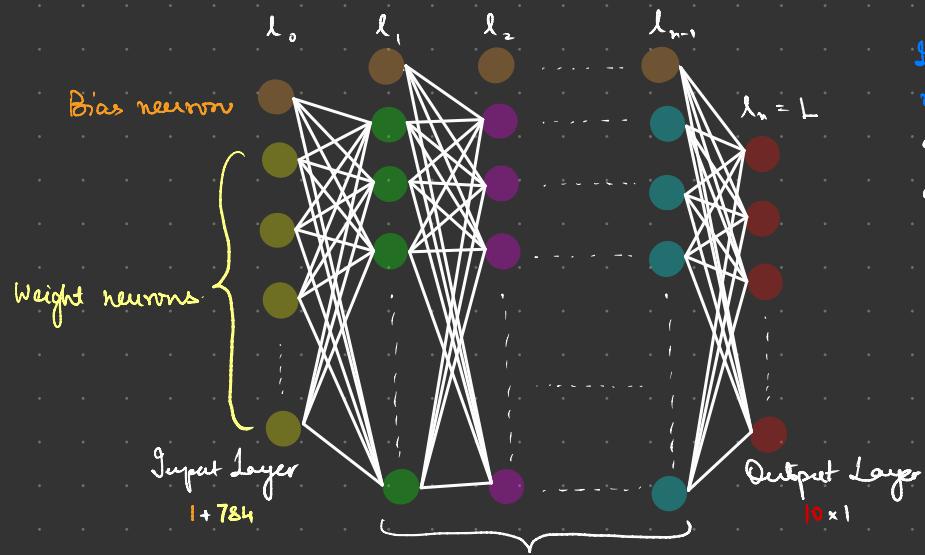


Image flattened into a row vector concatenated with a bias term

Hidden layers
 $1+n_i$
 $1 \leq i \leq n-1$

→ which only has a forward connection with the next layer.

Initialising a base architecture (before activations, etc) require us to specify the number of neurons in each layer.

Eg: $[784, 40, 30, 20, 10]$, qualifies for us to build the weights & biases list. The bias neuron in l_1 will have 40 connections with l_2 , while the bias neuron in l_2 will have 30 connections with l_3 , and so on.. Therefore the biases look like $\rightarrow [40 \times 1, 30 \times 1, 20 \times 1, 10 \times 1]$. Now each of the 784 weight neurons in l_1 will have 40 connections with l_2 , while each of the 40 weight neurons in l_2 will have 30 connections with l_3 , and so on.. Therefore weights look like $\rightarrow [40 \times 784, 30 \times 40, 20 \times 30, 10 \times 20]$

↑
Softmax over the output which are digits from 0 to 9.

① Constant

② Random - Uniform \Rightarrow $U[\min, \max]$

↳ Normal $\Rightarrow N[\text{mean}, \text{std}]$

$$\textcircled{3} \text{ Xavier-Glorot - Uniform } \Rightarrow U\left[-\sqrt{\frac{6}{n_{in} + n_{out}}}, +\sqrt{\frac{6}{n_{in} + n_{out}}}\right]$$

$$\text{Normal} \Rightarrow N\left[0, \frac{2}{n_m + n_{out}}\right]$$

$$\textcircled{5} \text{ He} \text{ --- Uniform} \Rightarrow U\left[-\sqrt{\frac{3}{n_{in}}}, +\sqrt{\frac{3}{n_{in}}}\right]$$

$$\text{Normal} \Rightarrow N\left[0, \sqrt{\frac{2}{n}}\right]$$

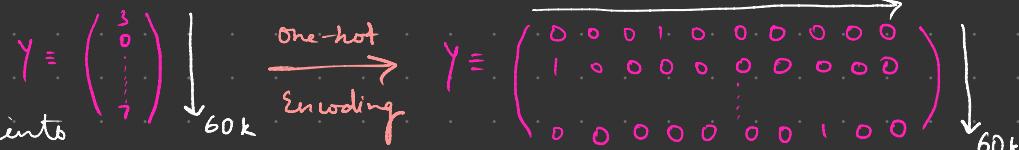
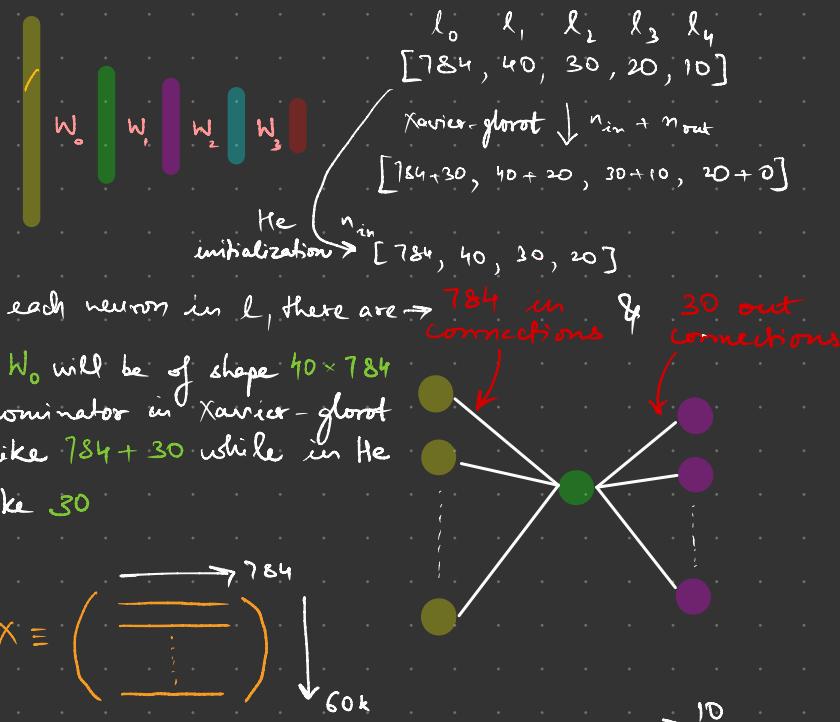
Images: Because pixels range from $[0, 255]$

Divide by 255 $\rightarrow [0, 1]$

Digits (output). range from 0 to 9

These class values will be encoded into

a vector of length = number of classes, which is greater than equal to 3. This is done to support matrix product during "softmax" at the final layer. If we have 2 classes, then using sigmoid at last layer doesn't require us to encode the output classes.



Activation Functions

① Sigmoid — $s(z) = \frac{1}{1 + \exp(-z)}$

$1 - s(z) = 1 - \frac{1}{1 + \exp(-z)} = \frac{\exp(-z)}{1 + \exp(-z)}$

$\frac{\partial s(z)}{\partial z} = \frac{\partial}{\partial z} \left(\frac{1}{1 + \exp(-z)} \right) = (-1) \left(\frac{1}{1 + \exp(-z)} \right)^2 \frac{\partial (1 + \exp(-z))}{\partial z} = (-1) \left(\frac{1}{1 + \exp(-z)} \right)^2 (0 + (-1) \exp(-z)) = \frac{\exp(-z)}{(1 + \exp(-z))^2}$

$$= \left(\frac{1}{1 + \exp(-z)} \right) \left(\frac{\exp(-z)}{1 + \exp(-z)} \right) = s(z) \odot (1 - s(z))$$

② Leaky ReLU (leak = 0 \Rightarrow ReLU) — $\text{relu}(z) = \begin{cases} z & \text{if } z \geq 0 \\ \alpha z & \text{if } z < 0 \text{ where } \alpha \geq 0 \end{cases}$

$\frac{\partial}{\partial z} \text{relu}(z) = \begin{cases} \frac{\partial}{\partial z} z & \text{if } z \geq 0 \\ \frac{\partial}{\partial z} \alpha z & \text{otherwise} \end{cases} = \begin{cases} 1 & \text{if } z \geq 0 \\ \alpha & \text{otherwise} \end{cases}$

③ Tanh — $\tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)}$

$\frac{\partial}{\partial z} \tanh(z) = \frac{[\exp(z) + \exp(-z)] \frac{\partial}{\partial z} [\exp(z) - \exp(-z)] - [\exp(z) - \exp(-z)] \frac{\partial}{\partial z} [\exp(z) + \exp(-z)]}{[\exp(z) + \exp(-z)]^2} = \frac{[\exp(z) + \exp(-z)]^2 - [\exp(z) - \exp(-z)]^2}{[\exp(z) + \exp(-z)]^2}$

$$= 1 - \tanh^2(z)$$

$$\frac{\partial}{\partial z} \tanh(z) = [1 - \tanh^2(z)] \odot [1 + \tanh(z)]$$

④ Softmax — $\phi(z) = \frac{\exp(z)}{\sum_i^c \exp(z_i)} = \frac{\exp(z - \max(z))}{\sum_i^c \exp(z_i - \max(z))}$

out of all logits
Numerically stable
(when $\exp(z)$ is large)

where $c = \text{Number of classes}$

Derivative of Softmax

Given that softmax depends on all of the input raw logits, taking the derivative becomes not so straightforward. Therefore, the derivative is partial when calculating for a single logit, so we have to consider the **Jacobian Matrix** while talking about the first order partials.

↑ constant, so ignoring this for derivation of the derivative

$$J = \begin{pmatrix} \frac{\partial z_1}{\partial z_i} & \frac{\partial z_2}{\partial z_i} & \cdots & \frac{\partial z_c}{\partial z_i} \\ \frac{\partial z_1}{\partial z_2} & \frac{\partial z_2}{\partial z_2} & & \vdots \\ \vdots & & & \ddots \\ \frac{\partial z_1}{\partial z_c} & & \cdots & \frac{\partial z_c}{\partial z_c} \end{pmatrix} \text{ where } q_i^l = \frac{\exp(z_i - \max(z))}{\sum_k^c \exp(z_k - \max(z))}$$

$$\text{Derivative of 'Log'} \Rightarrow \frac{\partial z_j}{\partial z_i} \log(q_i^l) = \frac{1}{q_i^l} \frac{\partial z_j}{\partial z_i} q_i^l$$

$$\frac{\partial z_j}{\partial z_i} q_i^l = q_i^l \cdot \frac{\partial z_j}{\partial z_i} \log(q_i^l)$$

$$\begin{aligned} \frac{\partial}{\partial z_j} (\exp z_1 + \dots + \exp z_j + \dots + \exp z_c) \\ = (0 + \dots + \exp z_j + \dots + 0) \\ = \exp z_j \end{aligned}$$

$$\text{Now, } \log(q_i^l) = z_i - \log(\sum_k^c \exp(z_k))$$

$$\rightarrow \frac{\partial z_j}{\partial z_i} \log(q_i^l) = \frac{\partial z_j}{\partial z_i} z_i - \frac{\partial z_j}{\partial z_i} \log(\sum_k^c \exp(z_k)) = \delta_{ij} - \frac{1}{\sum_k^c \exp(z_k)} (\partial z_j \sum_k^c \exp(z_k)) = \delta_{ij} - \frac{e^{z_i}}{\sum_k^c e^{z_k}} q_j^l$$

$$\therefore \frac{\partial z_j}{\partial z_i} q_i^l = q_i^l \cdot (\delta_{ij} - \frac{e^{z_i}}{\sum_k^c e^{z_k}} q_j^l)$$

This will come in handy while backpropagation, when the softmax layer stands trial against the actual encoded output classes by means of **Categorical Cross Entropy** as the loss/cost function.

Loss Functions

① Mean Squared Error — $C(a^L) = \frac{1}{2n} \sum_i [(y - a^L)^2]_i$

\downarrow
Derivative

$$\partial_{a^L} C = \frac{1}{2n} \sum_i [2(y - a^L)(-1)]_i = -\frac{1}{n} \sum_i [y - a^L]_i$$

$$\delta^L = \partial_{a^L} C \circ \partial_{z^L} a^L = -\frac{1}{n} \sum_i [(y - a^L) \circ \partial_{z^L} a^L]_i$$

② Binary Cross Entropy — $C(a^L) = -\frac{1}{n} \sum_i [y \ln a^L + (1-y) \ln(1-a^L)]_i$, where n = number of training instances

\downarrow
Derivative

$$\begin{aligned} \partial_{a^L} C &= -\frac{1}{n} \sum_i \left[\frac{y(1)}{a^L} + \frac{1-y(-1)}{1-a^L} \right]_i = -\frac{1}{n} \sum_i \left[\frac{y}{a^L} - \frac{1-y}{1-a^L} \right]_i = -\frac{1}{n} \sum_i \left[y - \frac{y a^L - a^L + a^L y}{a^L(1-a^L)} \right]_i \\ &= -\frac{1}{n} \sum_i \left[\frac{y - a^L}{a^L(1-a^L)} \right]_i, \text{ here } a^L \equiv \text{Sigmoid} \end{aligned}$$

$$\delta^L = \partial_{a^L} C \circ \partial_{z^L} a^L = -\frac{1}{n} \sum_i [y - a^L]_i$$

③ Categorical Cross Entropy — $C(a^L) = -\frac{1}{n} \sum_i \left[\sum_k y_k \ln a_k^L \right]_i$

\downarrow
Derivative

$$\partial_{a_k^L} C = -\frac{1}{n} \sum_i \left[\sum_k \frac{y_k}{a_k^L} \partial_{a_k^L} a_k^L \right]_i, \text{ here } a^L \equiv \text{Softmax}$$

$$\delta_j^L = \partial_{a_j^L} C \cdot \partial_{z_j^L} a_j^L = -\frac{1}{n} \sum_i \left[\sum_k \frac{y_k}{a_k^L} \partial_{a_j^L} a_k^L \cdot \partial_{z_j^L} a_j^L \right]_i = -\frac{1}{n} \sum_i \left[\sum_k \frac{y_k}{a_k^L} \partial_{z_j^L} a_k^L \right]_i = -\frac{1}{n} \sum_i \left[\sum_k \frac{y_k}{a_k^L} a_k^L (\delta_{kj} - a_j^L) \right]_i = -\frac{1}{n} \sum_i \left[\sum_k y_k \delta_{kj} - a_j^L \sum_k y_k \right]_i = -\frac{1}{n} \sum_i [y_j - a_j^L]_i$$

$$\delta^L = -\frac{1}{n} \sum_i [y - a^L]_i$$

Initialization Step

Given architecture $\rightarrow [784, 40, 30, 20, 10]$

$$\text{biases shape} = [40 \times 1, 30 \times 1, 20 \times 1, 10 \times 1]$$

$$\text{weights shape} = [40 \times 784, 30 \times 40, 20 \times 30, 10 \times 20]$$

Assume a batch contains 100 instances of flattened images

Forward Propagation

After processing the data & getting the batch, the forward pass starts. We need to store the logits & activations of each layer (to be used in BP) & update this memory at each epoch.

$$\text{logits_history} = [-, -, -, -]$$

$$\text{activations_history} = [-, -, -, -]$$

We have the input, weights & biases, choice of activation and a place to store the computations.

"Better Shape is all you need !!"

$$\text{Input } x \rightarrow (100, 784) \xrightarrow{\text{shape}} (784, 100)$$

$$\text{Weight } w_0 \rightarrow (40, 784)$$

$$\text{Bias } b_0 \rightarrow (40, 1)$$

$$\text{Logit } z_1 = w_0 a_0 + b_0 \rightarrow (40, 100)$$

$$\text{Activation } a_1 = \sigma(z_1) \rightarrow (40, 100)$$

$$\text{activations \& logits history} = [(40, 100), -, -, -]$$

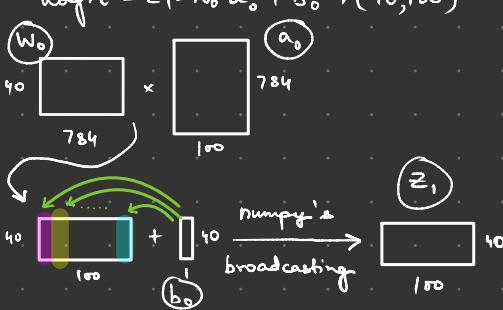
$$z_2 = w_1 a_1 + b_1 \rightarrow (30, 100)$$

$$(30, 40) \times (40, 100) + (30, 1)$$

$$(30, 100) + (30, 1) \xrightarrow[\text{for all layers}]{\text{100 on}} [(40, 100), (30, 100), (40, 100), (10, 100)]$$

= logit_history

This is useful in
completing last step of
Backpropagation



Activation is just an element-wise operation over the logit tensors.

$$[(784, 100)] + [(40, 100), (30, 100), (20, 100), (10, 100)]$$

= activation_history

At the last layer, softmax will be used

$$\sigma(z_i) \rightarrow (10, 100)$$

$$\sigma(z) = \frac{\exp(z - \text{max}(z))}{\sum \exp(z - \text{max}(z))}$$

$$\sum \exp(z - \text{max}(z))$$

This operation would reduce the dimension if used by default
So set $\text{keepdims} = \text{True}$

$\text{axis} = 0$

The 'max' & 'sum' operations are performed for each column separately, because each column is a different data vector in the instance space of all images

Backpropagation

→ (Element-wise) Hadamard product

$$\text{Four Equations of Backpropagation} \rightarrow \cdot \delta^L = \nabla_a C \odot \sigma'(z^L) \quad \cdot \delta^L = ((W^{L+1})^T \delta^{L+1}) \odot \sigma'(z^L) \quad \cdot \frac{\partial C}{\partial b^L} = \delta^L \quad \cdot \frac{\partial C}{\partial W^L} = \delta^L (a^{L-1})^T$$

To store the quantities: $\frac{\partial C}{\partial b^L}$ & $\frac{\partial C}{\partial W^L}$, we need to initialize empty lists while making sure to generate tensors of shape just like the bias & weight so that we can cross check any shape errors in matrix operations.

grad_bias = $[40 \times 1, 30 \times 1, 20 \times 1, 10 \times 1]$ of zeros

grad_weight = $[40 \times 784, 30 \times 40, 20 \times 30, 10 \times 20]$ of zeros

(100, 10)

(10, 100)

Now, the process starts in reverse, ∴ we start with actual output (encoded) vs last activation (prediction)

$$\text{delta} = -(\text{actual_output.T} - \text{activation_history}[-1]) \rightarrow (10, 100)$$

$$\text{grad_bias}[-1] = \begin{pmatrix} \sum \text{delta}[0] \\ \vdots \\ \sum \text{delta}[9] \end{pmatrix} \rightarrow (10, 1) \equiv \text{bias}[-1]$$

$$\text{grad_weights}[-1] = \text{delta} * \text{activation_history}^T[-2].T \rightarrow (10, 20) = \text{weights}[-1]$$

update delta to update at [-2]

$$\text{delta} = (\text{weights}[-1].T * \text{delta}) \odot \sigma(\text{logits}[-2]) \rightarrow (20, 100)$$

(20, 10)

(10, 100)

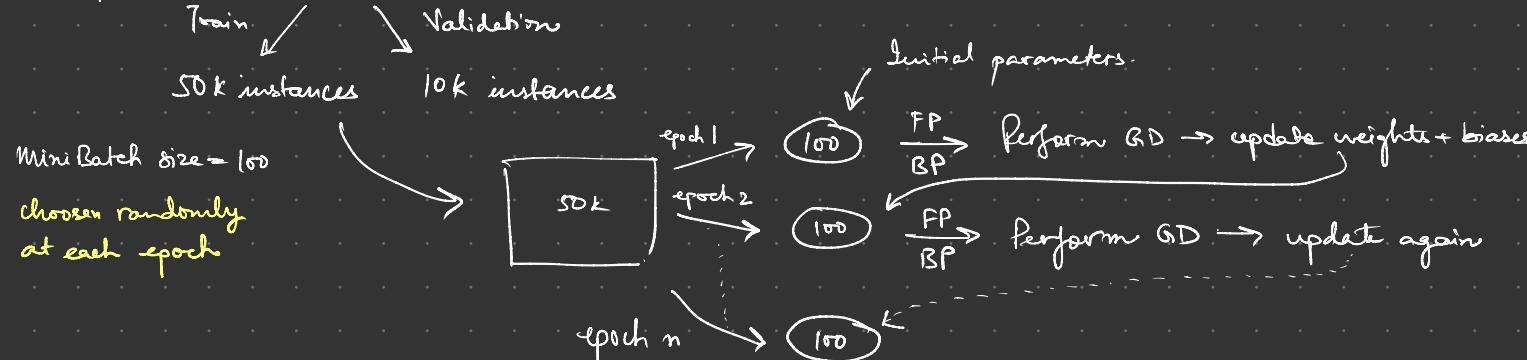
(20, 100)

→ & so on till the input layer

Normalize these gradients by the Batch-size or the Accuracy during fit would not improve !!

Data Splitting and Mini-Batch Gradient Descent

$$X = (60k, 784) \text{ , } y = (60k,)$$



To add momentum to the optimizer, we need to initialize a new list that remembers what were the gradients (scaled with momentum) in the previous epoch

$$\rightarrow \text{velocity} = [(\text{momentum}) * \text{velocity}] + [(\text{1 - momentum}) * \text{grad}]$$

$$\text{weight/bias} = \text{weight/bias} - [\text{lr} * \text{velocity}]$$