

Varying Gradient Descent

↳ Backpropagation

80's & 90's

Gradient Descent

Learning Rate ≈ 0.001

Gradient

$$\underline{w_{new}} = \underline{w_{old}} - \eta \left[\frac{\delta L}{\delta w_{old}} \right]$$

Really small

Learning gets stopped

$w_{new} \approx w_{old} \rightarrow$ Varying Gradient Descent

$$\eta_{gradient} = \frac{\delta L}{\delta w_{old}} = \frac{w_{old} - w_{new}}{\eta}$$

Multiple Hidden Layers

0 to 0.5

Activation function

derivative

Sigmoid

tanh

chain rule

$$\frac{\delta L}{\delta w'_{11}} = \frac{\delta L}{\delta \hat{y}} * \frac{\delta \hat{y}}{\delta z} * \frac{\delta z}{\delta o_{11}} * \frac{\delta o_{11}}{\delta w'_{11}}$$

$$0 < \frac{\delta L}{\delta w} < 1$$

$$= 0.1 * 0.1 * 0.1 * 0.1$$

$$= 0.0001$$

↳ Really small

Value

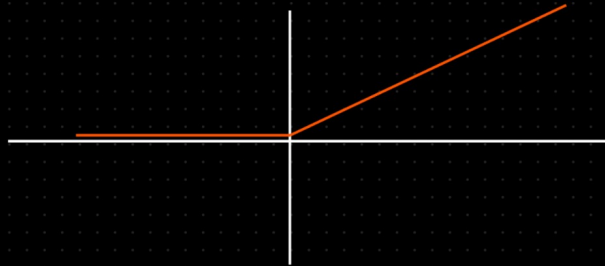
The vanishing gradient problem is a challenge that occurs during the training of artificial neural networks when gradients become very small and diminish as they are propagated back through the network.

In machine learning, the vanishing gradient problem is encountered when training neural networks with gradient-based learning methods and backpropagation. In such methods, during each training iteration, each neural network weight receives an update proportional to the partial derivative of the loss function with respect to the current weight.[1]

The problem is that as the network depth or sequence length increases, the gradient magnitude typically is expected to decrease (or grow uncontrollably), slowing the training process.[1] In the worst case, this may completely stop the neural network from further learning

* ReLU \rightarrow avoid issue of vanishing
gradient issue

$$\max(0, x_i)$$

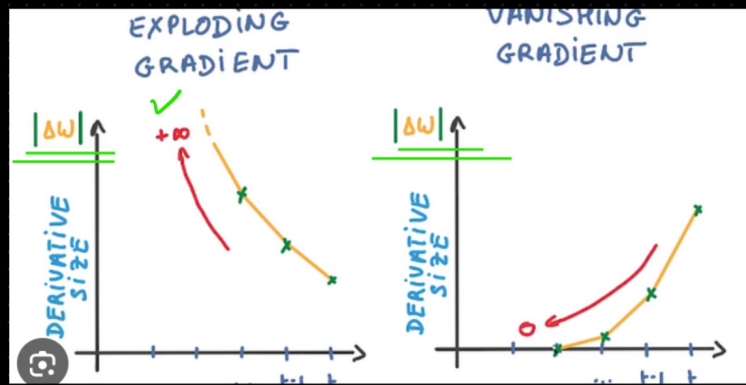


Exploding Gradient Descent

$$\left\{ \begin{aligned} \frac{\delta L}{\delta w'_{11}} &= \frac{\delta L}{\delta \hat{y}} * \frac{\delta \hat{y}}{\delta z} * \frac{\delta z}{\delta o_{11}} * \frac{\delta o_{11}}{\delta w'_{11}} \\ &\approx 10 * 10 * 10 * 10 \\ &= 10000 \end{aligned} \right.$$

\swarrow drastic change \searrow big

$$\underline{w_{\text{new}}} = \underline{w_{\text{old}}} - \eta \boxed{\frac{\delta L}{\delta w_{\text{new}}}}$$



Epoch = 20

$$\text{Loss} = 0.8$$

$$\text{Loss} = 0.7$$

$$\text{Loss} = 0.6$$

$$\text{Loss} = 0.5 \checkmark$$

$$\text{Loss} = 0.5 \checkmark$$

Vanishing α —

gradient

decay

