

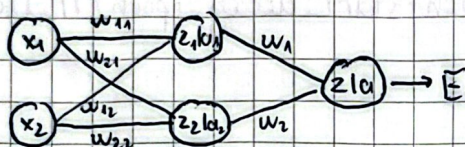
Section 7: Weight Initialization

- If we initialize all weights to 0, knowing that the updates to the weights is dependant on previous weights. And since ~~none~~ all of our weights is 0, no learning will occur

Ex: Our NN in Section 1, $\frac{\partial E}{\partial w_{11}} = \left\{ \sum_{j=1}^2 \left[\frac{\partial E}{\partial \delta(y_j)} \cdot \frac{\partial \delta(y_j)}{\partial y_j} \cdot \frac{\partial y_j}{\partial \delta(z_1)} \right] \right\} \cdot \frac{\partial \delta(z_1)}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{11}}$

we have $\frac{\partial z_1}{\partial w_{11}} = \frac{\partial y_1}{\partial \delta(z_1)} = w_{21} = 0, \frac{\partial y_2}{\partial \delta(z_1)} = w_{22} = 0$

- Random initialization: maybe vanish and explode.
- Random initialization from Normal Distribution - 1st way
- If we initialize all weights to the same value, for example



$$\begin{aligned} z_1 &= x_1 w_{11} + x_2 w_{12} \\ z_2 &= x_1 w_{21} + x_2 w_{22} \\ a_1 &= \sigma(z_1), a_2 = \sigma(z_2) \\ z_3 &= a_1 w_1 + a_2 w_2, a = \sigma(z) \end{aligned}$$

Weight update for w_1 : $\frac{dE}{dw_1} = \frac{dE}{da} \cdot \frac{da}{dz} \cdot \frac{dz}{dw_1}$

Weight update for w_2 : $\frac{dE}{dw_2} = \frac{dE}{da} \cdot \frac{da}{dz} \cdot \frac{dz}{dw_2}$

$$\frac{dz}{dw_1} = a_1 = \sigma(z_1) = \sigma(z_2) = a_2 = \frac{dz}{dw_2}$$

Both weight will get updated with the same value and they will be equal. This is called

Symmetry Breaking Problem

- Xavier Initialization:

- 2nd Way

- The output of a linear layer: $Y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n$
- Assume x_i and w_i are all independent and identically distributed

$$\begin{aligned} \text{Var}(w_i x_i) &= E[x_i^2] \text{Var}(w_i) + E[w_i^2] \text{Var}(x_i) + \text{Var}(w_i) \text{Var}(x_i) \\ &= \text{Var}(w_i) \text{Var}(x_i) \end{aligned}$$

$$\text{Then, } \text{Var}(Y) = \text{Var}\left(\sum_i w_i x_i\right) = n \text{Var}(w_i) \text{Var}(x_i)$$

$$\text{We want } \text{Var}(Y) = \text{Var}(x_i), \text{ then } \text{Var}(w_i) = \frac{1}{n} = \frac{1}{n_{in}}$$

Go through same steps for backpropagating, we also have $\text{Var}(w_i) = \frac{1}{n_{out}}$

Compromising: $\text{Var}(w_i) = \frac{2}{n_{in} + n_{out}} \Rightarrow \text{std} = \sqrt{\text{Var}} = \sqrt{\frac{2}{n_{in} + n_{out}}}$

• He Norm Initialization:

A ReLU is zero for half of its input (negative inputs), so we need to double the size of the weight variance to keep the signal's variance constant. Therefore, we multiply our previous equation of $\frac{1}{n_{in}}$ by 2.

$$\text{Var}(W) = \frac{2}{n_{in}} \quad \text{std} = \sqrt{\text{Var}} = \sqrt{\frac{2}{n_{in}}}$$

So,

+ When the layers have ReLU activation \rightarrow He initialization

+ When the layers have Sigmoid activation \rightarrow Xavier initialization.