

## Assignment #2 – Theoretical Part

1. Given a neural network with two input neurons,  $x_1$  and  $x_2$ , three hidden neurons,  $h_1, h_2$  and  $h_3$ , and two output neurons,  $y_1$  and  $y_2$ . The network acts as a classifier (2 classes). The following equations govern the network operation:

$$h_1 = \sigma(x_1 + 1) \quad (1)$$

$$h_2 = \sigma(x_2 + 1) \quad (2)$$

$$h_3 = \sigma(1 - x_1 - 2x_2) \quad (3)$$

$$y_1 = \sigma(2.5 - h_1 - h_2 - h_3) \quad (4)$$

$$y_2 = \sigma(h_1 + h_2 + h_3 - 2.5) \quad (5)$$

where  $\sigma(x) = 1$  if  $x \geq 0$  and else  $\sigma(x) = 0$  otherwise.

Draw the decision region for each class and the decision boundary.

2. Consider a three layer neural network whose structure is shown in Figure 1. You are required to calculate the sensitivity  $\delta_k = -\frac{\partial J}{\partial \text{net}_k}$  at the output node  $k$ , where  $J$  is the objective function to be minimized and  $\text{net}_k$  is the net activation of the output node  $k$ . We consider two cases, where the objective function  $J$  and the nonlinear activation function at the output layer are chosen differently.
- In the first case,  $J$  is chosen as the squared error  $J(W) = \frac{1}{2} \|t - z\|^2$  where  $z_k$  is the prediction at the output node  $k$  and  $t_k$  is the corresponding target value. In the classification problem, only one  $t_k$  equals to 1 (corresponding the ground truth class) and all the other  $t_k$ 's are all zeros. Sigmoid  $f(\text{net}_k) = 1/(1 + e^{-\text{net}_k})$  is chosen as the activation function at the output layer. Calculate the sensitivity  $\delta_k$  in terms  $t_k, z_k$  and  $\text{net}_k$ . Show that all the  $\delta_k$  could be close to zero even if the prediction error is large and explain why this is bad.
  - In the second case, the objective function is chosen as cross entropy,  $J(W) = -\sum_{k=1}^C t_k \log(z_k)$  and the nonlinear activation function at the output layer is chosen as softmax  $f(\text{net}_k) = e^{\text{net}_k} / \sum_{j=1}^C e^{\text{net}_j}$ . Calculate the sensitivity  $\delta_k$  again. Prove that if the prediction error is large, at least one of the  $\delta_k$  will be large.

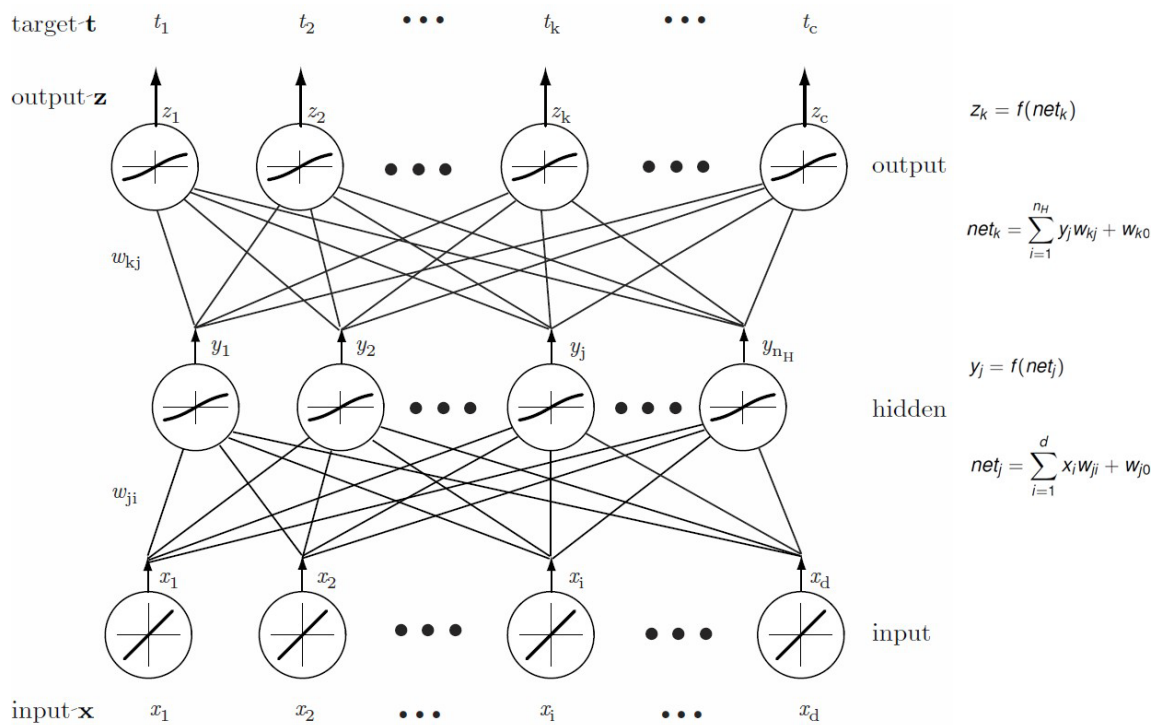


Figure 1

$$1. \quad h_1 = \sigma(x_1 + 1)$$

$$h_2 = \sigma(x_2 + 1)$$

$$h_3 = \sigma(1 - x_1, -2x_2)$$

$$y_1 = 5(2.5 - h_1 - h_2 - h_3)$$

$$y_2 = \sigma(h_1 + h_2 + h_3 - 2.5)$$

$$\Rightarrow h_1 = \begin{cases} 1 & \text{if } x_1 \geq -1 \\ 0 & \text{if } x_1 < -1 \end{cases}$$

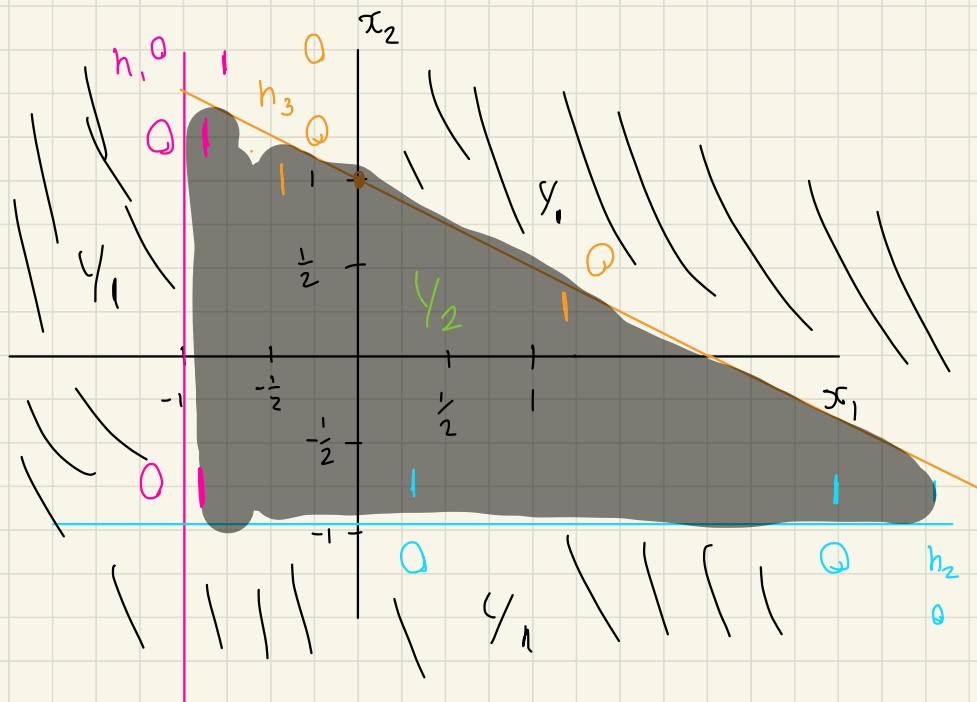
$$h_2 = \begin{cases} 1 & \text{if } x_2 \geq -1 \\ 0 & \text{if } x_2 < -1 \end{cases}$$


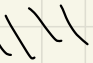
$$h_3 = \begin{cases} 1 & | \geq x_1 + 2x_2 \\ 0 & | < x_1 + 2x_2 \end{cases}$$

$$y_i = \begin{cases} 1 & 2.5 \geq h_1 + h_2 + h_3 \\ 0 & 2.5 < h_1 + h_2 + h_3 \end{cases}$$

$$y_2 = \begin{cases} 1 & n_1 + h_2 + h_3 \geq 2.5 \\ 0 & n_1 + h_2 + h_3 < 2.5 \end{cases}$$

$$x_2 = -\frac{1}{2}x_1 + 1$$



We see that the  shaded region is the region in which  $h_1 + h_2 + h_3 \geq 2.5$  is of class 2 (or  $\gamma_2$ ). The  shaded region outside of the triangle is of class 1 (or  $\gamma_1$ ).

$$2. a) \quad J(w) = \frac{1}{2} \|t - z\|^2$$

where  $z_k$  = prediction at output node  $k$

$t_k$  = desired value at node  $k$

$$f(\text{net}_k) = \frac{1}{1 + e^{-\text{net}_k}} \quad (\text{recall } \sigma'(x) = \sigma(x)(1 - \sigma(x)))$$

$$\begin{aligned} \delta_k &= -\frac{\partial J}{\partial \text{net}_k} = -\frac{\partial}{\partial \text{net}_k} \left( \frac{1}{2} \|t_k - z_k\|^2 \right) \\ &= -\frac{\partial}{\partial \text{net}_k} \left( \frac{1}{2} \|t_k - f(\text{net}_k)\|^2 \right) \end{aligned}$$

$$\Rightarrow -\frac{\partial J}{\partial \text{net}_k} = -\frac{\partial J}{\partial z_k} \cdot \frac{\partial z_k}{\partial \text{net}_k} = -\|t_k - z_k\| \cdot z_k(1 - z_k)$$

For  $t_k = 0$

$$\delta_k = -z_k^2(1 - z_k) = z_k^3 - z_k^2 = \sigma(n_k)^3 - \sigma(n_k)^2$$

Now to find max sensitivity we solve for  $\delta'_k = 0$

$$\frac{d}{dx} (\sigma(x) \sigma(x))$$

$$= 2 \sigma(x) \sigma(x) (1 - \sigma(x))$$

$$= 2 \sigma(x)^2 (1 - \sigma(x))$$

$$\frac{d}{dx} (\sigma(x) \sigma(x) \sigma(x))$$

$$= 3 \sigma(x)^3 (1 - \sigma(x))$$

$$\Rightarrow \delta'_k = -2z_k^2(1-z_k) + 3z_k^3(1-z_k) = 0$$

$$\Rightarrow 2z_k^2 = 3z_k^3 \Rightarrow$$

$$\boxed{\frac{2}{3} = z_k \text{ maximum value}}$$

$$\delta_k = -\left(\frac{2}{3}\right)^2 \left(\frac{1}{3}\right) = \frac{4}{27} = \underline{0.148} \Rightarrow \text{sensitivity of } -0.148 \text{ is the greatest magnitude @ } t_k = 0.$$

$$\boxed{t_k = 1}$$

$$\delta_k = -z_k(1-z_k)^2 = z_k(1-2z_k+z_k^2) = z_k - 2z_k^2 + z_k^3$$

$$\delta'_k = -z_k(1-z_k) + 4z_k^2(1-z_k) - 3z_k^3(1-z_k)$$

$$\underline{\text{solve for } \delta'_k = 0}$$

$$z_k - 4z_k^2 + 3z_k^3 = 0$$

$$\Rightarrow 3z_k^2 - 4z_k + 1 = (1-3z_k)(1-z_k)$$

$$\Rightarrow z_k = 1 \text{ or } z_k = \frac{1}{3}$$

$$\delta_k(1) = 0 \text{ (hence } z_k = 1 \text{ is a minimum)} \quad \delta_k\left(\frac{1}{3}\right) = -\frac{2}{9} = -0.22$$

so we see that the greatest magnitude of  $\delta_k$  is 0.22

For a large prediction error we see:

$$t_k = 1, z_k \approx 0: \delta_k = 0 \cdot (1-0)^2 \approx 0$$

$$t_k = 0, z_k \approx 1: \delta_k = 1^2(1-1) \approx 0$$

i.e.  $\delta_k \rightarrow 0$  for large prediction error

Hence for large prediction error we get sensitivity,  $\delta_k$ , close to zero.

This is bad because when applying backpropagation to a neural network, if our gradient vanishes then our weight updates will essentially stop which in turn causes the network to stop learning.

$$6) \quad J(W) = \sum_{k=1}^C t_k \log(z_k) \quad f(\text{net}_k) = \frac{e^{\text{net}_k}}{\sum_{j=1}^C e^{\text{net}_j}}$$

derivative of softmax: ( $x_k = \text{net}_k$ ) ( $j \neq k$ )

$$g(x_k) = e^{x_k} + e^{x_1} + \dots + e^{x_C} \quad f(x_k) = e^{x_k} \quad g(x_i) = e^{x_k} + e^{x_1} + \dots + e^{x_C} \quad f(x_i) = e^{x_i}$$

$$g'(x_k) = e^{x_k} \quad f'(x_k) = e^{x_k} \quad g'(x_i) = e^{x_k} \quad f'(x_i) = 0$$

$$\Rightarrow \frac{\partial z_k}{\partial x_k} = \frac{e^{x_k} \sum_{j=1}^C e^{x_j} - e^{x_k} \cdot e^{x_k}}{\left( \sum_{j=1}^C e^{x_j} \right)^2}$$

$$= \frac{e^{x_k}}{\sum_{j=1}^C e^{x_j}} \cdot \frac{\sum_{j=1}^C e^{x_j} - e^{x_k}}{\sum_{j=1}^C e^{x_j}}$$

$$= z_k \cdot (1 - z_k)$$

$$\Rightarrow \frac{\partial z_i}{\partial x_k} = \frac{0 - e^{x_i} e^{x_k}}{\left( \sum_{j=1}^C e^{x_j} \right)^2} = -z_i z_k$$

$$= - \frac{\partial J}{\partial \text{net}_k} = - \frac{\partial J}{\partial z_i} \cdot \frac{\partial z_i}{\partial \text{net}_k} = \sum_{i=1}^C \frac{t_i}{z_i} \cdot \frac{\partial z_i}{\partial \text{net}_k}$$

$$= - \sum_{i \in [1, C] \setminus \{i\}} t_i \cdot z_k + t_k (1 - z_k)$$

$$\text{For } z_k \approx 0, t_k = 1 \Rightarrow \delta_k = t_k = 1$$

$$z_k \approx 1, t_k = 0 \Rightarrow \delta_k = \sum_{i=1}^C t_i \approx 1$$

This shows us that for large prediction error we get a large  $\delta_k$