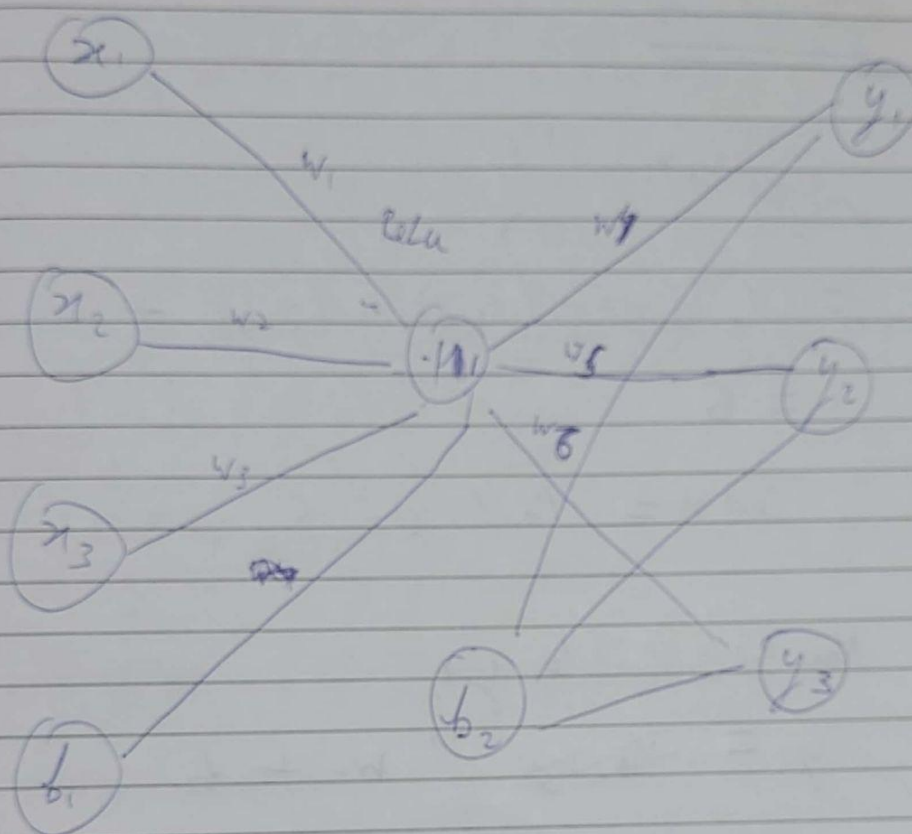


Section - C

- 1) These are the optimal hyper parameters for this model
`best_params = {'batch_size': 64, 'hidden_layer_sizes': (25, 25), 'learning_rate_init': 0.01, 'max_iter': 20}`
- 2) The best activation function was =ReLU and it gave the accuracy of 65%
- 3) Possible reason for misclassification is due to overfitting of model on training data.
Another reason could be the possible noise in the data as several samples had multiple numbers in the image.

Section - A

(a)



$$M = w_1 x_1 + w_2 x_2 + w_3 x_3 + b_1$$

Activation Function = ReLU

$$\text{Out } M_{01} = \text{Max}(0, M)$$

$$\lambda = 0.01 \quad (\text{Given})$$

$$\text{Let } w_1 = w_2 = w_3 = w_4 = w_5 = w_6 = 1, b_1 = b_2 = 0$$

$$x_1 = 1.2$$

$$x_2 = 0.8$$

$$x_3 = 2.0$$

$$T_1 \text{ target} = 2.0$$

$$T_2 \text{ target} = 2.5$$

$$T_3 \text{ target} = 4.0$$

Forward Pass

$$n_1 = \sum_{i=1}^3 w_i x_i + b_1$$

$$= 1 \cdot (1.2) + 1 \cdot (0.8) + 1 \cdot (2.0) + 2$$
$$= 6$$

$$\text{Out } n_1 = \text{Max}(0, 6) = 6$$

Now let's calculate y_1 ,

$$y_1 = \cancel{n_1 \cdot w_4} + b_2$$

$$= 6 \cdot 1 + 2$$

$$= 8$$

$$y_2 = n_1 \cdot w_5 + b_2$$

$$= 6 \cdot 1 + 2$$

$$= 8$$

$$y_3 = n_1 \cdot w_6 + b_2$$

$$= 6 \cdot 1 + 2$$

$$= 8$$

$$E_{\text{total}} = \sum_i \frac{1}{2} (t_i - y_i)^2$$

$$= \frac{1}{2} \left((8-3)^2 + (8-2.5)^2 + (8-4)^2 \right)$$

$$= \frac{1}{2} (25 + 30 + 25 + 6)$$

$$= 35 \quad (25)$$

$$\text{Error at } w_4 = \frac{\partial E_{\text{total}}}{\partial w_4}$$

~~$$= \frac{\partial E_{\text{total}}}{\partial w_4}$$~~

$$E_{\text{total}} = \frac{1}{2} (t_1 - \text{out } y_1)^2 + \frac{1}{2} \sum_i (T_i - \text{out } y_i)^2$$

$$\text{Error at } w_4 = \frac{\partial E_{\text{total}}}{\partial \text{out } y_1} \times \frac{\partial \text{out } y_1}{\partial y_1} \times \frac{\partial y_1}{\partial w_4}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out } y_i} = -(t_i - \text{out } y_i)$$

$$\frac{\partial \text{out } y_i}{\partial y_i} = \text{out } y_i - T_i$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out } y_1} = 8 - 3 = 5$$

$$\frac{\partial y_1}{\partial y_2} = 1$$

$$\frac{\partial y_1}{\partial w_4} = \text{out } H_1 = 6 \quad , \quad \frac{\partial y_2}{\partial w_5} = 0, \quad \frac{\partial y_2}{\partial w_6} = 5$$

$$\frac{\partial E_{\text{total}}}{\partial w_4} = (5) \times 1 \times 6$$

$$\partial w_4 = 30$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = (8 - 2 \cdot 5) \times 1 \times 6$$

$$\partial w_5 = 33$$

$$\frac{\partial E_{\text{total}}}{\partial w_6} = (8 - 4) \times 1 \times 6$$

$$\partial w_6 = 24$$

$$w_4 = w_4 - 2 \frac{\partial E_{\text{total}}}{\partial w_4}$$

$$= 1 - 0.01 \times 30$$

$$= 0.7$$

$$w_5 = 1 - 0.01 \times 33 = 0.67$$

$$w_g = 1 - 0.01 \times 24$$

$$= 0.76$$

~~$$\frac{\sum E_{total}}{\sum w_1} = \frac{\sum E_{total}}{\sum out\ n_1} \times \frac{\sum out\ n_1}{\sum n_1} \times \frac{\sum n_1}{\sum w_1}$$~~

$$\frac{\sum E_{total}}{\sum b_1} = \frac{\sum E_{total}}{\sum out\ y_1} \times \frac{\sum out\ y_1}{\sum y_1} \times \frac{\sum y_1}{\sum b_2}$$

$$+ \frac{\sum E_{total}}{\sum out\ y_2} \times \frac{\sum out\ y_2}{\sum y_2} \times \frac{\sum y_2}{\sum b_2}$$

$$+ \frac{\sum E_{total}}{\sum out\ y_3} \times \frac{\sum out\ y_3}{\sum y_3} \times \frac{\sum y_3}{\sum b_3}$$

$$= 5 \times 1 \times 1$$

$$+ 5.5 \times 1 \times 1$$

$$+ 9 \times 1 \times 1$$

$$= 14.5$$

$$b_2 = b_2 \times b_2 = 1 - 0.145 = 0.855$$

B)

(b) Here is the expression for soft-margin rule in this case.

$$y = \begin{cases} 1, & \text{if } \sum_{i=1}^n d_i y_i e^{\frac{-|x_i - x_j|^2}{2\sigma^2}} + b \geq 0 \\ -1, & \text{otherwise} \end{cases}$$

Derivation:

Let's say we have $\{x_i, y_i\}$, $i \in [0, n]$,
as our samples.

After solving optimization problem, we get,

$$w = \sum_{i=1}^n d_i y_i \phi(x_i)$$

When we do prediction,

$$y = \begin{cases} 1, & \text{if } w^T \phi(x_i) + b \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

The gaussian kernel is like as given below:

$$k(\phi(x_i), \phi(x_j))$$

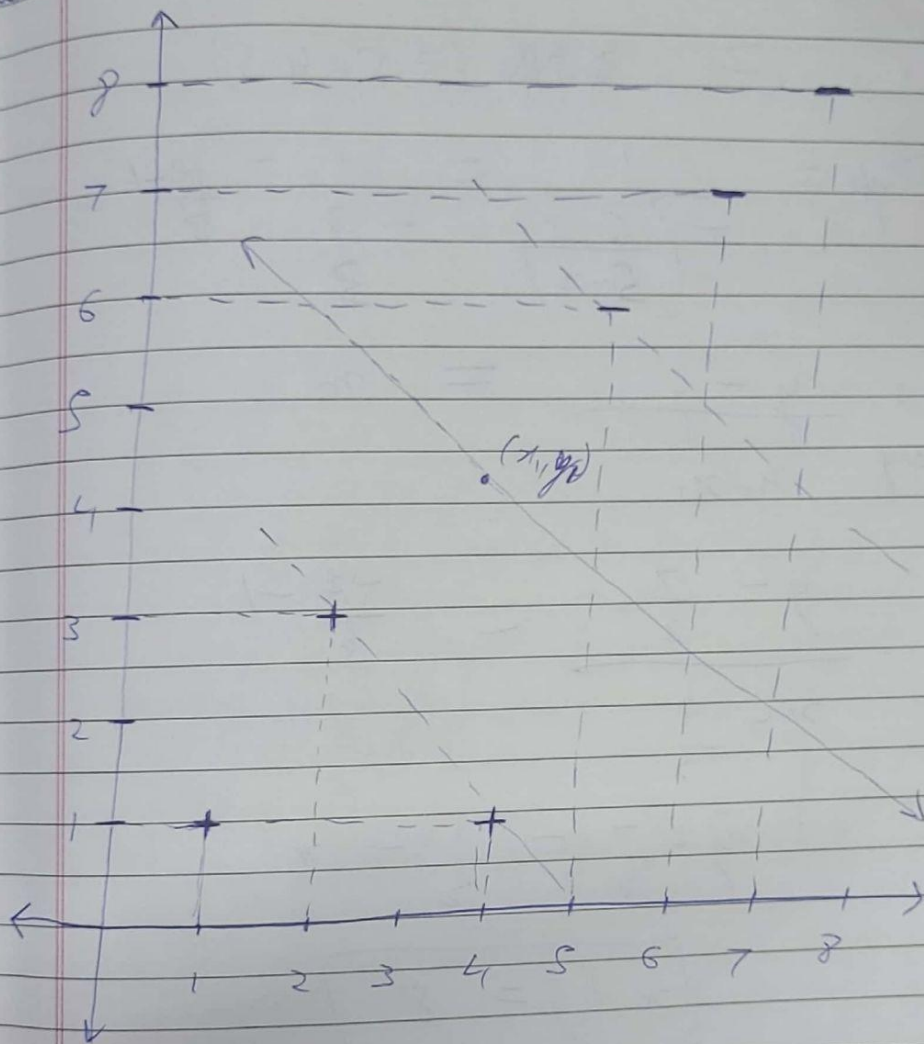
$$= k(x_i, x_j)$$

$$= e^{-\frac{|x_i - x_j|^2}{2\sigma^2}}$$

$$\therefore y = \begin{cases} 1, & \text{if } \sum \text{and } y_i \leq e^{-\frac{|x_i - x_j|^2}{2\sigma^2}} + b \\ -1, & \text{otherwise} \end{cases}$$

C)

(c)



Yes the data is linearly separable as we can clearly see from the plot

(d) We can determine that the plot that the decision boundary will be the perpendicular bisector of $(2, 3)$ and $(5, 6)$ as it will have ~~minimum~~ max. distance from them.

$$(x_1, y_1) = \left(\frac{2+5}{2}, \frac{3+6}{2} \right)$$

$$= (3.5, 4.5)$$

$$\text{Slope} = \frac{1}{-\left(\frac{6-3}{5-2}\right)} = \frac{1}{-\frac{3}{3}} = -1 = m$$

$$\text{Decision Boundary: } \frac{y - y_1}{x - x_1} = m$$

$$\frac{y - \frac{7}{2}}{x - \frac{9}{2}} = -1$$

$$y - \frac{7}{2} = \frac{9}{2} - x$$

$$x + y = 8$$

(ii) The support vectors are $(2, 3)$ and $(5, 6)$. They are the closest datapoints on either side.

$(4, 1)$, $(2, 3)$ lie in || to the decision boundary

(d) Margin = Distance between support vectors and hyperplane

= ~~for~~ distance between opposite support vectors.

= ~~for~~ distance b/w $(2, 3)$ and $(5, 6)$

$$= \sqrt{(5-2)^2 + (6-3)^2}$$

$$= \sqrt{3^2 + 3^2}$$

$$= 3\sqrt{2}$$

~~$$= 3 \times 1.414$$~~

~~$$= 3 \times 1.414 = 4.242$$~~

~~$$= 2.121$$~~

(a) If we remove any one of the support vectors, the optimal margin would change and we would need to minimize distance from new support vectors.

(a) If we remove either $(2, 3)$ and $(4, 1)$,
the optimal margin will not change.

If we remove $(5, 6)$ optimal margin will
change.