

$$X = \text{np.array}([1, x_1, x_2]) \quad x = [1, x_1, x_2]$$

$$w_1 = \text{np.array}([w_1, w_2], [w_3, w_4], [w_5, w_6]) \quad \begin{bmatrix} w_1, w_2 \\ w_3, w_4 \\ w_5, w_6 \end{bmatrix}$$

$$w_2 = \text{np.array}([v_1], [v_2], [v_3]) \quad \begin{bmatrix} v_1 \\ v_2 \\ v_3 \end{bmatrix}$$

If we denote $g(x \cdot w_1) = \text{np.array}([1, z_1, z_2])$, then
 $g(g(x \cdot w_1), w_2)$ can be written: $g(v_1 \cdot 1 + v_2 \cdot z_1 + v_3 \cdot z_2)$
 Let $g(x)$ be a linear activation fn s.t. $g(x) = x$

a. $g(v_1 \cdot 1 + v_2 \cdot z_1 + v_3 \cdot z_2)$
 $= g(g(x \cdot w_1), w_2) \quad x \cdot w_1 = [w_1, w_2] + [w_3 x_1, w_4 x_1] + [w_5, w_6] x_2$
 $x \cdot w_1 = [w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6]$
 $\therefore g(v_1 \cdot 1 + v_2 \cdot z_1 + v_3 \cdot z_2) = g(g([w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6], w_2))$

b. Assuming $g(x) = x$, we can simplify our expression from above:
 start with: $g(g([w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6], w_2))$

$$\rightarrow g([w_1, w_2] + [w_3, w_4] x_1 + [w_5, w_6] x_2, w_2)$$

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

Clearly we have written $g(v_1 \cdot 1 + v_2 \cdot z_1 + v_3 \cdot z_2)$ as
 $g(m \cdot 1 + n \cdot x_1 + o \cdot x_2)$ where $m = v_1 [w_1, w_2]$,
 $n = v_2 [w_3, w_4]$, $o = v_3 [w_5, w_6]$

c.

Since the result of this MLP using linear activation functions can be written in terms of the original input (only with modified weights), our MLP isn't able to learn anything more complicated than a single perceptron (which would do the exact same thing of assigning weights to each of the input features).

d. Now let $g(x) = (1 + \exp(-x))^{-1}$
 start with: $g(g([w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6], w_2))$

$$\rightarrow g\left(\frac{1}{1 + \exp(-[w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6] \cdot w_2)}\right)$$

$$\rightarrow g\left(\frac{1}{1 + \exp(-[w_1 + x_1 w_3 + x_2 w_5, w_2 + x_1 w_4 + x_2 w_6])} \cdot w_2\right)$$

Clearly the exponential term cannot be written as a linear combination of 1, x_1 , and x_2 (due to the sigmoid function being nonlinear itself), so we cannot write its dot product with w_2 as a linear combination of 1, x_1 , and x_2 either.

e.

Since a combination of sigmoid perceptrons cannot be reduced to a single sigmoid perceptron with altered weights, a MLP using sigmoid activation functions can learn more complicated decision boundaries than a single sigmoid perceptron.

f.

An activation must be nonlinear and differentiable in order to be useful in a neural network. Nonlinearity allows the network to learn arbitrarily complicated decision boundaries, while differentiability allows backpropagation through the network to work properly by allowing the use of gradient descent to update the weights of each node.