

Feed-forward neural network

binary classification

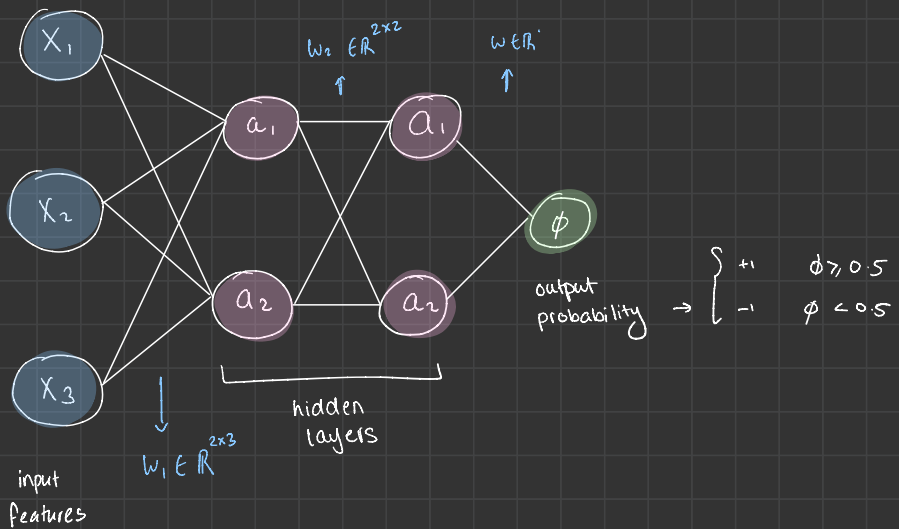
$X \in \mathbb{R}^{3 \times 1}$ is a datapoint w/ 3 features: $X = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$

$W \in \mathbb{R}^{1 \times 3}$ is a weight vector representing node connections

$b \in \mathbb{R}$ is a bias term

↑ activation function

$$a = f(WX + b) \in \mathbb{R}$$



$$X = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

$$W_1 = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \end{bmatrix}$$

$$b_1 = \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$$

$$W_2 = \begin{bmatrix} 0.7 & 0.8 \\ 0.9 & 1 \end{bmatrix}$$

$$b_2 = \begin{bmatrix} 0.3 \\ 0.4 \end{bmatrix}$$

$$W_3 = \begin{bmatrix} 1.1 & 1.2 \end{bmatrix}$$

$$b_3 = 0.5$$

$$f(z) = \sigma(z) = \frac{1}{1 + e^{-z}} \rightarrow \text{sigmoid activation}$$

hidden layer #1

P.S. i did not bother calculating, GPT did.
so if theres any mistakes, not my fault fr.

$$W_1 x + b_1 = \begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.5 & 0.6 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix}$$

$$= \begin{bmatrix} 1.4 \\ 3.2 \end{bmatrix} + \begin{bmatrix} 0.1 \\ 0.2 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 3.4 \end{bmatrix}$$

$$f(W_1 x + b_1) = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sigma(1.5) \\ \sigma(3.4) \end{bmatrix} = \begin{bmatrix} 0.818 \\ 0.968 \end{bmatrix}$$

hidden layer #2

$$\text{"new"} \quad x = \begin{bmatrix} 0.818 \\ 0.968 \end{bmatrix}$$

$$W_2 x + b_2 = \begin{bmatrix} 0.7 & 0.8 \\ 0.9 & 1 \end{bmatrix} \begin{bmatrix} 0.818 \\ 0.968 \end{bmatrix} + \begin{bmatrix} 0.3 \\ 0.4 \end{bmatrix}$$

$$= \begin{bmatrix} 2.226 \\ 2.734 \end{bmatrix}$$

$$f(w_2 x + b_2) = \begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} \sigma(2.226) \\ \sigma(2.434) \end{bmatrix} = \begin{bmatrix} 0.903 \\ 0.939 \end{bmatrix}$$

output layer

$$w_3 x + b_3 = \begin{bmatrix} 1.1 & 1.2 \end{bmatrix} \begin{bmatrix} 0.903 \\ 0.939 \end{bmatrix} + 0.5$$

$$= 2.1145 + 0.5$$

$$= 2.6195$$

$$f(w_3 x + b_3) = \hat{y} = \sigma(2.6195) = 0.932 \geq 0.5$$

$$\Rightarrow y_{\text{pred}} = +1$$

back propagation

goal: $\min_w L(y, \hat{y})$

loss function \uparrow

w \downarrow minimize w.r.t. w

y \downarrow actual label

\hat{y} \downarrow predicted label

• stochastic gradient descent

update rules \swarrow

$$w_j := w_j - \alpha \frac{\partial}{\partial w_j} L(y, \hat{y})$$

$$b_j := b_j - \alpha \frac{\partial}{\partial b_j} L(y, \hat{y})$$

\nwarrow learning rate

$\nearrow \hat{y} = f(wx + b)$

assume we're working w/ a regression problem:

$$L(y, \hat{y}) = \frac{1}{2n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

\downarrow number of samples \downarrow true \downarrow predicted

let's just do this for one datapoint:

$$L(y, \hat{y}) = \frac{1}{2} (y - f(z))^2 \quad \hat{y} = f(z) = e^z \quad z = wx + b$$

$$\Rightarrow \frac{dL}{dw} = \frac{dL}{df(z)} \cdot \frac{df(z)}{dz} \cdot \frac{dz}{dw}$$

\downarrow activation function \downarrow $wx + b$

$$\frac{dL}{df(z)} = (y - f(z)) \quad \frac{df(z)}{dz} = e^z \quad \frac{dz}{dw} = x$$

$$\Rightarrow \frac{dL}{dw} = (y - f(z)) \cdot e^z \cdot x$$

$$= (y - (wx + b)) \cdot e^{wx + b} \cdot x$$

\downarrow layer input, whatever it is

• SGD with momentum: $w_t := w_t - \alpha \overset{\text{velocity}}{\underset{\uparrow}{V_t}} \rightarrow$ accumulates prev. gradients

$$V_t = \beta \underset{\substack{\downarrow \\ \text{velocity at} \\ \text{previous time}}}{V_{t-1}} + (1 - \beta) \cdot \underset{\substack{\downarrow \\ \text{momentum}}}{\frac{dL}{dw_{t-1}}}$$

• Adam \rightarrow uses momentum & adaptive learning rate