

Approach A adds a false sense of distance between the classes of digits. For eg. If the desired output is say 2, it will think the output 3 is closer to 2 than the output 5 because of the definition of error function. And it will train the weights accordingly. while in reality output 5 is just as wrong as output 3. Hence learning is misaligned with the ultimate goal.

A linear activation frame simply returns a weighted sum of inputs without any non linear transformation. The given problem is that of classification, we require a probabilistic output.

Approach B uses one hot encoding which is used to represent categories as vectors. Here we use a vector of length 10 to represent each of the 10 possible digit class(0-9):

For eg if digit=2 $y=[0,0,1,0,0,0,0,0,0,0]$. This clearly indicates which class is correct without creating a false sense of distance between classes.

Softmax converts raw network outputs (called logits) into probabilities.

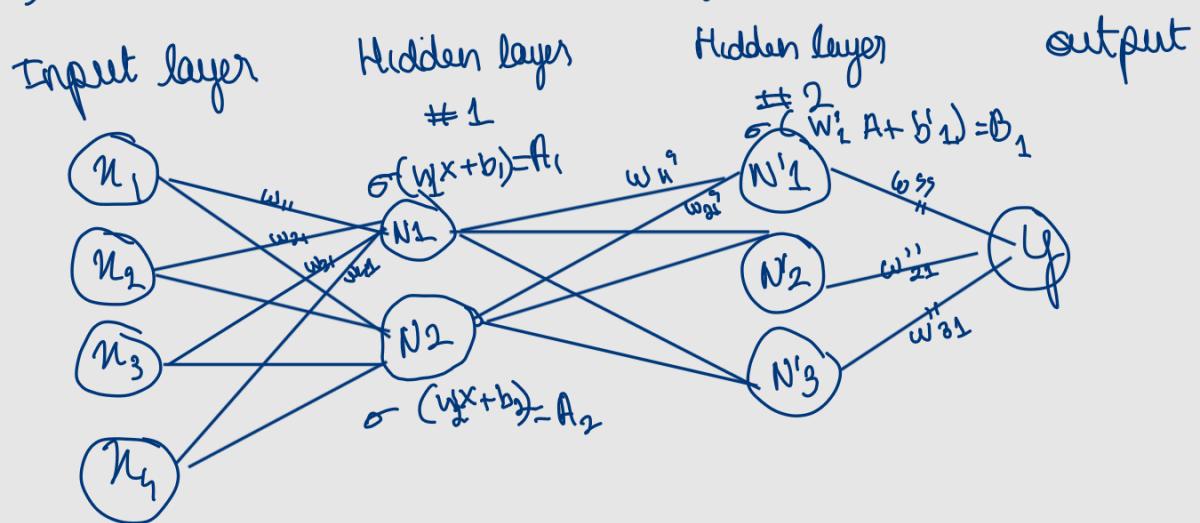
Example output $[0.01, 0.03, 0.92, 0.01 \dots]$

For outputs z_0, z_1, \dots, z_q :

$$\text{Softmax}(z_i) = \frac{e^{z_i}}{\sum_{j=0}^q e^{z_j}}$$

this approach treats all classes equally, produces probability output for each digit and final uses that to predict the most probable digit.

2) a) Network architecture diagram



b) Let input vector $X = [u_1, u_2, u_3, u_4]$
 weights are represented in matrix format
 Hidden layer #1

	N_1	N_2
u_1	w_{11}	w_{12}
u_2	w_{21}	w_{22}
u_3	w_{31}	w_{32}
u_4	w_{41}	w_{42}

Let this be column vector X this be row vector W_1 this be W_2

$\therefore \text{output of neuron } 1 = \sigma(z_1) = A_1 \text{ (let)}$
 where $z_1 = u_1 w_{11} + u_2 w_{21} + u_3 w_{31} + u_4 w_{41} + b_1$
 $= W_1 \times X + b_1$

$\text{output of neuron } 2 = \sigma(z_2) = A_2 \text{ (let)}$
 $z_2 = u_1 w_{12} + u_2 w_{22} + u_3 w_{32} + u_4 w_{42} + b_2$
 $= W_2 \times X + b_2$

and $\sigma(z) = \max(0, z) \rightarrow \text{ReLU fn}$

Output of Hidden Layer 1 = Input of H2

$$A = \text{vector } (A_1, A_2)$$

Hidden layer #2

	N'_1	N'_2	N'_3
A_1	w'_{11}	w'_{12}	w'_{13}
A_2	w'_{21}	w'_{22}	w'_{23}
A	w'_1	w'_2	w'_3

Output of $N'_1 = \sigma(z) = B_1$

$$z_1 = w'_{11} A_1 + w'_{21} A_2 + b'_1$$

$$= W'_1 \times A + b'_1$$

Output of $N'_2 = \sigma(z) = B_2$

$$z_2 = w'_{12} A_1 + w'_{22} A_2 + b'_2$$

$$= W'_2 \times A + b'_2$$

$$\text{output of } N_3 = \sigma(z_3) = B_3$$

$$z_3 = w_{13}^T A_1 + w_{23}^T A_2 + b_3^T$$

$$\text{output vector} = [B_1, B_2, B_3]$$

Output layer (single neuron)

$$y = w_{11}'' B_1 + w_{21}'' B_2 + w_{31}'' B_3 + b_1^r$$

c) Let $x = (-1, 0, 1, 2)$

$$\text{Let } w_1 = (1, 0, -1, 2), b_1 = 0$$

$$\begin{aligned} z_1 &= w_1 x + b_1 \\ &= -1 + 0 + -1 + 4 + 0 \\ &= 2 \end{aligned}$$

$$A_1 = \sigma(z_1) = 2$$

$$\text{Let } w_1^1 = (1, 1) \quad b_1^1 = 1$$

$$\begin{aligned} z_1^1 &= -2 + 0 + 1 \\ &= -1 \end{aligned}$$

$$B_1 = \sigma(z_1^1) = 0$$

$$\text{Let } w_2^1 = (0, 1) \quad b_2^1 = 0$$

$$\begin{aligned} z_2^1 &= 0 + 0 + 0 \\ &= 0 \end{aligned}$$

$$B_2 = \sigma(z_2^1) = 0$$

$$\text{Let } w_3^1 = (-1, 0) \quad b_3^1 = 1$$

$$\begin{aligned} z_3^1 &= -2 + 0 + 1 \\ &= -1 \end{aligned}$$

$$B_3 = \sigma(z_3^1) = 0$$

$$\text{Let } w_1''' = (1, 2, 3) \quad b_1''' = 0$$

$$y = 0 + 0 + 0 + 0 = 0$$

d) parameters

$$\text{Layer 1: } 4 \times 2 + 2 = 10$$

$$\text{Layer 2: } 2 \times 3 + 3 = 9$$

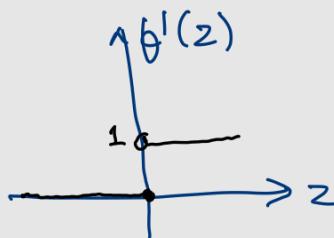
$$\text{Layer 3: } 3 + 1 = 4$$

3) a) Derivative of $f(z)$ w.r.t z

$$f(z) = \max(0, z) \therefore = \begin{cases} 0 & ; z \leq 0 \\ z & ; z > 0 \end{cases}$$

$$f'(z) = \begin{cases} 0 & ; z < 0 \\ \text{undefined} & ; z = 0 \\ 1 & ; z > 0 \end{cases}$$

* most libraries define $f'(z)=0$ at $z=0$



b) given pre-activation sum $z < 0$ for all inputs

$$\therefore f(z) = 0 \quad \therefore f' = 0$$

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial w_i} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_i} = 0$$

c) No, there is no change in Loss L^n due to weights or bias
bc $\frac{\partial L}{\partial w_i} \& \frac{\partial L}{\partial b} = 0$ bc of the ReLU function.

during training using gradient descent

$$w_f = w_i - (\eta \cdot r) \frac{\partial L}{\partial w_i} = w_i \quad b_f = b_i - (\eta \cdot r) \frac{\partial L}{\partial b} = b_i$$

\therefore there is no change in weights or bias during training.

\therefore we can say that bc of ReLU fn, any neuron that produces a negative z for all its inputs, doesn't affect the loss L^n and hence its weights and bias will never get trained in a way neuron becomes dead.

4) a) $J_{\text{total}}(w) = J_{\text{data}}(w) + \frac{\lambda}{2} w^2$

↓
total loss ↓
(MSE cross entropy loss) $\lambda > 0$ regularisation strength
 w = single weight parameter

$$w_{\text{new}} = w_{\text{old}} - \eta \left(\frac{\partial J_{\text{total}}(w)}{\partial w} \right)$$

$$\frac{\partial J_{\text{tot}}}{\partial w} = \frac{\partial J_{\text{data}}}{\partial w} + \frac{\lambda}{2} 2w$$

$$\therefore w_{\text{new}} = w_{\text{old}} - \eta \left(\frac{\partial J_{\text{data}}}{\partial w} + \lambda w_{\text{old}} \right)$$

b) $w_{\text{new}} = w_{\text{old}} \underbrace{(1 - \eta \lambda)}_{\text{factor}} - \eta \frac{\partial J_{\text{data}}}{\partial w}$

slightly $< 1 \quad \therefore \eta > 0 \quad \lambda > 0$

c) L2 regularisation is called weight decay because at every gradient descent step the weights are multiplied by a shrinking factor, even before applying the data gradient.

Even if the data gradient is 0, $w_{\text{new}} = (1 - \eta \lambda) w_{\text{old}}$

the weight still decays towards 0.

This discourages large weights and reduces overfitting

$$\nabla \left(\frac{\lambda}{2} w^2 \right) = \lambda w \rightarrow \text{points away from origin}$$

since its $-\lambda w \rightarrow$ points towards origin

\therefore the shrink step pulls the point radially inward towards origin by contracting its weight vector

and the data gradient moves in the direction that reduces data error.