

Let L be the loss function.

z_1 is the pre-activation at layer 1

B_1 is the parameter matrix (weights and biases) for layer 1.

Then the gradient of loss with respect to parameters at layer 1

$$\frac{dL}{dB_1} = \frac{dL}{dz_1} \cdot \frac{dz_1}{dB_1}$$

Chain rule for gradients state that

$$\frac{dL}{dB_1} = \frac{dL}{dz_L} \cdot \frac{dz_L}{dz_{L-1}} \cdot \frac{dz_{L-1}}{dz_{L-2}} \cdot \dots \cdot \frac{dz_{L+1}}{dz_1} \cdot \frac{dz_1}{dB_1}$$

Considering all weights are initialized to 0

$$\frac{dz_1}{dB_1} = 0 \text{ when multiplied in the chain}$$

rule will result to 0 for all layers following the gradient descent.

For the last layer L , the error is directly dependent on the final loss.

Therefore, $\frac{dL}{dB_L} \neq 0$ for the final layer.

The chain rule for gradient descent will not equate to zero as $\frac{dL}{dB_L} \neq 0$.

Therefore, during back propagation from final loss the parameters will update.