

04 Oct 2024

WEIGHT INITIALIZATION

Why weight initialization is important?

- (i) initialize the parameters
- (ii) Choose an optimization algorithm
- (iii) Repeat these steps.

- Forward propagate an input
- Compute the cost function.
- Compute the gradients of the cost with respect to parameters using backpropagation
- update each parameter using the gradients according to the optimization algorithm.

problems if wrong initialization of weights:

- vanishing gradient problem
- Exploding gradient problem
- slow convergence.

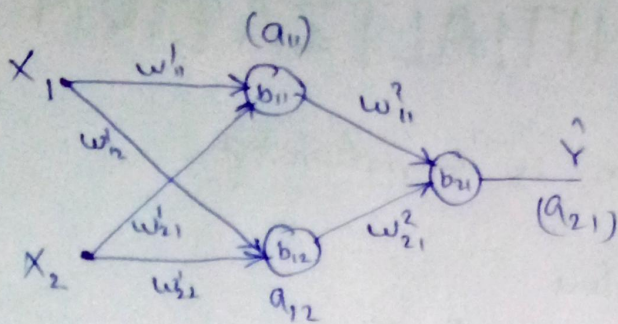
what not to do in weight initialization

case-1 \rightarrow zero initialization.

Question: why we not set the initial value of weight of zero ('0')?

let us take an example, regression example
assume. $w_0 = 0$ and $b = 0$

(i) ReLU in hidden layer:



Here $a_{11} + a_{12} + a_{21}$ is activation function

$a_{11} + a_{12} \rightarrow \text{ReLU}$

$a_{21} \rightarrow \text{linear}$

$$a_{11} = \max(0, z_{11})$$

$$z_{11} = w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}$$

$$a_{12} = \max(0, z_{12})$$

$$z_{12} = w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}$$

because $w=0$ and $b=0$

Hence,

$$a_{11} = 0, z_{11} = 0 \text{ and } a_{12} = 0, z_{12} = 0$$

$$\text{Hence, } a_{11} = a_{12} = 0$$

for weight update

$$w_{11}^{\text{new}} = w_{11}^{\text{old}} - \eta \frac{\partial L}{\partial w_{11}} \rightarrow 0$$

$$\therefore w_{11}^{\text{new}} = w_{11}^{\text{old}}$$

Here no update

thus there is no training ~~that~~ will take place

(ii) Tanh in hidden layer:

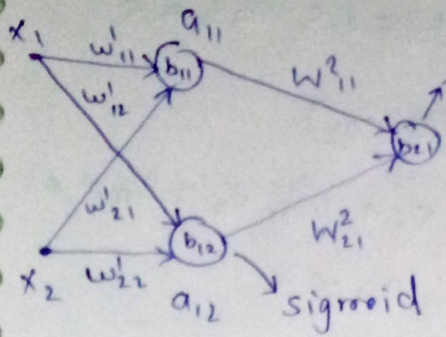
$$a_{11} = \frac{e^{z_{11}} - e^{-z_{11}}}{e^{z_{11}} + e^{-z_{11}}}$$

$$\therefore z_{11} = 0$$

$$\therefore a_{11} = 0 \text{ and } a_{12} = 0 \text{ also}$$

similar case as in previous one.

(iii) Sigmoid in hidden Layer:



$$a_{11} = \sigma(z_{11}) = 0.5$$

linear activation similarly $a_{12} = 0.5$

but $a_{11} = a_{12}$

also $z_{11} = z_{12}$

weight upd update.

$$w = w - \eta \frac{\partial L}{\partial w}$$

$$\frac{\partial L}{\partial w'_{11}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{11}} \cdot \frac{\partial a_{11}}{\partial z_{11}} \right] \cdot \left(\frac{\partial z_{11}}{\partial w'_{11}} \right) \rightarrow x_1$$

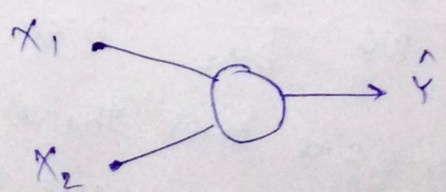
$$\frac{\partial L}{\partial w'_{12}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial z_{12}} \right] \cdot \left(\frac{\partial z_{12}}{\partial w'_{12}} \right) \rightarrow x_1$$

$$\frac{\partial L}{\partial w'_{21}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{11}} \cdot \frac{\partial a_{11}}{\partial z_{11}} \right] \cdot \left(\frac{\partial z_{11}}{\partial w'_{21}} \right) \rightarrow x_2$$

$$\frac{\partial L}{\partial w'_{22}} = \left[\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial z_{12}} \right] \cdot \left(\frac{\partial z_{12}}{\partial w'_{22}} \right) \rightarrow x_2$$

Hence $\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial w'_{12}}$ and $\frac{\partial L}{\partial w'_{21}} = \frac{\partial L}{\partial w'_{22}}$

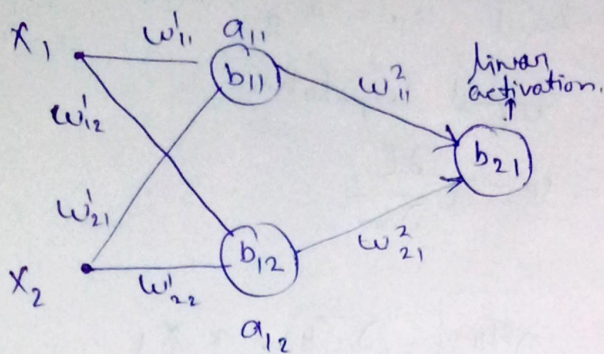
→ The neural network becomes



not capture non-linearity because. It work as a perceptron

Case-2 → Non-zero initialization (constant value)

let us assume $W=0.5$ $b=0.5$



$$a_{11} = \max(0, z_{11})$$

$$z_{11} = w'_{11} x_1 + w'_{21} x_2 + b_{11} \neq 0$$

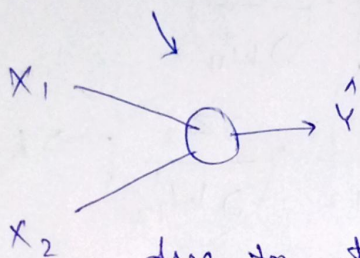
$$\text{hence } a_{11} \neq 0$$

$$a_{12} = \max(0, z_{12})$$

$$z_{12} = w'_{12} x_1 + w'_{22} x_2 + b_{12} \neq 0$$

$$\text{hence } a_{12} \neq 0$$

$$z_{12} = z_{11} \rightarrow a_{11} = a_{12}$$



due to this,

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{11}} \cdot \frac{\partial a_{11}}{\partial z_{11}} x_1$$

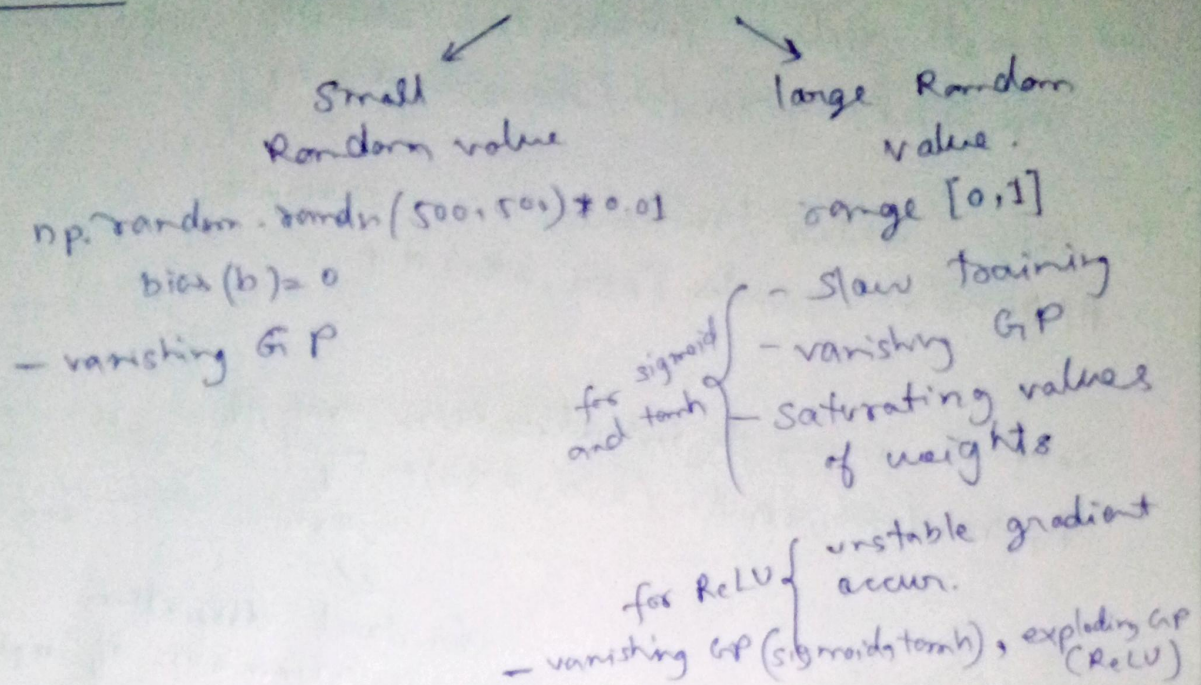
$$\frac{\partial L}{\partial w'_{12}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial z_{12}} x_1$$

$$\frac{\partial L}{\partial w'_{21}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{11}} \cdot \frac{\partial a_{11}}{\partial z_{11}} x_2$$

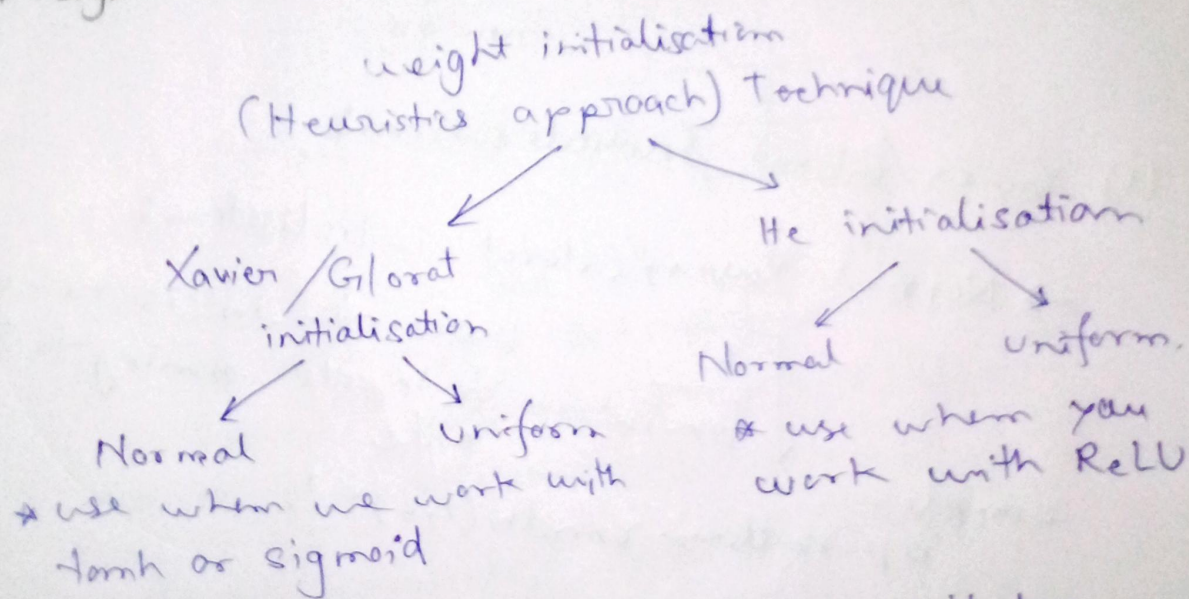
$$\frac{\partial L}{\partial w'_{22}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial a_{12}} \cdot \frac{\partial a_{12}}{\partial z_{12}} x_2$$

~~then~~ here we see that there is similar case is occur as previous one. (zero initialization of weigh and bias with hidden layer activation function is sigmoid)

case-3 → Random Initialization.



Weight Initialisation Techniques:



- till now we know the weight is initialize with Random value but we don't know what it's range. because small random value and large random value can be face a problem.

→ small value initialization example
`np.random.randn(250, 250) * 0.01`

→ large value initialization example
`np.random.randn(250, 250) * 1`

correct way to initialization example

$$\text{np.random.randn}(250, 250) * \sqrt{\frac{1}{250}}$$

↙
standard deviation
250 is number of input
for given particular
neuron.

(i) Xavier/Glorot Initialisation:

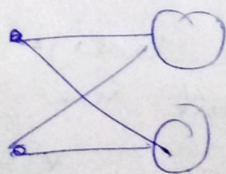
→ Normal Xavier/Glorot initialisation:

$$\sqrt{\frac{1}{\text{fan-in}}}$$

here fan-in is number
of inputs coming to the
node

example.

$$\text{np.random.randn}(2, 2) * \frac{1}{2}$$



or we can also use $\sqrt{\frac{2}{\text{fan-in} + \text{fan-out}}}$

It is used in the case of tanh and sigmoid.

→ uniform distribution Xavier/Glorot initialisation:

$$[-limit, limit]$$

$$limit = \sqrt{\frac{6}{fan-in + fan-out}}$$

(ii) He initialisation:

→ Normal He initialisation:

$$\sqrt{\frac{2}{fan-in}}$$

→ uniform distributed He initialisation:

$$[-limit, limit]$$

$$\text{Here, } limit = \sqrt{\frac{6}{fan-in.}}$$

* all these things are done in the code.
as a parameter set ~~for~~ using.

`kernel_initializer = 'initialize_type'`,
but in the default case.

`kernel_initializer = 'glorot_uniform'`