

Activation-Based Saliency Maps: Two-Neuron Example

We extend the single-neuron example to a network with two neurons in the hidden layer.

Parameters:

- **Hidden Neuron 1 (h1):** Weights: $w_{11} = 2$, $w_{12} = -1$, $w_{13} = 0.5$, Bias: $b_1 = 0.5$
- **Hidden Neuron 2 (h2):** Weights: $w_{21} = -1$, $w_{22} = 1.5$, $w_{23} = -0.5$, Bias: $b_2 = -0.5$
- **Output Neuron:** Weights: $w_{o1} = 1$, $w_{o2} = -2$, Bias: $b_o = 0.5$
- **Input Values:** $x_1 = 1$, $x_2 = 2$, $x_3 = -1$

Step 1: Compute Activations of Hidden Neurons

$$h_1 = \text{sigmoid}(2 \cdot 1 + (-1) \cdot 2 + 0.5 \cdot (-1) + 0.5) = \text{sigmoid}(0) = 0.5$$

$$h_2 = \text{sigmoid}(-1 \cdot 1 + 1.5 \cdot 2 + (-0.5) \cdot (-1) - 0.5) = \text{sigmoid}(2) \text{ approximately } 0.88$$

Step 2: Compute Output

$$y = \text{sigmoid}(1 \cdot h_1 + (-2) \cdot h_2 + 0.5) = \text{sigmoid}(0.5 - 1.76 + 0.5) = \text{sigmoid}(-0.76) \text{ approximately } 0.32$$

Step 3: Compute Gradients (Saliency Map)

$$dy/dz_o = 0.32 \cdot (1 - 0.32) = 0.217$$

Gradients w.r.t. hidden neurons:

$$- dy/dh_1 = 0.217 \cdot 1 = 0.217$$

$$- dy/dh_2 = 0.217 \cdot (-2) = -0.434$$

Gradients w.r.t. inputs:

$$- dy/dx_1 = (0.217 \cdot 0.25 \cdot 2) + (-0.434 \cdot 0.1056 \cdot -1) = 0.1544$$

$$- dy/dx_2 = (0.217 \cdot 0.25 \cdot -1) + (-0.434 \cdot 0.1056 \cdot 1.5) = -0.1229$$

$$- dy/dx_3 = (0.217 \cdot 0.25 \cdot 0.5) + (-0.434 \cdot 0.1056 \cdot -0.5) = 0.05$$

Step 4: Interpretation

- x_1 (0.1544) has a strong positive effect on y .
- x_2 (-0.1229) has a negative effect on y .
- x_3 (0.05) has a small positive effect.

This saliency map helps interpret how each input affects the model's decision.