

ANN - Artificial Neural Networks for regression and classification



gradient descent is good for training data
but on new data, learns specific
noise and trends that may not
be general to all data.

$$f(w_1x_1 + w_2x_2 + b) = y \quad \text{feed forward equation}$$

↑
activation function

* carried out using backpropagation

Backpropagation : propagate error backwards

adjusts the weights and bias by measuring the error b/w the true output and output. The gradient tells us how the error would change, and we use gradient descent, moving in the direction that reduces error most.

* allows nonlinear mapping b/w inputs and outputs

Learning Rate : step size of optimisation algorithm.

too high → oscillate around minimum, overshoot

too low → get stuck in local optimum, never converge.

LAB 1

$$\text{sigmoid } f : \frac{1}{1+e^{-(w^T x + b)}} = \sigma$$

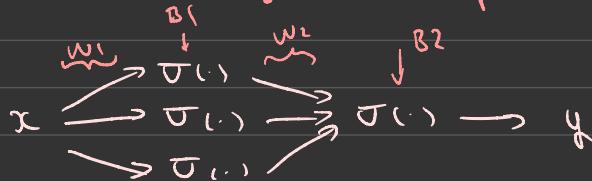
$$\text{squared error: } E = (\hat{y} - o)^2 \quad \begin{matrix} \hat{y} & \text{true output} \\ o & \text{output} \end{matrix} \quad x = w^T x + b$$

gradient descent using chain rule

$$\left\{ \begin{array}{l} \frac{dE}{d\sigma} = -2(\hat{y} - o) \quad \frac{d\sigma}{dx} = \sigma(1-\sigma) \quad \frac{dx}{dw} = x \quad \frac{dx}{db} = 1 \\ \frac{dE}{dw} = \frac{dE}{d\sigma} \times \frac{d\sigma}{dx} \times \frac{dx}{dw} \quad \frac{dE}{db} = \frac{dE}{d\sigma} \times \frac{d\sigma}{dx} \times \frac{dx}{db} \end{array} \right.$$

multiplied by learning rate

LAB 2 - adding hidden layer



Different dimensions? depends on the # of nodes.
first layer (hidden) is $(3, 2)$ because it connects 3 sigmoid nodes for 2 inputs (1 or 0)
second layer is $(1, 3)$ as it is reduced to 1 sigmoid node for the 3 inputs (outputs from prior layer)

Backpropagation Derivation

W2 : $w = w + \alpha \frac{dE}{dw}$, B2, Both same as lab 1

$$W1 : \frac{dL}{dw_1} = \frac{dL}{dH} \times \frac{dH}{dz} \times \frac{dz}{dw}, \quad \frac{dL}{dH} = \text{Output dot}(W2)$$

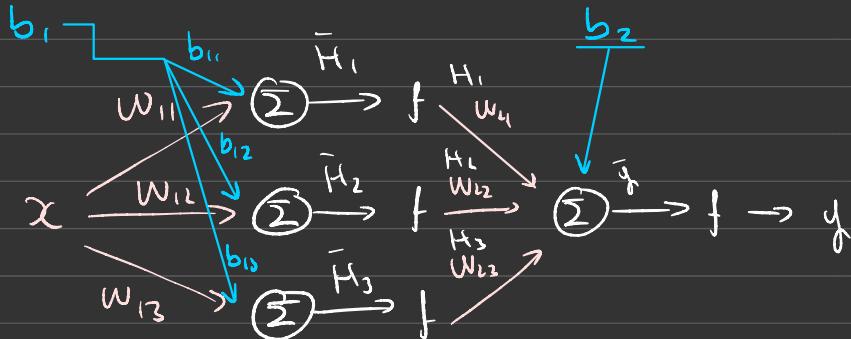
$$\frac{dL}{db_1} = \frac{dL}{dH} \cdot \frac{dH}{\partial z} \cdot \frac{dz}{db_1}$$

Sigmoid function : Classification Problem



$\rightarrow o_i = f_i(s) = \frac{e^{s_i}}{\sum\limits_{j=1}^k e^{s_j}}$ rescales output to range [0, 1] and sets sum=1, so output is a probability.

- for softmax, use cross-entropy loss for error function
- Measures performance at classification, a perfect model would have a loss of 0.
 - for backpropagation: $\frac{dH}{ds_i} = o_i - z_i$
- $$H = -\sum_{j=1}^n z_j \log(o_j)$$



① work backwards and find diff-BL & diff-BI

$$y = f(\bar{y}) \quad \bar{y} = \sum w_i H_i + b_i$$

$$\therefore \bar{y} = w_{11}H_1 + w_{12}H_2 + w_{13}H_3 + b_1$$

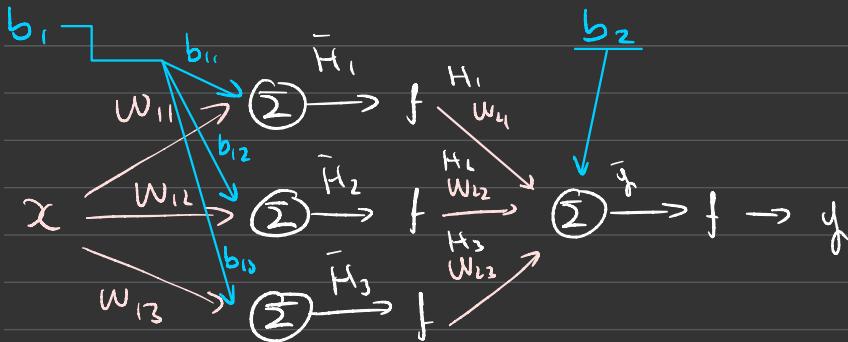
$$E = -\frac{1}{2}(\hat{y} - y)^2 \quad \frac{dE}{dy} = \hat{y} - y \quad * \hat{y} = \text{real } y = \text{calculated}$$

$$\textcircled{2} \quad \frac{dE}{dw_2} = \begin{bmatrix} \frac{dE}{dy} & \frac{dy}{d\bar{y}} & \frac{d\bar{y}}{dw_2} \end{bmatrix} = (\hat{y} - y) \underbrace{(f'(\bar{y}))}_{\text{if } f \text{ is sigmoid}} (H)$$

it f is sigmoid

$$f'(\bar{y}) = \bar{y}(1 - \bar{y})$$

$$\textcircled{3} \quad \frac{dE}{db_1} = \begin{bmatrix} \frac{dE}{dy} & \frac{dy}{d\bar{y}} & \frac{d\bar{y}}{dw_2} \end{bmatrix} = (\hat{y} - y) (f'(\bar{y})) (1)$$



① work backwards and find diff_BL & int-B1

$$y = f(\bar{y}) \quad \bar{y} = \sum w_i H_i + b_i$$

$$\therefore \bar{y} = w_{11}H_1 + w_{12}H_2 + w_{13}H_3 + b_2$$

$$E = -\frac{1}{2}(\hat{y} - y)^2 \quad \frac{dE}{dy} = 0 \quad , \quad * \hat{y} = \text{real } y = \text{calculated}$$

$$\textcircled{2} \quad \frac{dE}{dw_2} = \left[\frac{dE}{d\hat{y}} \cdot \frac{d\hat{y}}{d\bar{y}} \cdot \frac{d\bar{y}}{dw_2} \right] = (\hat{y} - y) \underbrace{(f'(\bar{y}))}_{\text{it } f \text{ is sigmoid}} (H)$$

$$\textcircled{3} \quad \frac{dE}{db_2} = \left[\frac{dE}{d\hat{y}} \cdot \frac{d\hat{y}}{d\bar{y}} \cdot \frac{d\bar{y}}{db_2} \right] = (\hat{y} - y) (f'(\bar{y})) (1)$$

④ find $\frac{dE}{dw_1}$ and $\frac{dE}{db_1}$

$$= f(\bar{H}) \quad \bar{H} = \sum w_i X + b_i \quad E = -\frac{1}{2}(\hat{H} - H)^2$$

$$\frac{dE}{dH} = \hat{H} - H$$

$$\frac{dE}{dw_1} = \frac{dE}{dH} \cdot \frac{dH}{d\bar{H}} \cdot \frac{d\bar{H}}{dw_1} = \underbrace{\frac{dE}{dH}}_{\hat{y} - y} \cdot f'(\bar{H}) \times X = (\hat{y} - y)(y(1-y))w_2(H(1-H))X$$

$$\frac{dE}{db_1} = \frac{dE}{dH} \cdot \frac{dH}{d\bar{y}} \cdot \frac{d\bar{y}}{db_1} = (\hat{y} - y) \cdot y(1-y) \cdot w_2$$

$\frac{dE}{db_1}$ is the same w/o w_2