

Let x be the input with f features.

We have two layers.

Hidden layer: w_1 weight, b_1 bias

Output layer: w_2 weight, b_2 bias

We have,

$$z_1 = w_1 x + b_1$$

$$a_1 = \sigma(z_1) \quad [\text{Sigmoid function}]$$

$$z_2 = w_2 a_1 + b_2$$

$$\hat{y} = z_2$$

We have linear function for the output because this is the appropriate activation function for continuous output (regression).

$$L = \text{MSE} = \frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2$$

$$\frac{\partial L}{\partial z_2} = \frac{\partial L}{\partial \hat{y}} = (\hat{y} - y)$$

$$\begin{aligned}\frac{\partial L}{\partial w_2} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial w_2} \\ &= (\hat{y} - y) \cdot a_1\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial b_2} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial b_2} \\ &= (\hat{y} - y) \cdot 1 \\ &= (\hat{y} - y)\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial z_1} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \\ &= (\hat{y} - y) \cdot w_2 \cdot \sigma'(z_1)\end{aligned}$$

where, $\sigma'(z_1)$ is the derivative of $\sigma(z_1)$ function,

$$\sigma'(z_1) = \sigma(z_1) (1 - \sigma(z_1))$$

$$\therefore \frac{\partial L}{\partial z_1} = (\hat{y} - y) w_2 \cdot \sigma(z_1) (1 - \sigma(z_1))$$

$$\begin{aligned}\frac{\partial L}{\partial w_1} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_1} \\ &= (\hat{y} - y) \cdot w_2 \cdot x \cdot \sigma'(z_1)\end{aligned}$$

$$\begin{aligned}\frac{\partial L}{\partial b_1} &= \frac{\partial L}{\partial z_2} \cdot \frac{\partial z_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial b_1} \\ &= (\hat{y} - y) \cdot w_2 \cdot \sigma'(z_1)\end{aligned}$$

Update rules

$$\begin{aligned}w_1 &= w_1 - \alpha \frac{\partial L}{\partial w_1} \\ &= w_1 - \alpha \cdot (\hat{y} - y) \cdot w_2 \cdot x \cdot \sigma'(z_1)\end{aligned}$$

$$\begin{aligned}b_1 &= b_1 - \alpha \frac{\partial L}{\partial b_1} \\ &= b_1 - \alpha \cdot (\hat{y} - y) \cdot w_2 \cdot \sigma'(z_1)\end{aligned}$$

$$\begin{aligned}w_2 &= w_2 - \alpha \frac{\partial L}{\partial w_2} \\ &= w_2 - \alpha \cdot (\hat{y} - y) \cdot a_1\end{aligned}$$

$$b_2 = b_2 - \alpha \frac{\partial L}{\partial b_2}$$

$$= b_2 - \alpha (\hat{y} - y)$$

Using these rules, we will perform backpropagation and train the neural network.

This is different from binary classification, because we are using identity function in the output to ensure regression instead of any other activation function that works for discrete classification.