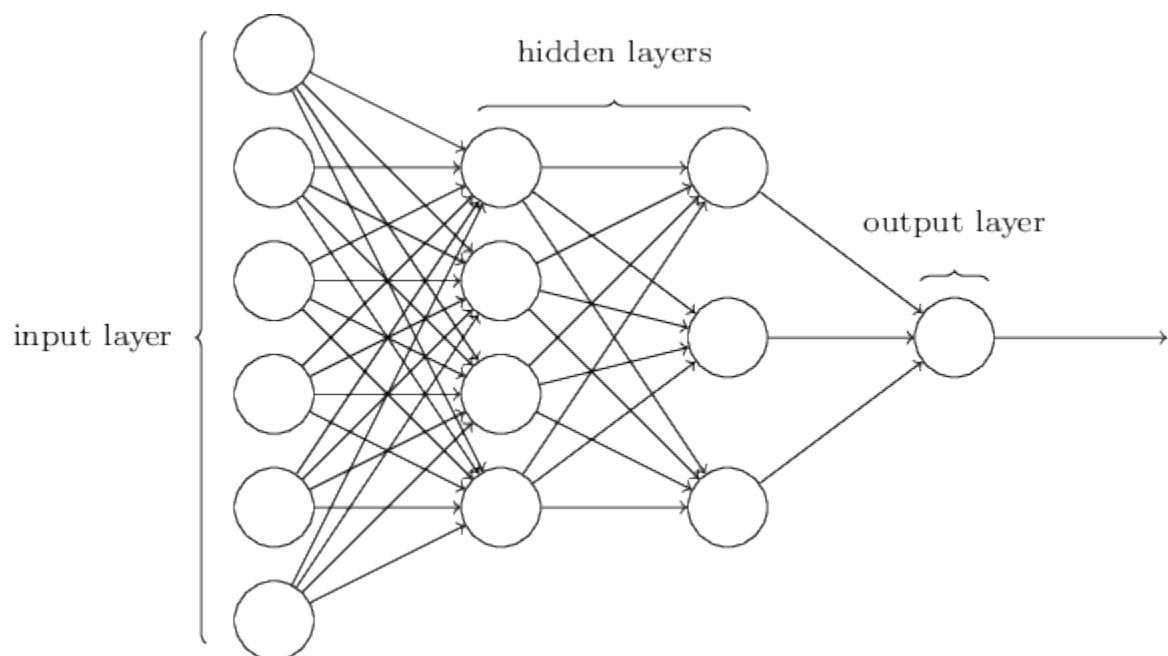


# Lecture Notes

## Part 1

### Neural Networks and Deep Learning



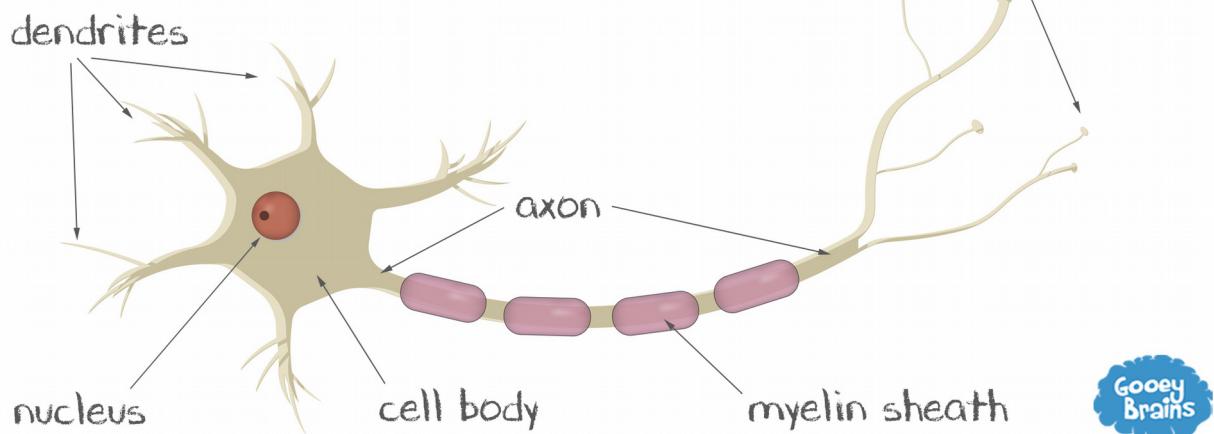
Prof. Ankita Pramanik  
IEST, Shibpur

Siladitya Manna  
IEST, Shibpur

## 1. What is a Neuron?

The fundamental units of brain and nervous system, responsible for receiving sensory information from different parts of an animal body and transmitting motor command to the muscles.

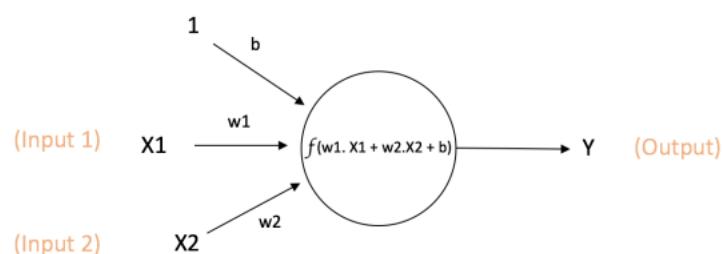
# NEURON



## 2. What is a perceptron?

Perceptron is the mathematical model of a biological neuron.

As in the biological neuron, the dendrites receive the electrical signals from the axons of other neurons, the perceptrons receive numerical inputs. The synapses are modeled by incorporating a weight to each input, before the inputs are added. The perceptron also exhibits the firing characteristics of a biological neuron when the input is above a certain threshold, by incorporating an activation function  $f$ , which represents the frequency of spikes along the axon.



$$\text{Output of neuron} = Y = f(w_1 \cdot X_1 + w_2 \cdot X_2 + b)$$

### 3. What is Neural Network?

Neural Network is computational system, formed by stacking several neurons/perceptrons together, used for supervised and unsupervised learning, and solving complex problems, and does not require to be hard-coded/programmed for any specific task, but works by training itself from given examples.

### 4. Binary Classification

Binary Classification is the task of classifying a given element into any one of given set of two groups.

The idea of Binary Classification can be conveyed using Logistic Regression, which is an algorithm for binary classification.

Suppose, we want to train a classifier, which will predict if the input is a '3' or not.



The image is taken from MNIST Dataset, which is a hand-written digit dataset, and contains images of numbers from 0-9. The images are grayscale images and have size 28X28.

The image is unwrapped and a feature vector  $\mathbf{x}$  is created. (Start from the top left corner pixel and move right until you reach the end, then start again from the first(leftmost) pixel of the next line)

If  $n_x$  be the number of elements in  $\mathbf{x}$  then  $n_x = 28*28 = 724$

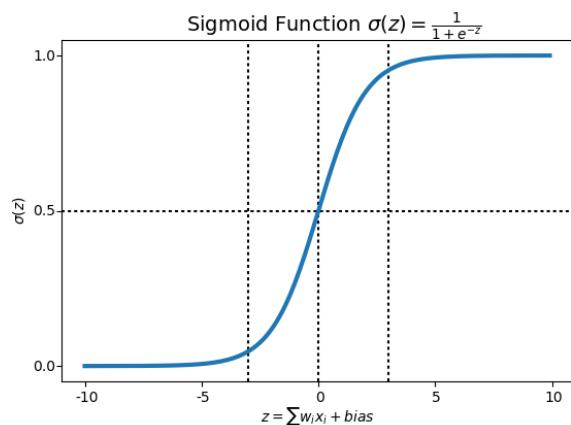
$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 235 \\ 100 \\ 210 \\ \vdots \end{bmatrix}$$

## 5. Logistic Regression

Hence, given input  $\mathbf{x}$ , we want  $\hat{y} = P(y=1|\mathbf{x})$

Input:  $\mathbf{x} \in \mathbf{R}^{nx}$  Parameters:  $\mathbf{W} \in \mathbf{R}^{nx}$ ,  $b \in \mathbf{R}$

Output:  $\hat{y} = \sigma(\mathbf{W}^T \mathbf{x} + b)$ , where  $\sigma(z) = 1/(1+e^{-z})$



Hence,  $0 \leq \hat{y} \leq 1$

So the job is to learn the parameters  $\mathbf{W}$  and  $b$ , such that,  $\hat{y}$  is close to 1 when  $z$  is large and close to 0 when  $z$  is very small(negative).

## 6. Cost Function

To train the Logistic Regression, a cost Function is required.

Given  $\{(x^{(i)}, y^{(i)})\}$  we want  $\hat{y}^{(i)} \approx y^{(i)}$ , where  $x^{(i)}$  is the  $i^{\text{th}}$  training example.

### Loss Function

Let us assume, the squared error function as the Loss function

$$L(\hat{y}^{(l)}, y^{(l)}) = \frac{1}{2}(\hat{y}^{(l)} - y^{(l)})^2$$

But, this loss function is not a good choice for logistic regression, because, the error surface is not always a convex one, and hence, may get stuck on multiple local minimas.

Hence, a more suitable choice for the loss function is

$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$

- If  $y^{(i)} = 1: L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$  where  $\log(\hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 1
- If  $y^{(i)} = 0: L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$  where  $\log(1 - \hat{y}^{(i)})$  and  $\hat{y}^{(i)}$  should be close to 0

The above two points explains the intuitive logic for which the above function is considered as a good choice for the Loss function.

### Cost Function

The cost function is the average of the loss function of the entire training set.

The objective is to find the parameters W and b that minimizes the overall cost function

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

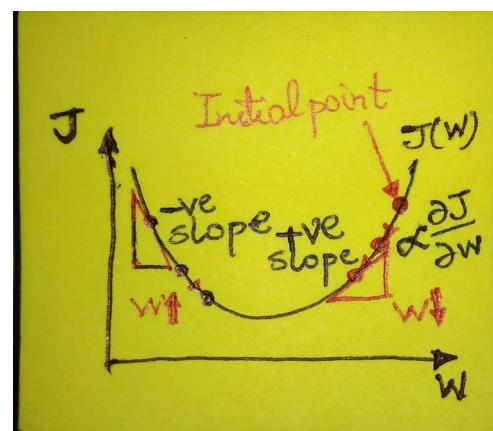
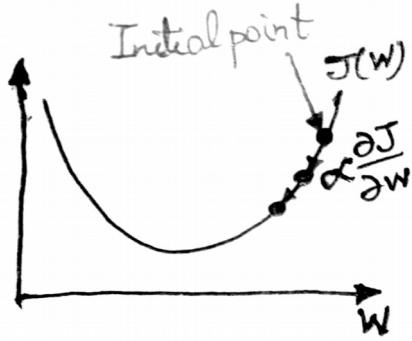
An important to note here is that the loss function computes the error for a single training example; the cost function is the average of the loss functions of the entire training set.

### 7. Gradient Descent

However, to train the neural network, the weights or parameters need to be updated.

Otherwise stated, we want to find the values of W and b for which  $J(W, b)$  is minimized.

First, W and b needs to be initialized to some values. Usually, in logistic regression, W and b are initialized to 0, however, random initialization is also done.



As shown in the figure, the slope  $\partial J(W,b)/\partial W$  is negative on the left of the minima and

positive on the right of the minima. So, the value of  $W$  increases when the point is on the left and increases when on the right.

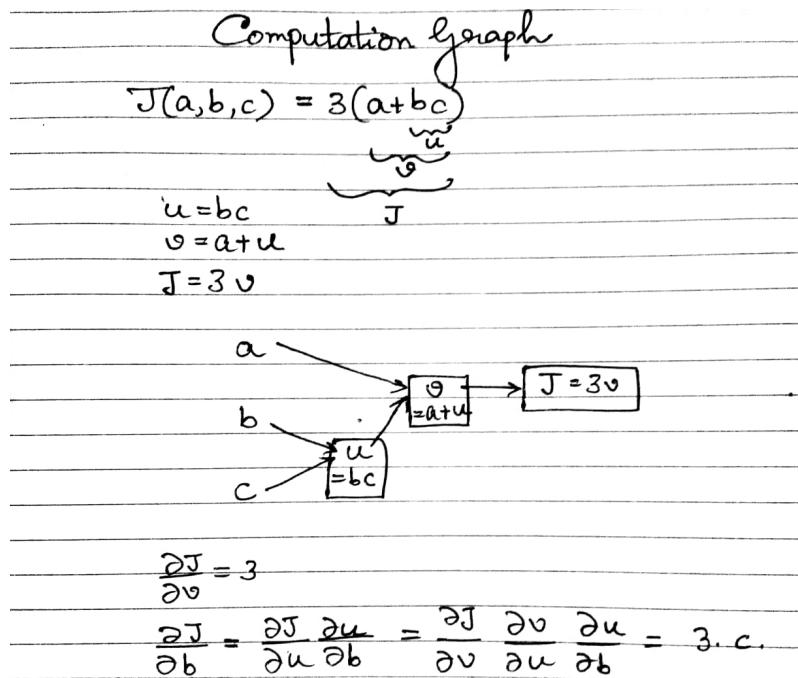
### Update Rule:

$$W = W - \alpha \cdot \partial J(W,b) / \partial W \quad b = b - \alpha \cdot \partial J(W,b) / \partial b, \text{ where } \alpha \text{ is the learning rate.}$$

$\partial J(W,b)/\partial W$  is the derivative of the cost function  $J(W,b)$  with respect to  $W$ , that is the

change of the cost function, in the direction of  $W$ .

### 8. Computation Graph



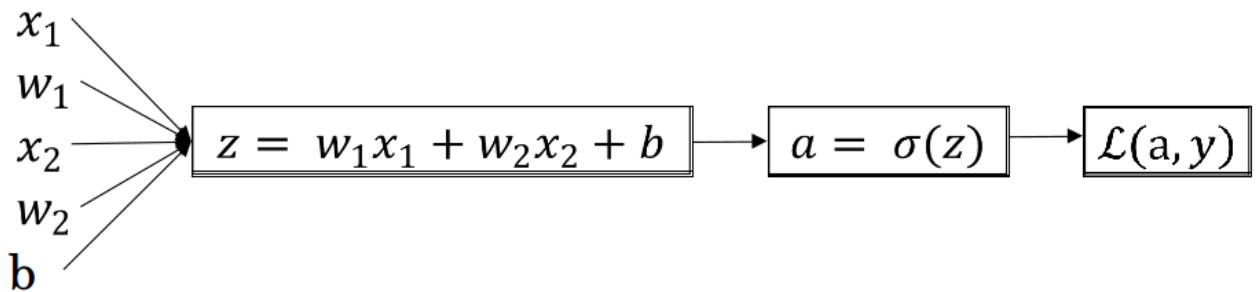
You can see that computing derivatives using computation graph is basically same as **Chain Rule**.

## 9. Logistic Regression Gradient Descent

Let, the feature vector  $\mathbf{x}$  contain only two values, i.e.  $\mathbf{x} = [x_1 \ x_2]^T$

Now, if we write the operation of the neural network as a computation graph then, in addition to the feature vector, we also need to feed the parameters as input.

The computation graph will look like this:



This shows the result of forward propagation. Now, by applying backward propagation, we can get the change in the values of the parameters

$$\begin{aligned}
 & \text{Forward Propagation: } z = w_1x_1 + w_2x_2 + b \rightarrow a = \sigma(z) \rightarrow L(a, y) \\
 & \text{Backward Propagation: } \frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial}{\partial a} [\ln(1+a) - y \ln a - (1-y) \ln(1-a)] \\
 & \quad = \frac{\partial L}{\partial a} = -\frac{y}{a} + \frac{1-y}{1-a} \\
 & \quad \frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_1} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1} = -\frac{y}{a} + \frac{1-y}{1-a} \cdot 1 \cdot x_1 \\
 & \quad \frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial z} \cdot \frac{\partial z}{\partial w_2} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_2} = -\frac{y}{a} + \frac{1-y}{1-a} \cdot 1 \cdot x_2
 \end{aligned}$$

Here,  $a$  is the output of the activation function **sigmoid** and  $y$  is the true label of the data.

## 10. Logistic Regression on m examples

For  $m$  training examples, we need to calculate the cost function value

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y^{(i)})$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^\top x^{(i)} + b)$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)})$$

$$= \frac{1}{m} \sum_{i=1}^m \frac{\partial z^{(i)}}{\partial w_1} \frac{\partial a^{(i)}}{\partial z^{(i)}} \frac{\partial L(a^{(i)}, y^{(i)})}{\partial a^{(i)}}$$

$$= \frac{1}{m} \sum_{i=1}^m x_1^{(i)} (a^{(i)} - y^{(i)})$$

$$\omega_1 = \omega_1 - \alpha \frac{\partial}{\partial w_1} J(w, b) \quad \text{Update Rule:}$$

$$= \omega_1 - \alpha \sum_{i=1}^m x_1^{(i)} (a^{(i)} - y^{(i)})$$

Now, if we implement this using a for loop in any programming language, the code will not be efficient, because as the dataset increases, the neural network will take more iterations to train. Hence, computational time will be huge.

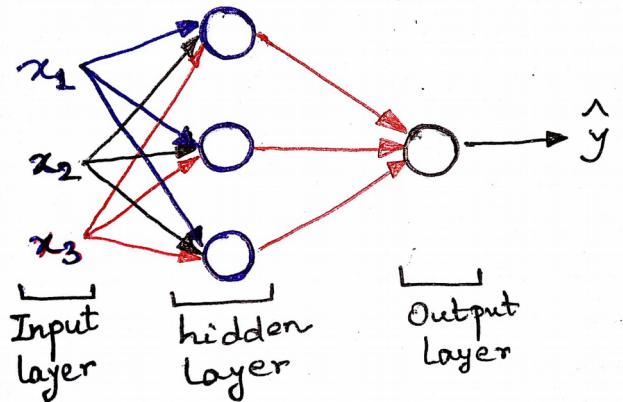
Hence, here comes the concept of vectorization.

## 11. Vectorization

Vectorization will be discussed, during **Part 3.a.**

## 12. Neural Network Representation

Neural Network with a single hidden layer



This is a **2-layer neural network** (because the input layer is not counted)

The term “**Hidden Layer**” refers to the fact that in the training set, the true values of the nodes in these layers are not observed in the training set, hence they are called hidden layers. The function of the Hidden Layer is to take the input and apply a linear transformation to the input, followed by a squashing non-linearity. However, we have values for the input layers and the output label in the training set.

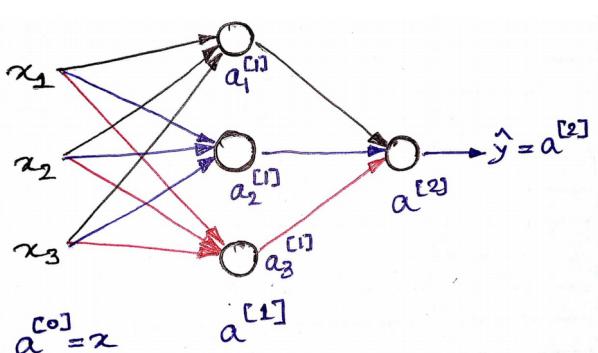
The letter “**a**” refers to the output of a layer which is passed on to the next layer as input.

The input layer passes on the input to the first hidden layer, hence the input layer is termed as the 0<sup>th</sup> layer  $\mathbf{a}^{[0]}$ .

**The output of the first hidden layer is denoted by  $\mathbf{a}^{[1]}$ .** For a hidden layer with  $n$

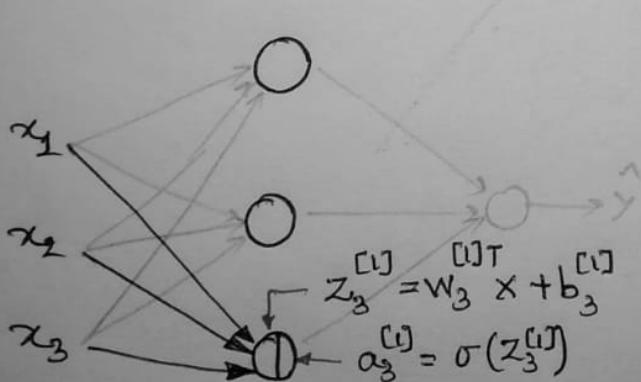
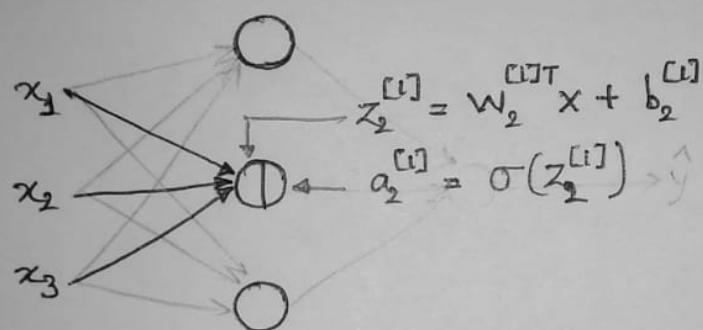
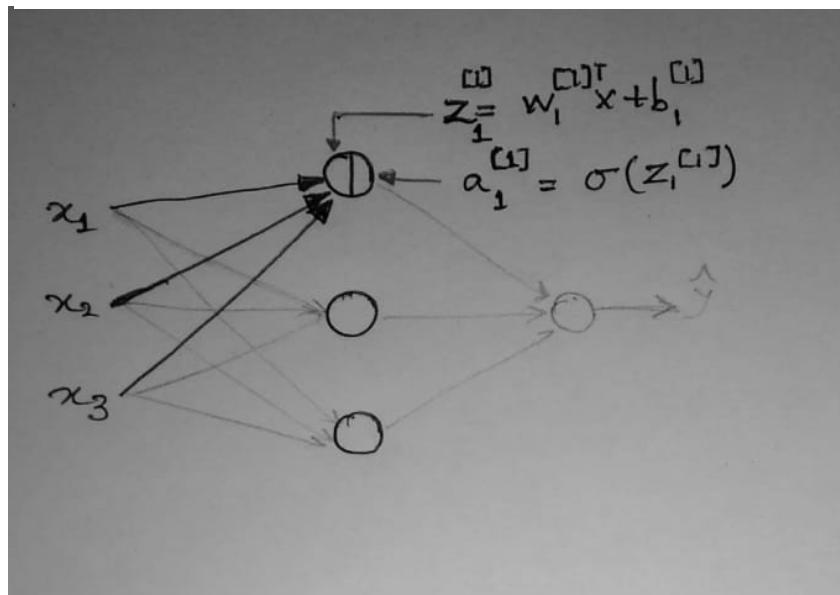
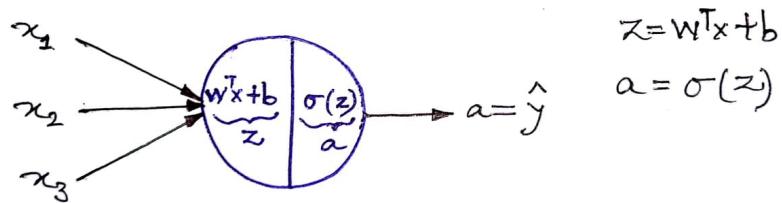
$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ \vdots \\ a_n^{[1]} \end{bmatrix}$$

nodes, the activation generated is

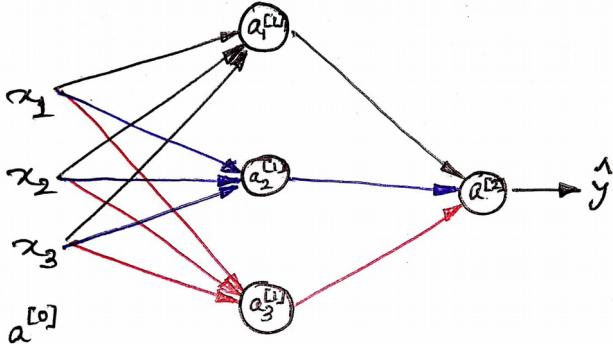


### 13. Computing a Neural Network's Output

Two steps of Computation in a perceptron



Writing all the equations together,



$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]}).$$

The equations can be written as

$$\begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix} = \underbrace{\begin{bmatrix} w_1^{[1]T} \\ w_2^{[1]T} \\ w_3^{[1]T} \end{bmatrix}}_{W^{[1]} (n_a^{[0]}, n_a^{[1]})} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \underbrace{\begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \end{bmatrix}}_{b^{[1]} (n_a^{[1]}, 1)}$$

$$z^{[1]} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \end{bmatrix}; \quad a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \end{bmatrix} = \sigma(z^{[1]})$$

The size of the resulting weight matrix  $W$  is thus  $(n_a^{[0]}, n_a^{[1]})$  and that of  $b$  is  $(n_a^{[0]}, 1)$

Thus given input  $x$

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

For n-layer neural net

$$z^{[n]} = W^{[n]}a^{[n-1]} + b^{[n]}$$

$$a^{[n]} = \sigma(z^{[n]})$$

### 13. Vectorizing across multiple examples

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

For  $m$  examples,

$$X = [x^{[1]} \ x^{[2]} \ \dots \ x^{[m]}]$$

Size of  $X$  is  $(n_x \times m)$

Hence,  $Z^{[1]}$  becomes of size  $(n_{a[1]} \times m)$ , and so does  $A^{[1]}$

Hence at the input of the second layer the input is of the size  $n_{a[1]} \times m$

#### For the second layer

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

Size of  $W^{[2]}$  is  $(n_{a[2]} \times n_{a[1]})$

Hence, size of  $Z^{[2]}$  is  $(n_{a[2]} \times m)$  and so is for  $A^{[2]}$ .

### 13. Activation Functions

Activation Functions are a very important part of neural nets. It is responsible for the non-linear mapping between the input and the output. Without the activation function, the Neural Network would simply be a linear function of degree 1. Hence, the neural network would be unable to learn complex functional relations between the input and the target. Therefore, we need to apply the non-linear activation functions, so as to empower the neural network to learn the non-linear complex functional mappings that exists between the inputs and outputs, and represent it too.

Furthermore, the activation functions should be differentiable, otherwise, back-propagation cannot be performed.

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$a^{[1]} = g(z^{[1]})$$

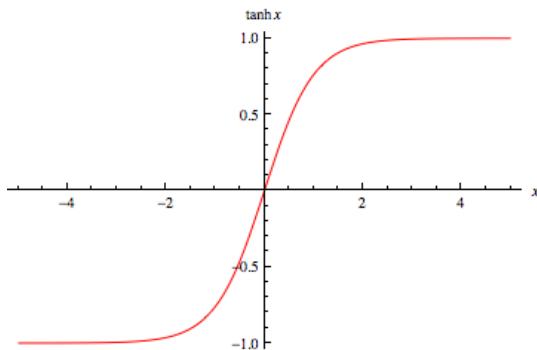
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = g(z^{[2]})$$

In more general case, we consider  $g$  to be the activation function

**The activation function which almost always works better than the sigmoid function is the tanh function.**

$$\tanh(z) = (e^z - e^{-z}) / (e^z + e^{-z})$$

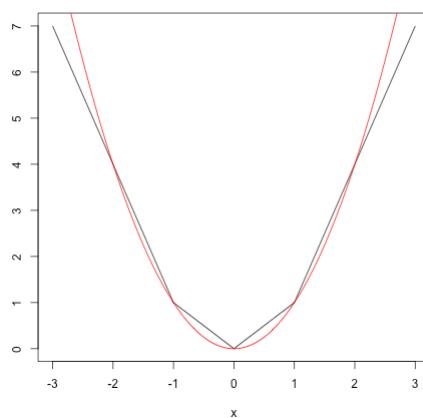
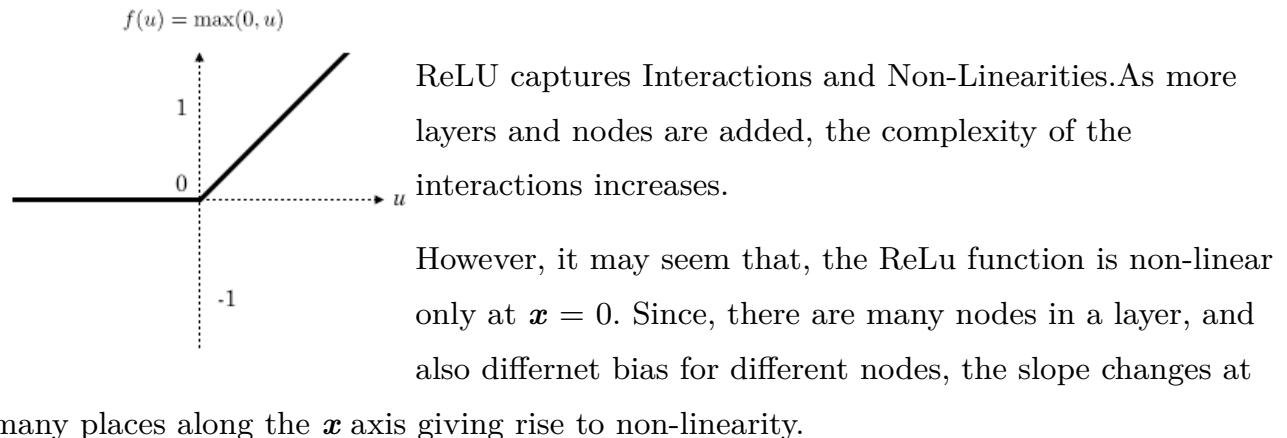


This function has the effect of centering the data nearer to 0 mean, and it makes learning a little bit easier, than in sigmoid, in which the data gets centered at 0.5.

However, for the output layer, the activation function is preferably sigmoid, because we want the output label to be in  $[0,1]$ .

Hence, the activation function can be different for different layers depending on the function.

Rectified Linear Unit( ReLU): It is the most commonly used activation function.



There are some disadvantages of these activation functions.

Firstly, for **tanh** function, the non-linear part of the tanh function lies roughly between -2.0 to 2.0 and the rest of the function is relatively flat. So, the change in the output will not be well reflected in the input. The derivative of the **tanh** function beyond this narrow range, is almost zero, which makes update of weights very difficult. This problem gets worse, if the number of layers increases. This problem is termed as **vanishing gradient problem**.

**ReLU however rectifies vanishing gradient problem.**

In a multilayer network, during back propagation, the derivatives are gradually multiplied by more and more terms arising from the derivatives of that layer. If tanh or sigmoid functions are used then the updates becomes close to zero, hence we get a vanishing gradient.

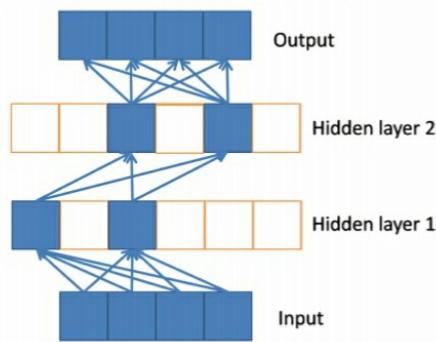
However, if we use ReLU, then the derivative of ReLU is always 1, if  $x > 0$ , i.e. the activation is positive. That is, the error signal is fully propagated to the input, and the

gradient value does not decrease as we move up the network layers. Thus the vanishing gradient problem is resolved. ReLU is actually more biologically plausible and enables the network to obtain sparse representations.

### Other advantages of ReLU:

Network easily obtains sparse representation. Refer to this [link](#) for the advantages of sparsity introduced by ReLU. The only, non-linearity that comes due to ReLU, is due to the path selection associated with the individual active neurons.

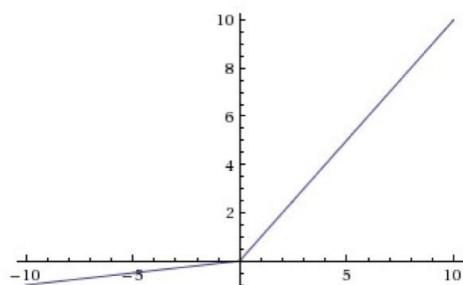
The function obtained by using ReLU as an activation function in the hidden layers, is linear by parts, due to which the gradient flows well on the active paths of the neurons.



ReLU however has a **disadvantage** too. Since ReLU gives zero output for inputs less than zero, the gradient can go to zero, resulting in the weights not being updated. Hence, neurons stop responding to changes in error/input. Thus, it can make a neuron never activate again, and thus making a **Dead Neuron**. To fix this problem **Leaky ReLU** is used.

$$\text{Leaky ReLU: } f(x) = x, x > 0$$

$$= 0.01x, x < 0$$

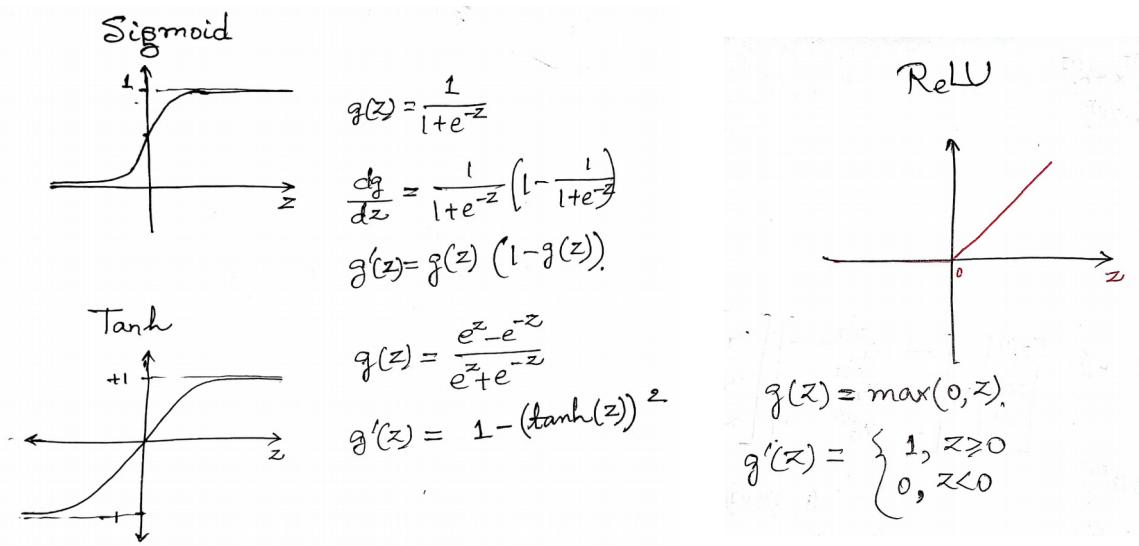


$$\text{Parametric ReLU: } f(x) = x; x > 0$$

$$= a \cdot x; x < 0$$

The parameter  $a$  is trainable in PReLU.

## Derivatives of activation Functions



## Gradient Descent for Neural Networks

### Forward Propagation

$$z^{[1]} = W^{[1]}x + b^{[1]}$$

$$A^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g(z^{[2]}).$$

### Back Propagation

$$dz^{[2]} = A^{[2]} - Y$$

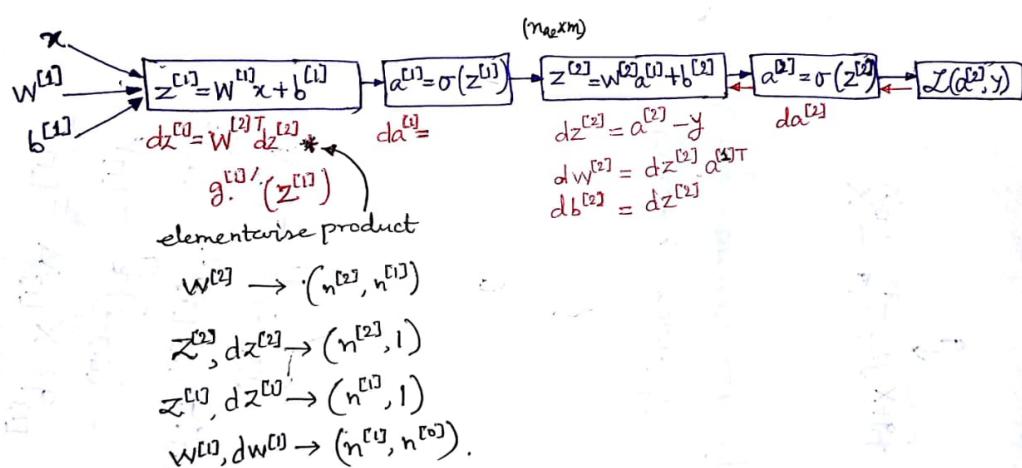
$$dW^{[2]} = \frac{1}{M} dz^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{M} \text{ [summation of } dz^{[2]} \text{ over axis } = 1 \text{].}$$

$$dZ^{[2]} = \frac{\partial \mathcal{L}}{\partial z^{[1]}} = W^{[2]T} dz^{[2]} g'(z^{[2]})$$

## Backpropagation Explanation using Computation Graph

$$dw^{[1]} = dz^{[1]} \cdot x^T$$



$$da^{[1]} = \frac{\partial L}{\partial a^{[1]}} = \frac{\partial L}{\partial z^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} = (a^{[2]} - y) W^{[2]}$$

To match the dimension of  $dW^{[2]}$  with the RHS,  $a^{[1]T}$  is taken.

## Random Initialization

### 1. Zero Initialization

The problem with zero initialization is that,  $a_1^{[1]}$  and  $a_2^{[1]}$  will be equal. Because both of the hidden units will be computing the same function. Also, during backpropagation,  $dz_1^{[1]}$  and  $dz_2^{[1]}$  will also be same, due to symmetry. That is both the hidden units are computing the same function. So, even if weight is updated,

the error signal propagating back through the network will also be same for both nodes. Hence, the updated weights:

$$W^{[1]} = W^{[1]} - \alpha dW$$

Hence, the updated weights will be same for the hidden units. Hence, even after several iterations, it will be computing the same function.

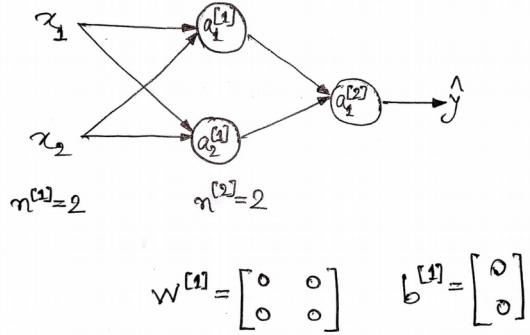
This happens if the weights are initialized to the same value, other than zero.

No matter what was the input - if all weights are the same, all units in hidden layer will be the same too. In other words, the neural network fails to break symmetry.

The solution to this is to initialize the weights randomly.

### 2. Xavier Initialization

Now for random initialization, if the weights are too small then the variance of the input signal propagating through the network, starts decreasing and if the weights are too large then the variance increases rapidly. Eventually it becomes very large and becomes useless. This is mainly observed, if sigmoid function is used.



$$dW = \begin{bmatrix} u & v \\ u & v \end{bmatrix}$$

Now, from the computation graph we have seen that

$$dz^{[i]} = W^{[i+1]} dz^{[i+1]} * g'(z^{[i]})$$

$$dW^{[i]} = dz^{[i]} a^{[i]T}$$

If  $n_i$  be the size of the layer  $i$  and  $x$  be the layer input, then

$$g'(z^{[i]}) \approx 1$$

$$\text{Var}[z^{[i]}] = \text{Var}[x] \prod_0^{i-1} n_i \text{Var}[W^{i'}]$$

$\text{Var}[W^{i'}]$  = shared scalar variance of all weights at layer  $i'$

Then for a network with  $d$  layers

$$\text{Var}[dz^{[i]}] = \text{Var}[dz^{[d]}] \prod_i^d n_{i+1} \text{Var}[W^{i'}]$$

$$\text{Var}[dW^{[i]}] = \prod_{i'=0}^{i-1} n_{i'} \text{Var}[W^{i'}] \prod_{i'=i}^{d-1} n_{i'+1} \text{Var}[W^{i'}] \times \text{Var}[x] \text{Var}[dz^{[d]}]$$

From forward propagation point of view, to keep information flowing,

$$\forall (i, i') \text{ Var}[z^{[i]}] = \text{Var}[z^{[i']}]$$

From Back propagation point of view,

$$\forall (i, i') \text{ Var}[dz^{[i]}] = \text{Var}[dz^{[i']}]$$

These two conditions transform to

$$\forall i, n_i \text{Var}[W^{[i]}] = 1 \text{ and } n_{i+1} \text{Var}[W^{[i]}] = 1$$

Combining these two constraints, we get

$$\forall i, \text{Var}[W^{[i]}] = \frac{2}{n_i + n_{i+1}}$$

However, the variance of the backpropagated gradient may still vanish or explode in deep networks.

Thus, the normalization factor is important in deep networks because of the multiplicative effect of the layers.

In their paper, Xavier Glorot et al. Suggests that the weights be initialized(**Normalized initialization**) with weights

such that

$$W \sim U\left[\frac{-\sqrt{6}}{\sqrt{n_i + n_{i+1}}}, \frac{\sqrt{6}}{\sqrt{n_i + n_{i+1}}}\right]$$

Ref.: Understanding the difficulty of training deep feedforward neural networks

### 3. He Initialization

Xavier Intialization was mainly based on the assumption that the activation functions are linear and works good for sigmoid activation function.

Later on, ReLU surpassed Sigmoid, and solved the problem of vanishing and exploding gradient. Hence, came a new initialization procedure proposed by Kaiming He et al and is referred to as He Initialization. This initialization method allows very deep models to converge while Xavier initialization does not.

The derivation follows the procedure as given in the paper by Xavier et al.

## Forward Propagation case

For a layer,  $\mathbf{z}^{[i]} = \mathbf{W}^{[i]} \mathbf{x} + \mathbf{b}^{[i]}$

$$\mathbf{x} = f(\mathbf{z}^{[i-1]})$$

The initialized elements in  $\mathbf{W}$  are assumed to be mutually independent and share the same distribution. Furthermore, it is assumed that the elements in  $\mathbf{x}$  are also mutually independent and share the same distribution and  $\mathbf{x}$  and  $\mathbf{W}$  are independent of each other.

Then we have:

$$\text{Var}[z] = n \text{Var}[w]$$

$z, w, x$  represents random variables in  $\mathbf{z}, \mathbf{W}$  and  $\mathbf{x}$  respectively..

$$\text{Var}[z] = n \cdot \text{Var}[w] E[x^2]$$

For ReLU

$$E[x^2] = \int_{-\infty}^{+\infty} \max(0, z^{[i-1]})^2 p(z) dz$$

Since, in ReLU, the part of the ReLU for  $z < 0$  will not contribute to the integral

$$E[x^2] = \int_0^{+\infty} (z^{[i-1]})^2 p(z) dz$$

The above equation can be written as half of the integral over the entire real domain (as  $(z^{[i-1]})^2$  is symmetric around 0 and  $p(z)$  is assumed to be symmetric around 0)

$$E[x^2] = \frac{1}{2} \int_{-\infty}^{+\infty} (z^{[i-1]})^2 p(z) dz$$

Since  $z^{[i-1]}$  is assumed to have zero mean, hence,

$$E[x^2] = \frac{1}{2} \int_{-\infty}^{+\infty} (z^{[i-1]} - E[z^{[i-1]}])^2 p(z) dz = \frac{1}{2} \text{Var}[z^{[i-1]}]$$

Therefore,

$$\text{Var}[z^{[i]}] = \frac{1}{2} n \text{Var}[w^{[i]}] \text{Var}[z^{[i-1]}]$$

With  $L$  layers put together,

$$\text{Var}[z^{[L]}] = \text{Var}[y^{[1]}] \left( \prod_{i=2}^L \frac{1}{2} n^{[i]} \text{Var}[w^{[i]}] \right)$$

A proper initialization should avoid reducing or magnifying the magnitudes of input signals exponentially. So, the above product should be a scalar(e.g. 1).

A sufficient condition is

$$\frac{1}{2} n^{[i]} \text{Var}[w^{[i]}] = 1 \quad \forall i$$

This results in a zero mean Gaussian distribution with standard deviation  $\sqrt{\frac{2}{n^{[i]}}}$

This is the way of initialization as mentioned in the paper by He et al.. They also initialized  $\mathbf{b} = 0$ . and the input layer with weights  $w^{[1]}$ , such that  $z^{[1]} \text{Var}[w^{[1]}] = 1$

But the factor (1/2) does not matter if it just exists on one layer.

### Backward propagation Case.

$$da^{[i]} = W^{[i]T} dz^{[i+1]}$$

As in the forward propagation case, we consider,  $W$  and  $dz^{[i+1]}$  to be independent of each other, then  $dz^{[i]}$  has zero mean for all  $i$ , when  $w$  is initialized by a symmetric distribution around zero.

We also have,

$$dz^{[i]} = g'(z^{[i]}) da^{[i]}, \text{ where } g' \text{ is the derivative of } g$$

For ReLU, derivative is either 1 or 0, with equal probability.

Assuming  $dz^{[i]}$  and  $g'(z^{[i]})$  are independent of each other.

We have,  $E[dz^{[i]}] = E[da^{[i]}]/2 = 0$  and also  $E[(dz^{[i]})^2] = \text{Var}[dz^{[i]}] = (1/2)\text{Var}[da^{[i]}]$

$$\text{Thus, } \text{Var}[da^{[i]}] = n_i \text{Var}[w^{[i]}] \text{Var}[dz^{[i]}] = (1/2)n_i \text{Var}[w^{[i]}] \text{Var}[da^{[i]}]$$

Though the derivations are different, the factor (1/2) is the result of using ReLU.

Putting together  $L$  layers,

$$\text{Var}[\text{da}^{[2]}] = \text{Var}[\text{da}^{[L+1]}] \left( \prod_{i=2}^L \frac{1}{2} n^{[i]} \text{Var}[w^{[i]}] \right)$$

To prevent the gradient from getting exponentially large or small

$$\frac{1}{2} n^{[i]} \text{Var}[w^{[i]}] = 1, \forall i$$

The above equation results in a Gaussian Distribution with zero mean and standard

$$\text{deviation } \sqrt{\frac{2}{n^{[i]}}}$$

The main difference with the Xavier initialization is that, the He initialization addresses the rectifier nonlinearities. For very deep networks, the Xavier initialization will stall learning, while the He initialization will be able to make the model converge.

[Ref.: Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification](#)