

a.

```
X = np.array([[1, x1, x2]])
W1 = np.array([[w1, w2], [w3, w4], [w5, w6]])
W2 = np.array([[v1], [v2], [v3]])
```

* $g(x) = x$

* $g(X @ W1) = np.array([z1, z2])$

* $g(g(X @ W1) @ W2) = g(w1 * 1 + v2 z1 + v3 z2)$

$$\begin{bmatrix} 1 & z1 & z2 \end{bmatrix}_{1 \times 3} @ \begin{bmatrix} v1 \\ v2 \\ v3 \end{bmatrix}_{3 \times 1} = \begin{bmatrix} v1 * 1 + v2 z1 + v3 z2 \end{bmatrix}_{1 \times 1}$$

$$X @ W1: \begin{bmatrix} 1 & x1 & x2 \end{bmatrix}_{1 \times 3} @ \begin{bmatrix} w1 & w2 \\ w3 & w4 \\ w5 & w6 \end{bmatrix}_{3 \times 2} = \begin{bmatrix} w1 + x1w3 + x2w5, w2 + x1w4 + x2w6 \end{bmatrix}_{1 \times 2}$$

$(X @ W1) @ W2:$

$$\begin{bmatrix} w1 + x1w3 + x2w5, w2 + x1w4 + x2w6 \end{bmatrix}_{1 \times 2} @ \begin{bmatrix} v1 \\ v2 \\ v3 \end{bmatrix}_{3 \times 1} \quad \text{* intercept term!}$$

$$= \begin{bmatrix} w1v1 + w3x1v1 + w5x2v1 + w2v2 + w4x1v2 + w6x2v2 \end{bmatrix}$$

$$\rightarrow X @ W1: \begin{bmatrix} 1 & z1 & z2 \end{bmatrix}_{1 \times 3} \dots \begin{bmatrix} w1 + x1w3 + x2w5, w2 + x1w4 + x2w6 \end{bmatrix}_{1 \times 2}$$

$1 \rightarrow 1$ intercept term

$z1 \rightarrow w1 + x1w3 + x2w5$

$z2 \rightarrow w2 + x1w4 + x2w6$

$$\rightarrow g(v1 * 1 + v2 z1 + v3 z2) = \begin{bmatrix} v1 * 1 + v2 z1 + v3 z2 \end{bmatrix}$$

$$= \begin{bmatrix} v1 + v2w1 + v2x1w3 + v2x2w5 + v3w2 + v3x1w4 + v3x2w6 \end{bmatrix}$$

b.

Handwritten derivation showing the simplification of a linear function $g(x) = x$ into a form $g(m + n \cdot x_1 + o \cdot x_2)$.

Initial expression:

$$g(x) = x$$

$$g(m + n \cdot x_1 + o \cdot x_2)$$

Intermediate expression (sum of products):

$$[v_1 + v_2 w_1 + v_2 x_1 w_3 + v_2 x_2 w_5 + v_3 w_2 + v_3 x_1 w_4 + v_3 x_2 w_6]$$

Simplified expression (grouped terms):

$$\rightarrow [v_1 + v_2 w_1 + v_3 w_2 + x_1 (v_2 w_3 + v_3 w_4) + x_2 (v_2 w_5 + v_3 w_6)]$$

Final simplified coefficients:

$$\therefore \begin{aligned} m &= v_1 + v_2 w_1 + v_3 w_2 \\ n &= v_2 w_3 + v_3 w_4 \\ o &= v_2 w_5 + v_3 w_6 \end{aligned}$$

c.

This shows that an MLP using linear activation functions is equivalent to a linear perceptron because the entire network can be expressed as a linear function. Though, it should be stated that using linear activation functions limits the expressive power of a neural network as it can only learn linear decision boundaries, while an MLP with one or more hidden layers and linear activation functions can approximate any continuous function (of course, given enough hidden units). In this case, using the linear activation function $g(x) = x$ allows one to express the entire network as a linear function, which is equivalent to a linear perceptron. However, this approach, as briefly mentioned, limits the network's ability to model complex relationships between features and outputs. Any other complex function should be expressed by an MLP.

d.

$$g(x) = \frac{1}{1+e^{-x}} ; g(m \times 1 + n \times 1 + 0 \times 2) = \frac{1}{1+e^{-m \times 1 - n \times 1 - 0 \times 2}}$$

$$x @ w1 = [w1 + x1w3 + x2w5 \quad w2 + x1w4 + x2w6]_{1 \times 2}$$

$$g(x @ w1) = \frac{1}{1+e^{-(w1 + x1w3 + x2w5 - w2 - x1w4 - x2w6)}} = \frac{1}{1+e^{-(x @ w1)}}$$

$$g(g(x @ w1) @ w2) = \frac{1}{1+e^{-\left(\frac{1}{1+e^{-(x @ w1)}}\right) \cdot w2}}$$

$$= \frac{1}{1+e^{-\left(\frac{1}{1+e^{-(w1 + x1w3 + x2w5 - w2 - x1w4 - x2w6)}}\right) \cdot \begin{bmatrix} v1 \\ v2 \\ v3 \end{bmatrix}}}$$

$$= \frac{1}{1+e^{-\left(\frac{v1}{1+e^{-(x @ w1)}} - \frac{v2}{1+e^{-(x @ w1)}} - \frac{v3}{1+e^{-(x @ w1)}}\right)}}$$

$$\therefore g(g(x @ w1) @ w2) = \frac{1}{1+e^{-m \times 1 - n \times 1 - 0 \times 2}}$$

$$-\frac{v1}{1+e^{-(x @ w1)}} - \frac{v2}{1+e^{-(x @ w1)}} - \frac{v3}{1+e^{-(x @ w1)}} = -m \times 1 - n \times 1 - 0 \times 2$$

$$\frac{-(v1 + v2 + v3)}{1+e^{-(x @ w1)}} = -m \times 1 - n \times 1 - 0 \times 2 \leftarrow$$

$$* e^A = 1 + A + \frac{1}{2!} A^2 + \frac{1}{3!} A^3 + \dots + \sum_n \frac{1}{n!} A^n \leftarrow$$

* computationally expensive!

Yes, we'd need to use numerical methods to solve for m , n , and o .

e.

This now shows that an MLP using sigmoid activation functions is more expressive than a single sigmoid perceptron because the non-linear activation function allows for more complex mappings between the input features and the output. Using such an activation function allows the MLP to learn nonlinear decision boundaries and thus model more complex relationships within the data, rather than what a single sigmoid perceptron can do.

f.

As discussed prior, non-linearity is crucial for an activation function characteristic. This is an important property as the model must capture complex relationships between the input features and the output within the data. Without the nonlinear property, a neural network is intensely limited to learning only linear relationships, which cannot capture the true nature of relationships in many cases. Another important feature is differentiability. The activation function must be differentiable, as the backpropagation algorithm in neural networks uses the chain-rule, which requires differentiation, to train neural networks. This then allows one to compute gradients of the loss with respect to the network parameters and update them using gradient descent.