

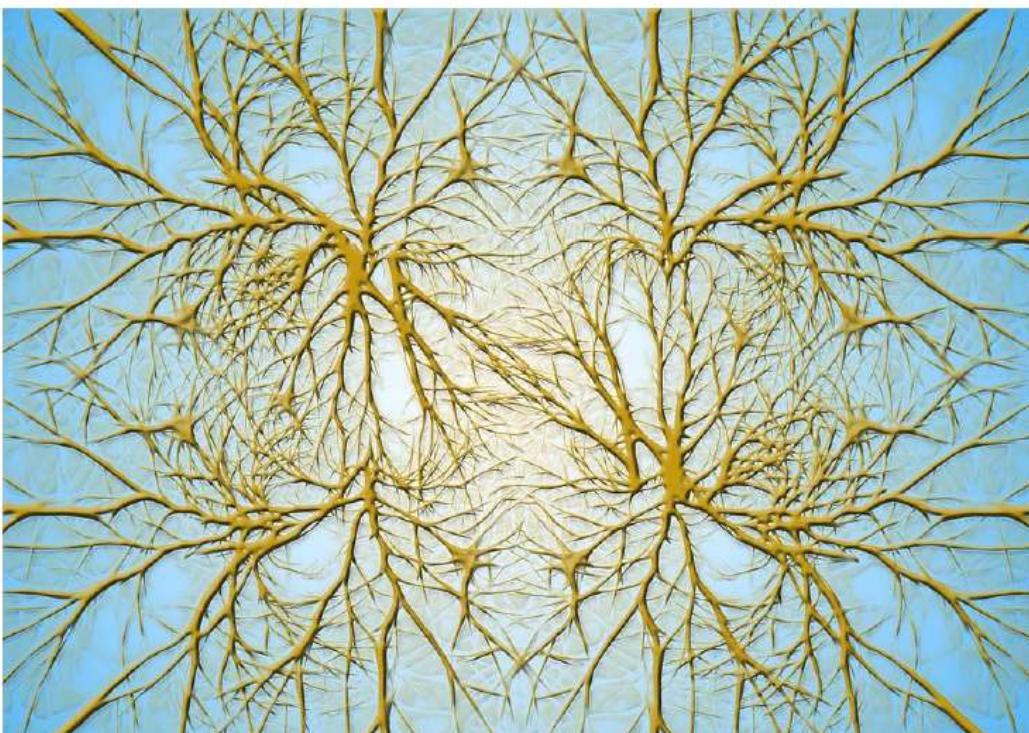
ANN : Multi Layered Perceptron (MLP)

Sourav Karmakar

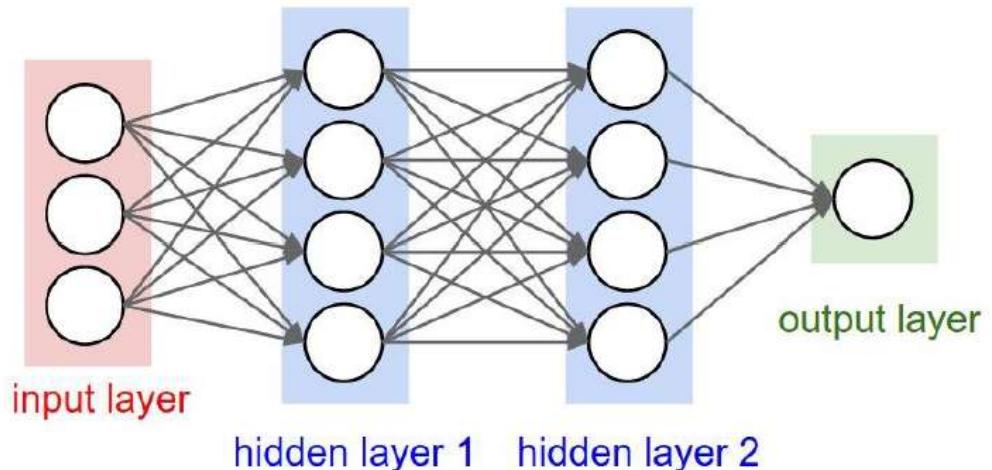
souravkarmakar29@gmail.com

Artificial Neural Networks

Biological Neurons:
Complex connectivity patterns

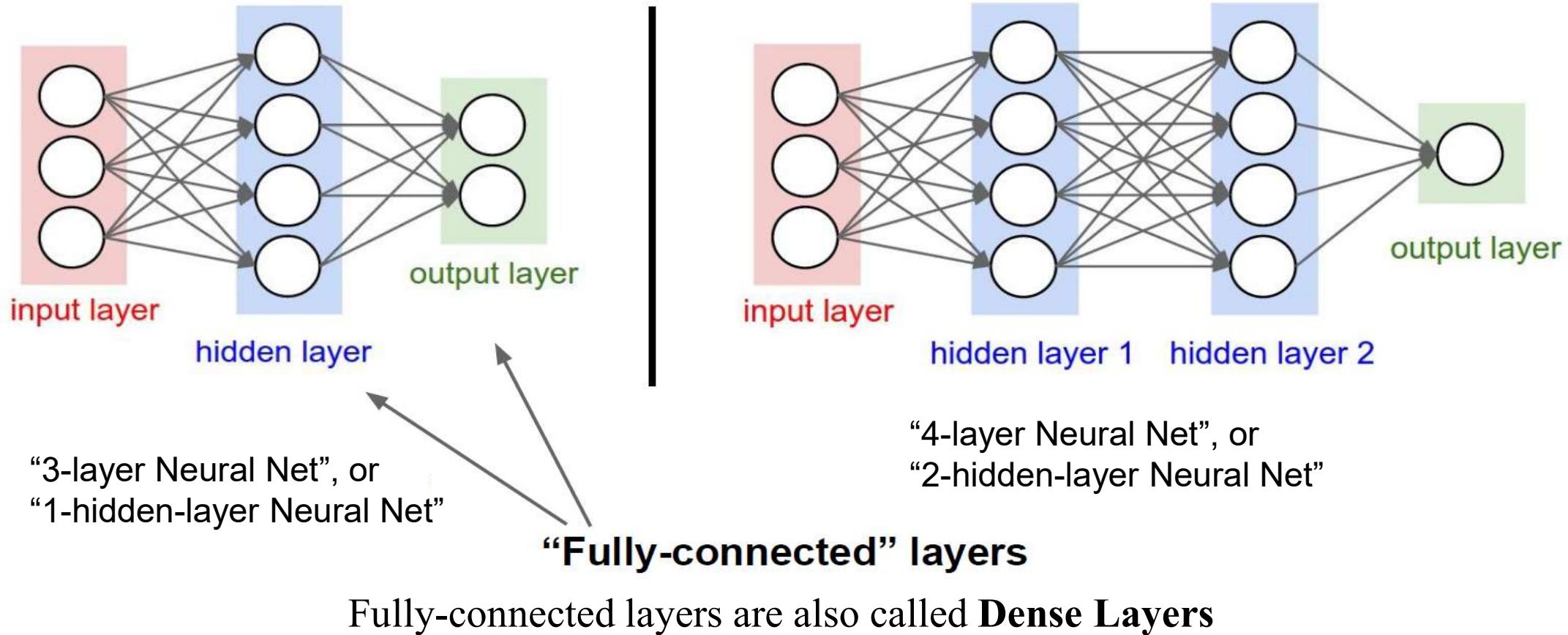


Neurons in a neural network:
Organized into regular layers for computational efficiency



This type of architectures are called **Multi-Layered Perceptron (MLP)**

Multi Layered Perceptron (MLP)



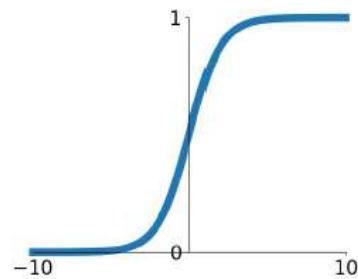
Sometimes, people don't consider input layer as a layer because it has no computation in it. In that case the first diagram will be called "2-layer Neural Net" and the second one "3-layer Neural Net".

But in our discussion in this deck, we will consider input layer as a separate non-computing layer.

Different Activation Functions

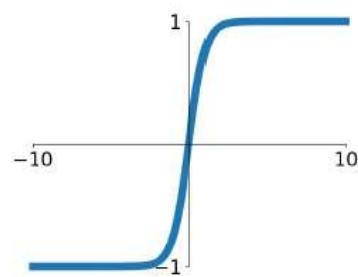
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



tanh

$$\tanh(x)$$

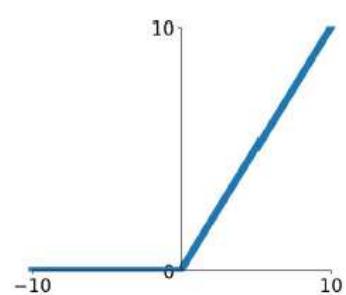


Rectified Linear Unit

ReLU

$$\max(0, x)$$

if $x > 0$, $\text{ReLU}(x) = x$
 $x \leq 0$, $\text{ReLU}(x) = 0$

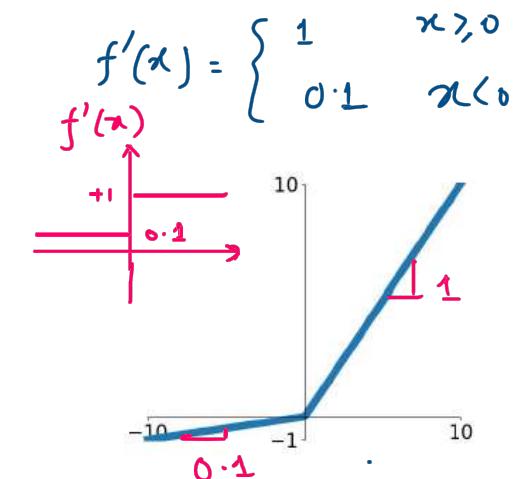


ReLU is a good practical choice for most problems

Leaky ReLU

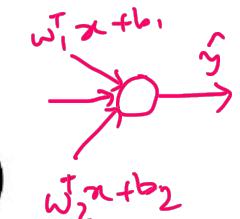
$$\max(0.1x, x)$$

if $x \geq 0$; $f(x) = x$
 $x < 0$; $f(x) = 0.1x$



Maxout (Rarely Used)

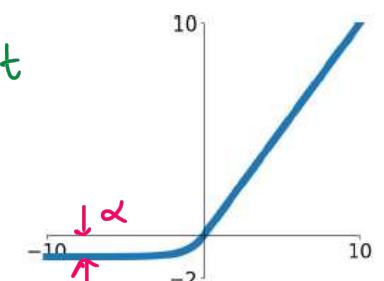
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$



Exponential Linear Unit

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Control parameter : saturation value
 $\text{as } x \rightarrow -\infty \quad \therefore \text{ELU}(x) = -\alpha$

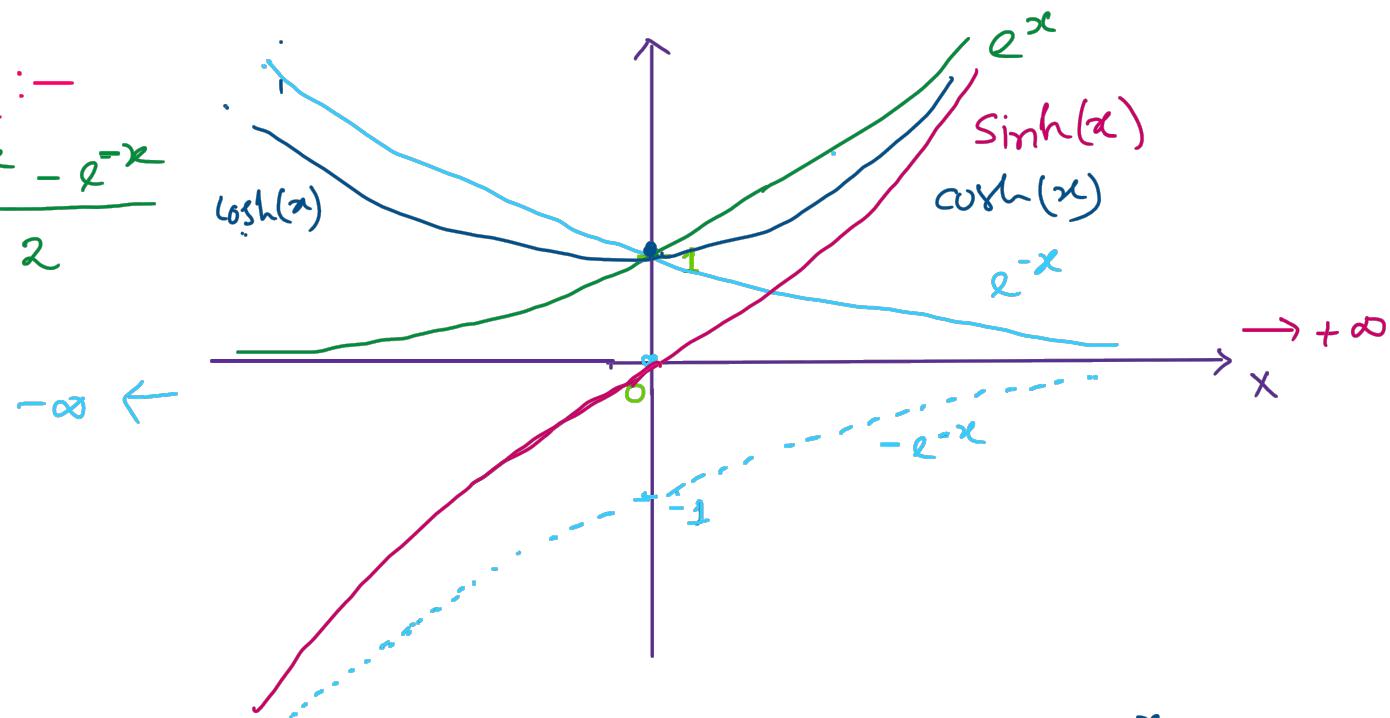
Hyperbolic Functions:-

$$\sinh(x) = \frac{e^x - e^{-x}}{2}$$

$$\sinh(0) = 0$$

$$\cosh(x) = \frac{e^x + e^{-x}}{2}$$

$$\cosh(0) = 1$$



$$\tanh(x) = \frac{\sinh(x)}{\cosh(x)} = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Range of $\tanh(x)$ is -1 to $+1$

$$\begin{aligned} x \rightarrow \infty & \quad \tanh(x) = +1 \\ x \rightarrow -\infty & ; \quad \tanh(x) = -1 \end{aligned}$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{-x}(e^x - e^{-x})}{e^{-x}(e^x + e^{-x})} = \frac{1 - e^{-2x}}{1 + e^{-2x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\begin{aligned}\frac{d}{dx} [\tanh(x)] &= \frac{(e^{2x}+1) \cdot (2e^{2x}) - (e^{2x}-1) \cdot (2e^{2x})}{(e^{2x}+1)^2} = \frac{\cancel{2e^{2x}} \cdot e^{2x} + 2e^{2x} - \cancel{2e^{2x}} \cdot e^{2x} + \cancel{2e^{2x}}}{(e^{2x}+1)^2} \\ &= \frac{4 \cdot e^{2x} \cdot 1}{(e^{2x}+1)^2} \Leftrightarrow \frac{(e^{2x}+1)^2 - (e^{2x}-1)^2}{(e^{2x}+1)^2} = 1 - \left(\frac{e^{2x}-1}{e^{2x}+1}\right)^2\end{aligned}$$

$$(a+b)^2 - (a-b)^2 = 4ab.$$

$$\frac{d}{dx} (\tanh(x)) = 1 - (\tanh(x))^2$$

$$\frac{d}{dx} (\sigma(x)) = \sigma(x) \cdot (1 - \sigma(x))$$

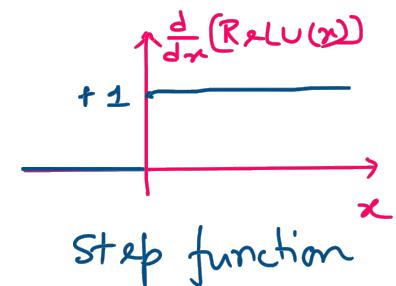
$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \in [-1, +1]$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

$$\left. \frac{d}{dx} \tanh(x) \right|_{x=0} = 1$$

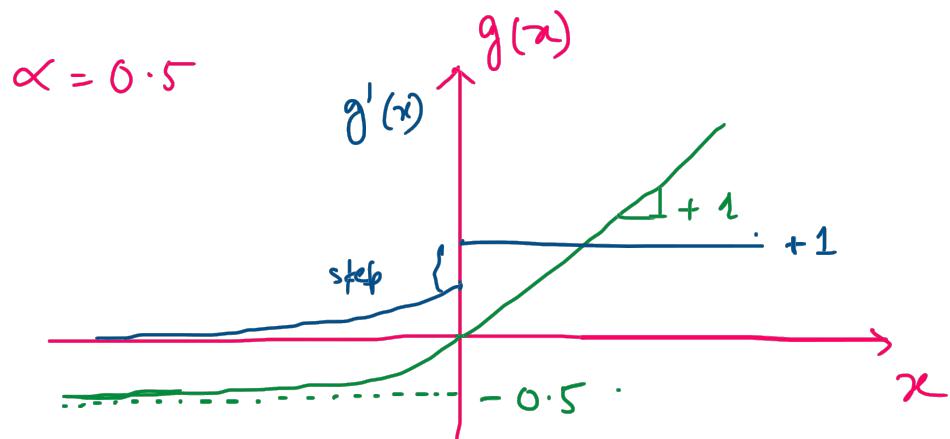
$$\text{ReLU}(x) = \begin{cases} 0 & ; x \leq 0 \\ x & ; x > 0 \end{cases} \Rightarrow \frac{d}{dx} [\text{ReLU}(x)] = \begin{cases} 0 & ; x \leq 0 \\ 1 & ; x > 0 \end{cases}$$



$$g(x) = \begin{cases} x & ; x \geq 0 \\ \alpha(e^x - 1) & ; x < 0 \end{cases}$$

$$g'(x) = \begin{cases} 1 & ; x \geq 0 \\ \alpha e^x & ; x < 0 \end{cases}$$

$x \rightarrow 0^-$
 $g'(x) = \alpha$



GELU(x) \rightarrow Gaussian Error Linear Unit

$$\text{GELU}(x) = x \cdot \Phi(x)$$

$\Phi(x)$ = CDF of standard normal distribution

The Role of Activation Functions

Without Activation Function:

- Suppose we only use **linear operations** in a network: $y = Wx + b$
- Even if we stack multiple layers, the output is still just a **linear transformation** of the input.
$$y = W_3(W_2(W_1x + b_1) + b_2) + b_3 = W'x + b'$$
- That means no matter how many layers you add, the whole network is just a linear model. Hence, without activation function the model can't capture the complex non-linear relationship in the data.

Purpose of Activation Function:

1. Introduce Non-linearity:

- Helps capture complex patterns in data by introducing non-linearity.
- With nonlinear activations, neural networks become **universal function approximators** (proved by the **Universal Approximation Theorem**)

2. Control Information Flow:

- Decide which neurons "fire" (activate) and which don't.
- Makes the network selective about what information passes forward.

3. Bounded Outputs / Normalization:

- Some activations (like sigmoid, tanh) squash outputs into a small range, making training stable.

4. Gradient Flow:

Differentiable activation functions allow gradient-based learning (backpropagation)

Softmax Regression

Softmax Regression is a generalization of Logistic Regression for multiclass classification problems.

- **Training Set:** $\{ \langle \mathbf{x}^{(1)}, y^{(1)} \rangle, \langle \mathbf{x}^{(2)}, y^{(2)} \rangle, \dots, \langle \mathbf{x}^{(m)}, y^{(m)} \rangle \}$,

Where m = Number of training data, the feature vector $\mathbf{x}^{(i)} \in \mathbb{R}^{n+1}$ ($x_0 = 1$) , and the labels $y^{(i)} \in \{1, 2, 3, \dots, K\}$ (K = Number of classes).

- **Hypothesis:**
$$h_{\theta}(\mathbf{x}^{(i)}) = \frac{1}{\sum_{j=1}^K e^{\boldsymbol{\theta}_j^T \mathbf{x}^{(i)}}} \begin{bmatrix} e^{\boldsymbol{\theta}_1^T \mathbf{x}^{(i)}} \\ e^{\boldsymbol{\theta}_2^T \mathbf{x}^{(i)}} \\ \vdots \\ e^{\boldsymbol{\theta}_K^T \mathbf{x}^{(i)}} \end{bmatrix}$$

Where $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_K \in \mathbb{R}^{n+1}$, are the parameters of our model. The denominator $\sum_{j=1}^K e^{\boldsymbol{\theta}_j^T \mathbf{x}^{(i)}}$ normalizes the distribution, so that it sums to one.

Usually Softmax regression is used at the output / classification layer of the ANN. Where it is also referred to as **Softmax Activation function**.

(logits) $\leftarrow z_1 = \omega_{11}h_1 + \omega_{12}h_2 + \omega_{13}h_3 + b_1$
 $\omega_1 \left\{ \begin{array}{c} \xrightarrow{\quad} \\ \xrightarrow{\quad} \\ \xrightarrow{\quad} \end{array} \right. z_1$
 $\omega_2 \left\{ \begin{array}{c} \xrightarrow{\quad} \\ \xrightarrow{\quad} \\ \xrightarrow{\quad} \end{array} \right. z_2 = \omega_{21}h_1 + \omega_{22}h_2 + \omega_{23}h_3 + b_2$
 $\omega_3 \left\{ \begin{array}{c} \xrightarrow{\quad} \\ \xrightarrow{\quad} \\ \xrightarrow{\quad} \end{array} \right. z_3 = \omega_{31}h_1 + \omega_{32}h_2 + \omega_{33}h_3 + b_3$

logits: $z_1 = -1, z_2 = 2, z_3 = 1$

$$\hat{y}_1 = \frac{e^{-1}}{e^{-1} + e^2 + e^1} = 0.035$$

$$\hat{y}_2 = \frac{e^2}{e^{-1} + e^2 + e^1} = 0.70$$

$$\hat{y}_3 = \frac{e^1}{e^{-1} + e^2 + e^1} = 0.265$$

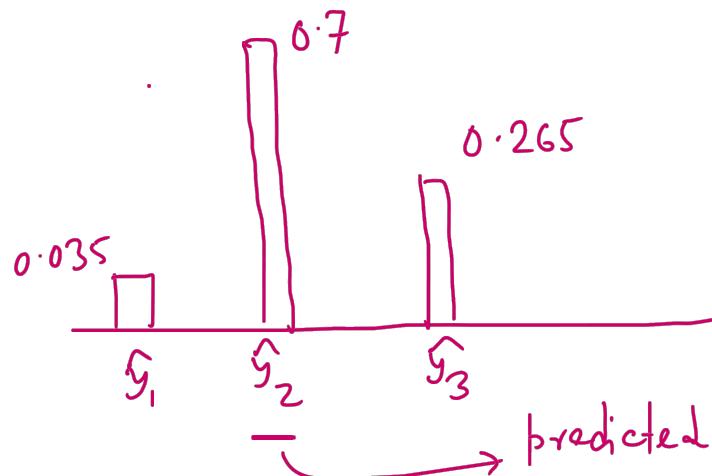
σ
 $+ \cdot \cdot \cdot$
 m
 a
 x

$$\hat{y}_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\hat{y}_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\left[\frac{0}{0.035}, \frac{1}{0.7}, \frac{2}{0.265} \right]$$



$\text{predicted class} = \arg \max (\hat{y}_1, \hat{y}_2, \hat{y}_3)$
 $= 2$

0	1	2
↑	↑	↑

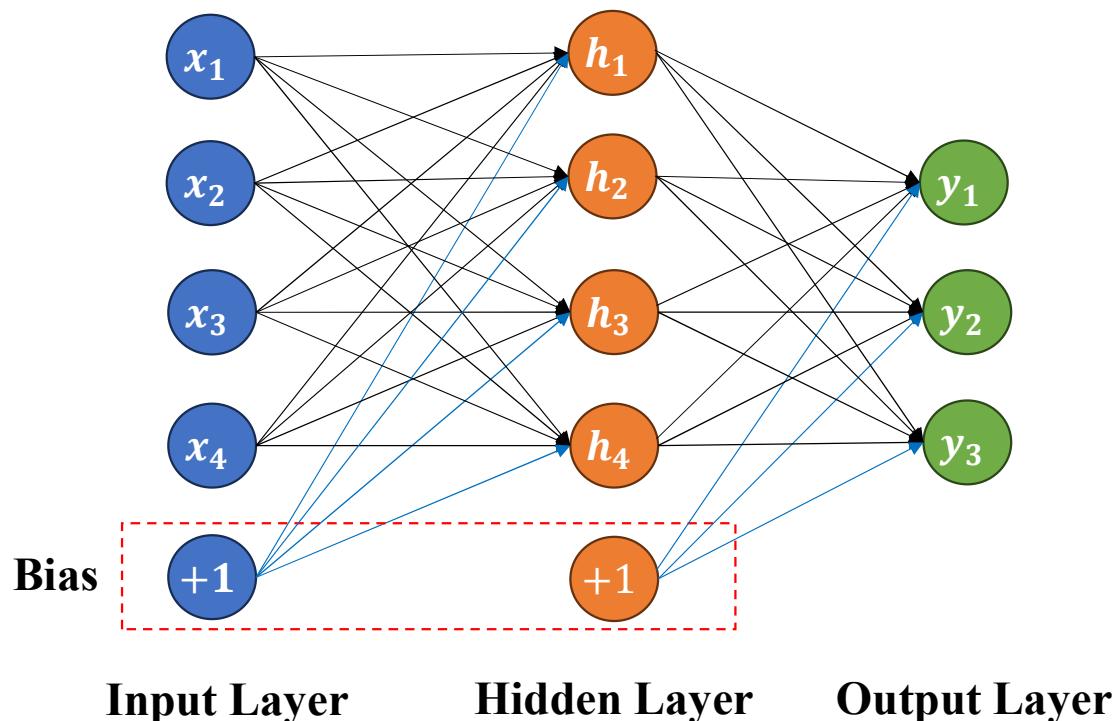
Multi Layered Perceptron (MLP): Architecture

- A layer of sensory nodes called input nodes and the layer is called input layer. No computation happens in input nodes.
- The number of nodes in the input layer = number of features or attributes of the data.
- Layers of computing nodes. These are called **hidden layers**. In successive layers, the nodes are all connected by learnable weights. This is called ‘dense architecture’.
- Number of Hidden layers and Number of nodes in each hidden layer are user choice (hyper-parameters) and there is no general guideline for the choice of number of hidden layers or the number of nodes in each hidden layers. However, there should be at least 1 hidden layer.
- Layer of final computing nodes is called output layer.
- The number of nodes in the output layer = number of classes of the dataset. (for a classification problem). For multi-class classification it is usually called ‘softmax’ layer due to the softmax regression. For binary classification only one node with sigmoid activation is needed. For regression problem, there will be only one output node without any activation function.
- Usually we use ‘**ReLU**’ activation for all nodes in hidden layer and ‘**softmax**’ activation in output layer for multi-class classification problem or ‘**sigmoid**’ activation in output layer for binary classification problem. However we can also use different activation functions in the hidden layer.

Multi Layered Perceptron (MLP): Architecture

For example: To classify the simple IRIS dataset using ANN

- The number of input nodes = 4 (as the dataset has four features / attributes)
- The number of output nodes = 3 (as there are three classes)
- We can build an ANN with one hidden layer containing 4 nodes (say)



- Number of connections from input layer to hidden layer:
$$(4+1) \times 4 = 20$$

(4 + 1) Bias Number of nodes in hidden layer
- Number of connections from hidden layer to output layer:
$$(4+1) \times 3 = 15$$

(4 + 1) Bias Number of nodes in output layer

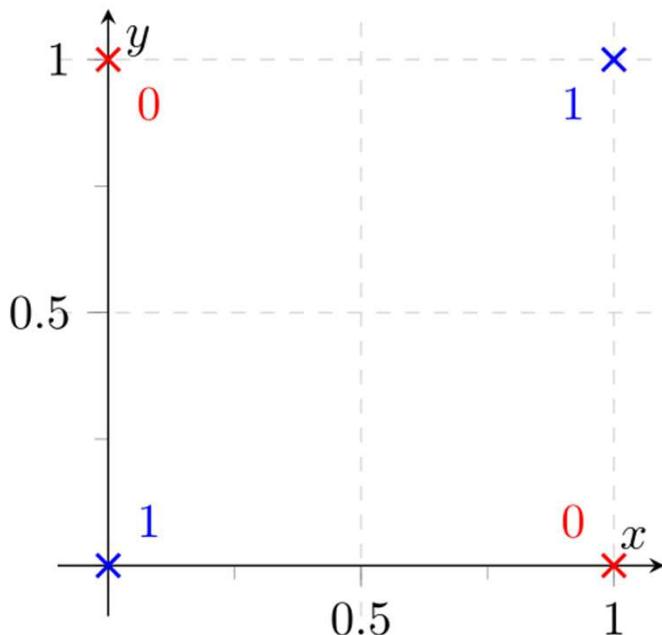
Multi Layered Perceptron (MLP): Advantages

- **Non-linearity:** MLP is very suitable to model a physical phenomenon which is in general non-linear.
- **Automatic Feature Extraction:** Each layer in MLP learns a representation of the data. The network learns features automatically during training by adjusting weights. The learned features are optimal for final task as MLP is an universal function approximator. The first few hidden layers learns low level patterns, whereas the deeper layers learn high-level representation. This is called **representation learning**.
- **Adaptivity:** It can cope with the change in dataset distribution.
- **Parallelism:** It can be implemented in multi-core parallel architecture (like GPU) for faster computation.
- **Robustness:** It has inherent ability to handle confusing / noisy / missing datapoints
- **Fault tolerance:** It has the ability to work to some extent even in case of component (neuron) failure.
- **Complex decision boundary:** It has the ability to have a complex decision boundary for classification.

We shall use MLP with supervised learning for complex classification tasks.

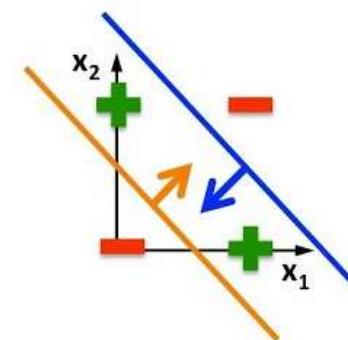
How MLP Solved The XOR Problem

XOR Problem

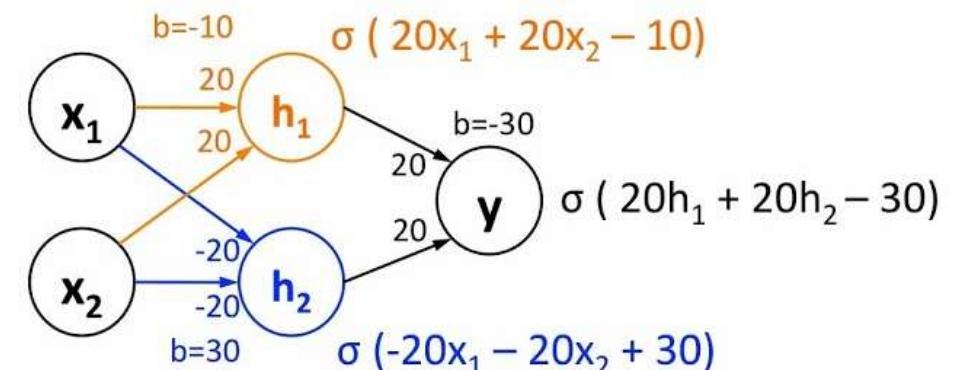


Input 1	Input 2	Output
0	0	0
0	1	1
1	1	0
1	0	1

Linear classifiers
cannot solve this



Multi Layered Perceptron



$$\sigma(20 \cdot 0 + 20 \cdot 0 - 10) \approx 0$$

$$\sigma(20 \cdot 1 + 20 \cdot 1 - 10) \approx 1$$

$$\sigma(20 \cdot 0 + 20 \cdot 1 - 10) \approx 1$$

$$\sigma(20 \cdot 1 + 20 \cdot 0 - 10) \approx 1$$

$$\sigma(-20 \cdot 0 - 20 \cdot 0 + 30) \approx 1$$

$$\sigma(-20 \cdot 1 - 20 \cdot 1 + 30) \approx 0$$

$$\sigma(-20 \cdot 0 - 20 \cdot 1 + 30) \approx 1$$

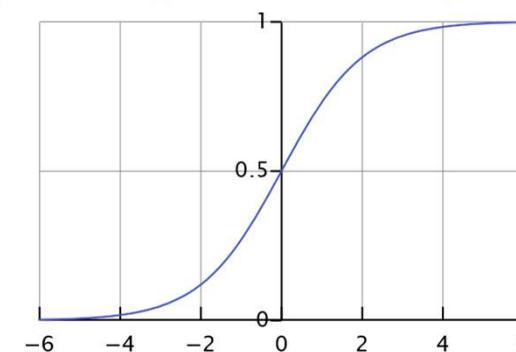
$$\sigma(-20 \cdot 1 - 20 \cdot 0 + 30) \approx 1$$

$$\sigma(20 \cdot 0 + 20 \cdot 1 - 30) \approx 0$$

$$\sigma(20 \cdot 1 + 20 \cdot 0 - 30) \approx 0$$

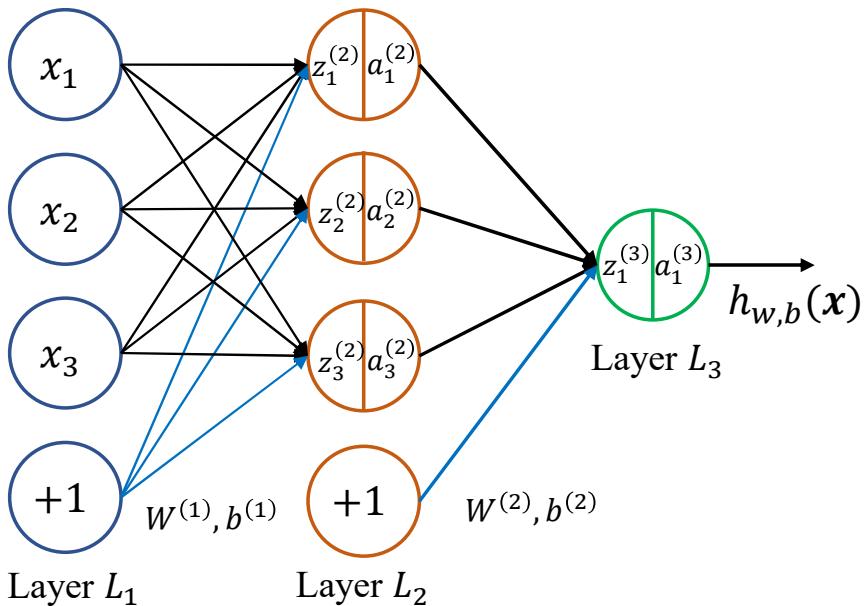
$$\sigma(20 \cdot 1 + 20 \cdot 1 - 30) \approx 1$$

$$\sigma(20 \cdot 1 + 20 \cdot 1 - 30) \approx 1$$



Multi Layered Perceptron (MLP)

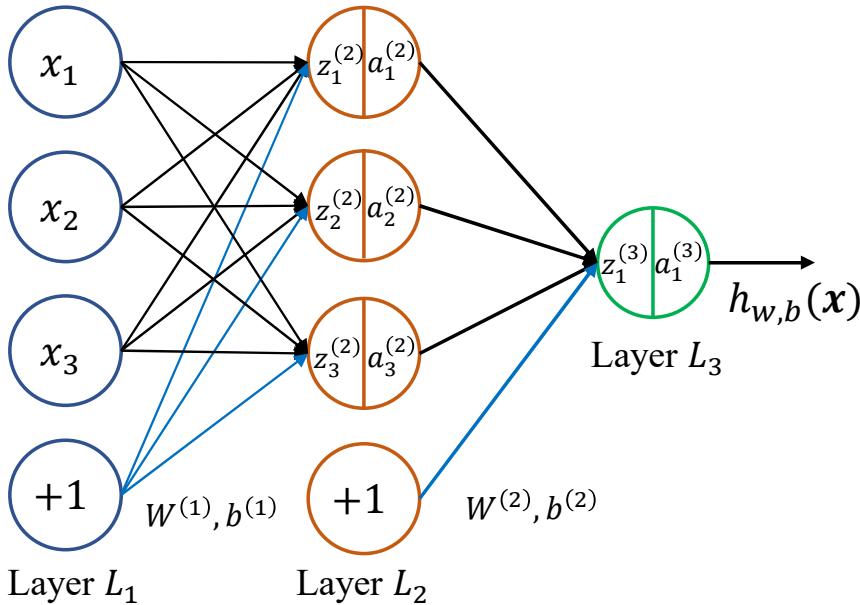
Key Notations:



- n_l denotes the number of layers in our network; thus in the figure beside $n_l = 3$.
 - s_l denotes the number of nodes in layer l (not counting the bias unit)
 - L_1 is the input layer and L_{n_l} is the output layer.
 - $b_i^{(l)}$ is the bias associated with unit i in layer $l + 1$
-
- The neural network as shown above has parameters $\langle W, b \rangle = \langle (W^{(1)}, b^{(1)}), (W^{(2)}, b^{(2)}) \rangle$, Where $W_{ij}^{(l)}$ denotes the parameter (or weight) associated with the connection between unit j in layer l , and unit i in layer $l + 1$. Note that bias units don't have inputs or connections going into them, since they always output value $+1$.
 - In this example we have $W^{(1)} \in \mathbb{R}^{3 \times 3}$, $b^{(1)} \in \mathbb{R}^{3 \times 1}$, $W^{(2)} \in \mathbb{R}^{1 \times 3}$, and $b^{(2)} \in \mathbb{R}$

Multi Layered Perceptron (MLP)

Forward Propagation:



- \$z_i^{(l)}\$ is the pre-activation of unit \$i\$ in layer \$l\$. This is the weighted summation of inputs connected to that unit.
- \$a_i^{(l)}\$ to denote the post activation output of unit \$i\$ in layer \$l\$. For \$l = 1\$, we also use \$a_i^{(1)} = x_i\$ to denote \$i^{th}\$ input.

$$a_1^{(2)} = f \left(W_{11}^{(1)} x_1 + W_{12}^{(1)} x_2 + W_{13}^{(1)} x_3 + b_1^{(1)} \right) = f \left(z_1^{(2)} \right)$$

$$a_2^{(2)} = f \left(W_{21}^{(1)} x_1 + W_{22}^{(1)} x_2 + W_{13}^{(1)} x_3 + b_2^{(1)} \right) = f \left(z_2^{(2)} \right)$$

$$a_3^{(2)} = f \left(W_{31}^{(1)} x_1 + W_{32}^{(1)} x_2 + W_{33}^{(1)} x_3 + b_3^{(1)} \right) = f \left(z_3^{(2)} \right)$$

Here \$f(\cdot)\$ denotes the activation function.

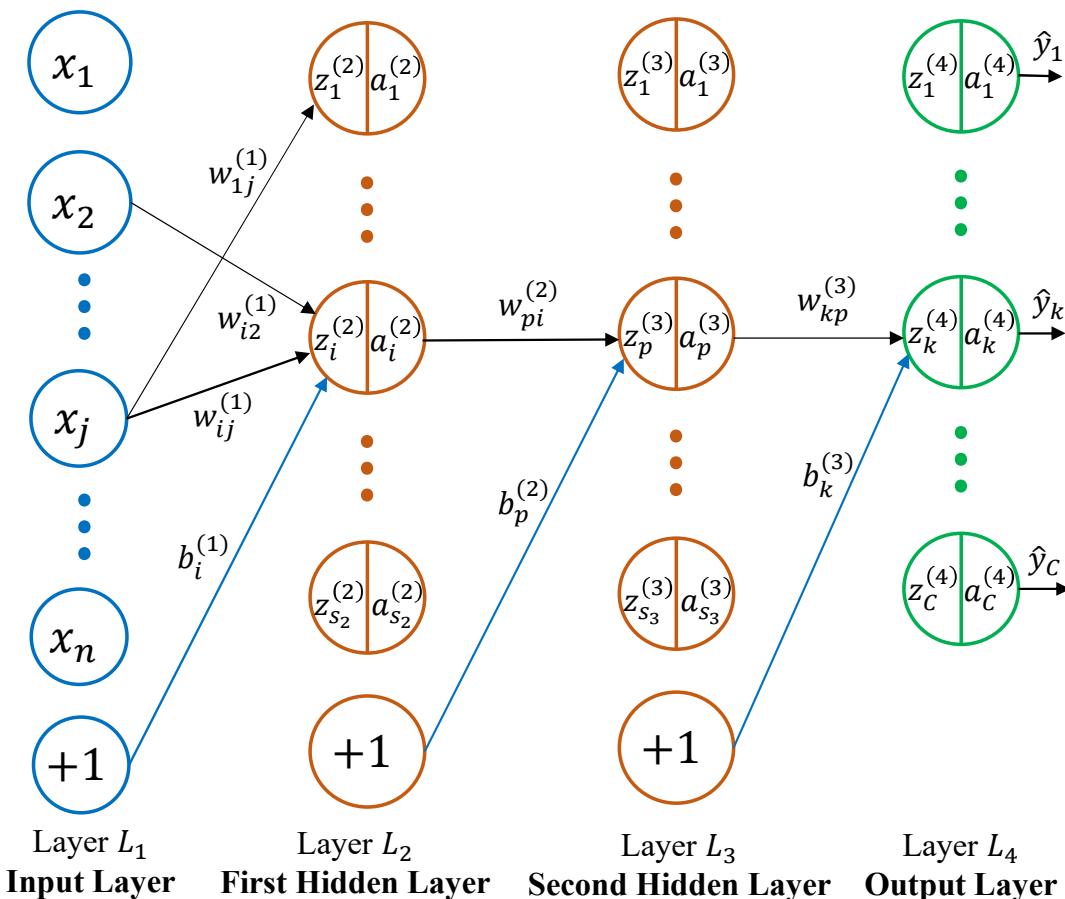
$$\text{The final output } h_{w,b}(x) = a_1^{(3)} = f \left(W_{11}^{(2)} a_1^{(2)} + W_{12}^{(2)} a_2^{(2)} + W_{13}^{(2)} a_3^{(2)} + b_1^{(2)} \right) = f \left(z_1^{(3)} \right)$$

All the equations shown above can be written in simplified fashion in vector matrix notation:

$$\mathbf{a}^{(l+1)} = f(\mathbf{z}^{(l+1)}) = f(W^{(l)} \mathbf{a}^{(l)} + \mathbf{b}^{(l)})$$

Multi Layered Perceptron (MLP)

- n_l denotes the number of layers in the network. s_l denotes the number of nodes in layer l .
- $W_{ij}^{(l)}$ denotes the parameter (or weight) associated with the connection between unit j in layer l , and unit i in layer $l + 1$.
- $b_i^{(l)}$ is the bias associated with unit i in layer $l + 1$.



$$\begin{bmatrix} z_1^{(2)} \\ z_i^{(2)} \\ \vdots \\ z_{s_2}^{(2)} \end{bmatrix}_{s_2 \times 1} = \begin{bmatrix} w_{11}^{(1)} & \dots & w_{1n}^{(1)} \\ \vdots & \ddots & \vdots \\ w_{i1}^{(1)} & \dots & w_{in}^{(1)} \\ \vdots & \ddots & \dots \\ w_{s_21}^{(1)} & \dots & w_{s_2n}^{(1)} \end{bmatrix}_{s_2 \times n} \begin{bmatrix} x_1 \\ x_j \\ \vdots \\ x_n \end{bmatrix}_{n \times 1} + \begin{bmatrix} b_1^{(1)} \\ b_i^{(1)} \\ \vdots \\ b_{s_2}^{(1)} \end{bmatrix}_{s_2 \times 1} \text{ and } \begin{bmatrix} a_1^{(2)} \\ a_i^{(2)} \\ \vdots \\ a_{s_2}^{(2)} \end{bmatrix}_{s_2 \times 1} = f \left(\begin{bmatrix} z_1^{(2)} \\ z_i^{(2)} \\ \vdots \\ z_{s_2}^{(2)} \end{bmatrix}_{s_2 \times 1} \right)$$

$$z^{(2)} = W^{(1)}x + \mathbf{b}^{(1)} \text{ and } \mathbf{a}^{(2)} = f(z^{(2)})$$

$$\begin{bmatrix} z_1^{(3)} \\ z_p^{(3)} \\ \vdots \\ z_{s_3}^{(3)} \end{bmatrix}_{s_3 \times 1} = \begin{bmatrix} w_{11}^{(2)} & \dots & w_{1s_2}^{(2)} \\ \vdots & \ddots & \vdots \\ w_{p1}^{(2)} & \dots & w_{ps_2}^{(2)} \\ \vdots & \ddots & \dots \\ w_{s_31}^{(2)} & \dots & w_{s_3s_2}^{(2)} \end{bmatrix}_{s_3 \times s_2} \begin{bmatrix} a_1^{(2)} \\ a_i^{(2)} \\ \vdots \\ a_{s_2}^{(2)} \end{bmatrix}_{s_2 \times 1} + \begin{bmatrix} b_1^{(2)} \\ b_p^{(2)} \\ \vdots \\ b_{s_3}^{(2)} \end{bmatrix}_{s_3 \times 1} \text{ and } \begin{bmatrix} a_1^{(3)} \\ a_p^{(3)} \\ \vdots \\ a_{s_3}^{(3)} \end{bmatrix}_{s_3 \times 1} = f \left(\begin{bmatrix} z_1^{(3)} \\ z_p^{(3)} \\ \vdots \\ z_{s_3}^{(3)} \end{bmatrix}_{s_3 \times 1} \right)$$

$$z^{(3)} = W^{(2)}\mathbf{a}^{(2)} + \mathbf{b}^{(2)} \text{ and } \mathbf{a}^{(3)} = f(z^{(3)})$$

$$\begin{bmatrix} z_1^{(4)} \\ z_k^{(4)} \\ \vdots \\ z_c^{(4)} \end{bmatrix}_{C \times 1} = \begin{bmatrix} w_{11}^{(3)} & \dots & w_{1s_3}^{(3)} \\ \vdots & \ddots & \vdots \\ w_{k1}^{(3)} & \dots & w_{ks_3}^{(3)} \\ \vdots & \ddots & \dots \\ w_{c1}^{(3)} & \dots & w_{cs_3}^{(3)} \end{bmatrix}_{C \times s_3} \begin{bmatrix} a_1^{(3)} \\ a_p^{(3)} \\ \vdots \\ a_{s_3}^{(3)} \end{bmatrix}_{s_3 \times 1} + \begin{bmatrix} b_1^{(3)} \\ b_k^{(3)} \\ \vdots \\ b_c^{(3)} \end{bmatrix}_{C \times 1} \text{ and } \begin{bmatrix} a_1^{(4)} \\ a_k^{(4)} \\ \vdots \\ a_c^{(4)} \end{bmatrix}_{C \times 1} = f \left(\begin{bmatrix} z_1^{(4)} \\ z_k^{(4)} \\ \vdots \\ z_c^{(4)} \end{bmatrix}_{C \times 1} \right)$$

$$z^{(4)} = W^{(3)}\mathbf{a}^{(3)} + \mathbf{b}^{(3)} \text{ and } \mathbf{a}^{(4)} = f(z^{(4)})$$

The predicted output, $\hat{\mathbf{y}} = \mathbf{a}^{(4)}$

Multi Layered Perceptron (MLP)

Loss / Cost Function:

Given a training set of ‘m’ samples, we can define the overall cost functions in following way:

- **Mean Squared Error Loss:** Suitable for regression tasks.

$$J(W, b) = \frac{1}{m} \left[\sum_{i=1}^m \frac{1}{2} (h_{W,b}(x^{(i)}) - y^{(i)})^2 \right]$$

- **Binary Cross-entropy Loss:** Suitable for binary classification.

$$J(W, b) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log(h_{W,b}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{W,b}(x^{(i)})) \right]$$

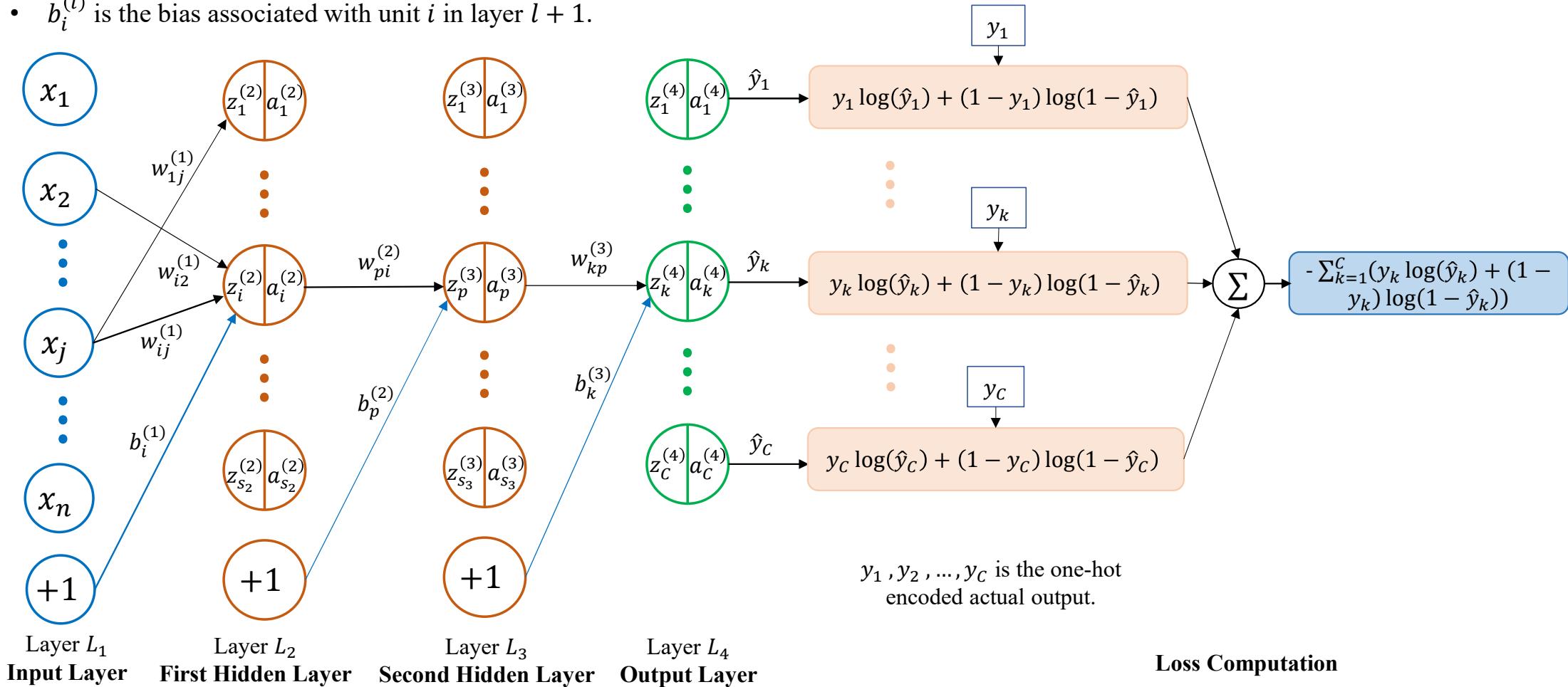
- **Categorical Cross-entropy Loss:** Suitable for multi-class classification (number of classes = C).

$$J(W, b) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^C y_k^{(i)} \log(h_{W,b}(x^{(i)})_k) + (1 - y_k^{(i)}) \log(1 - h_{W,b}(x^{(i)})_k) \right]$$

$h_{W,b}(x^{(i)})$ is the model’s prediction for input $x^{(i)}$. So, $h_{W,b}(x^{(i)}) = \hat{y}^{(i)}$ and the actual label is $y^{(i)}$.

Multi Layered Perceptron (MLP)

- n_l denotes the number of layers in the network. s_l denotes the number of nodes in layer l .
- $W_{ij}^{(l)}$ denotes the parameter (or weight) associated with the connection between unit j in layer l , and unit i in layer $l + 1$.
- $b_i^{(l)}$ is the bias associated with unit i in layer $l + 1$.



y_1, y_2, \dots, y_C is the one-hot
encoded actual output.

Loss Computation

Multi Layered Perceptron (MLP)

Back Propagation Learning Algorithm:

The weights and biases in the neural network are updated using Backpropagation Algorithm, where each weights are updated using Gradient Descent update rule (or some variant of it). For that we have to calculate the partial derivatives of the loss function with respect to the weights and biases.

User has to specify:

- Number of **layers**: n_l , or people often specify number of hidden layers which is $n_l - 2$.
- Number of **nodes** in each hidden layer: s_l for $l = 2, 3, \dots, n_l - 1$
- Number of **epochs** for which the network to be trained. An epoch = one complete pass of the entire dataset through the neural network. During an epoch, the model sees all training samples once and update weights depending on the **batch size**.
- The **loss function** to be used by the neural network.
- The optimization technique (also called **optimizer**) to be used to optimize the loss. The most basic optimizer is gradient descent.
- The **batch size**: The number of training samples processed before the model updates its weights.

Gradient Descent Optimization Technique

General idea of Gradient Descent:

We want to minimize a loss function $J(W, b)$ with respect to model parameters W and b .

Update rule: (showing only for w for simplification)

$$w \leftarrow w - \eta \nabla_w J ; \quad \eta \Rightarrow \text{learning rate}, \nabla J \Rightarrow \text{gradient of loss } J \text{ w.r.t. parameter } w.$$

(a) Batch Gradient Descent:

- In each epoch the error in the final layer is average of errors due to all the training samples in the dataset.
- Update the parameters / weights using backpropagation after all the training samples are fed to the network.
- This is stable and achieves smooth convergence. However, this is very slow for large dataset.

(b) Stochastic Gradient Descent (SGD):

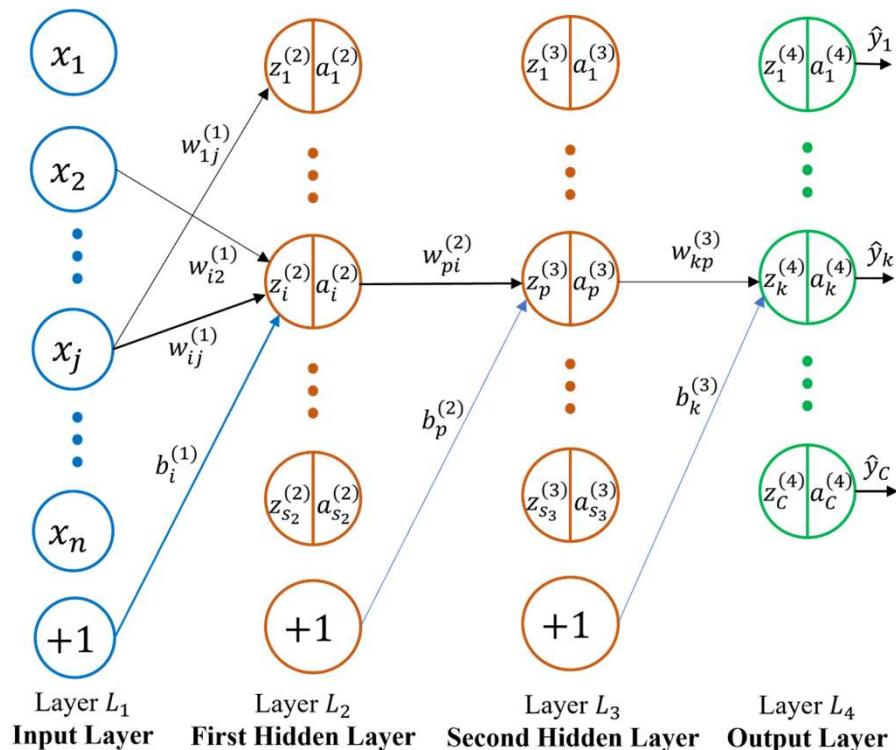
- Updates the weight after each training sample.
- Faster convergence but noisy updates due to randomness.
- It can escape local minima due to randomness.

(c) Mini-batch Gradient Descent:

- Uses small random subsets (batches) of data to compute gradients. Compromise between batch and SGD.
- For example: if dataset size = 5000 and batch size = 100, then one epoch = 50 batches. Each batch produces one weight update, so 50 iterations per epoch.
- It is most widely used in practice. The batch size depends on dataset size, GPU/RAM memory, and training stability. Most commonly used batch sizes are 32, 64, 128, 256 etc. [i.e. between 32 to 256]

Back Propagation Learning Algorithm

To better understand the flow of computations, let us draw a computation graph

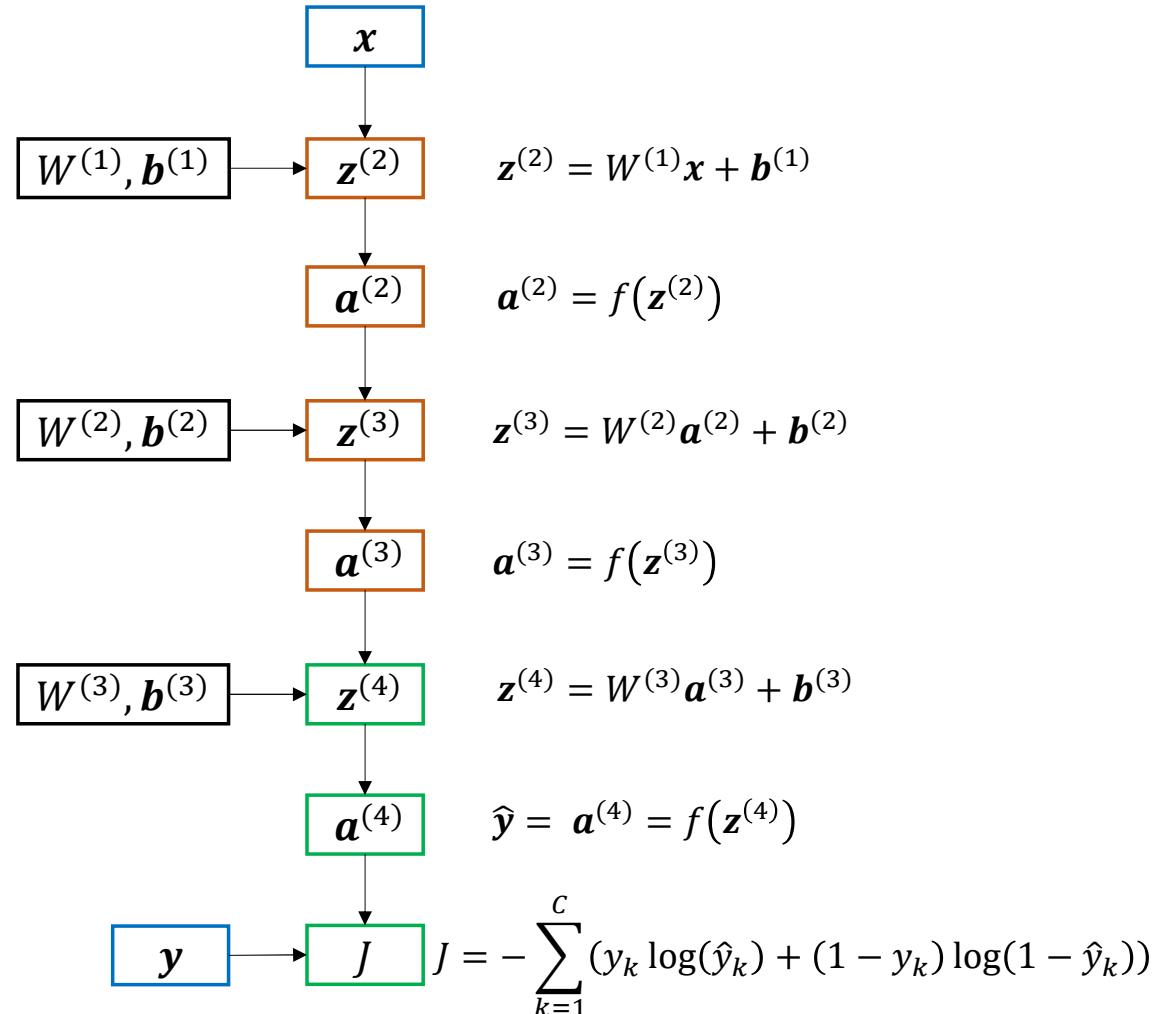


$$\mathbf{z}^{(2)} = W^{(1)}\mathbf{x} + \mathbf{b}^{(1)} \text{ and } \mathbf{a}^{(2)} = f(\mathbf{z}^{(2)})$$

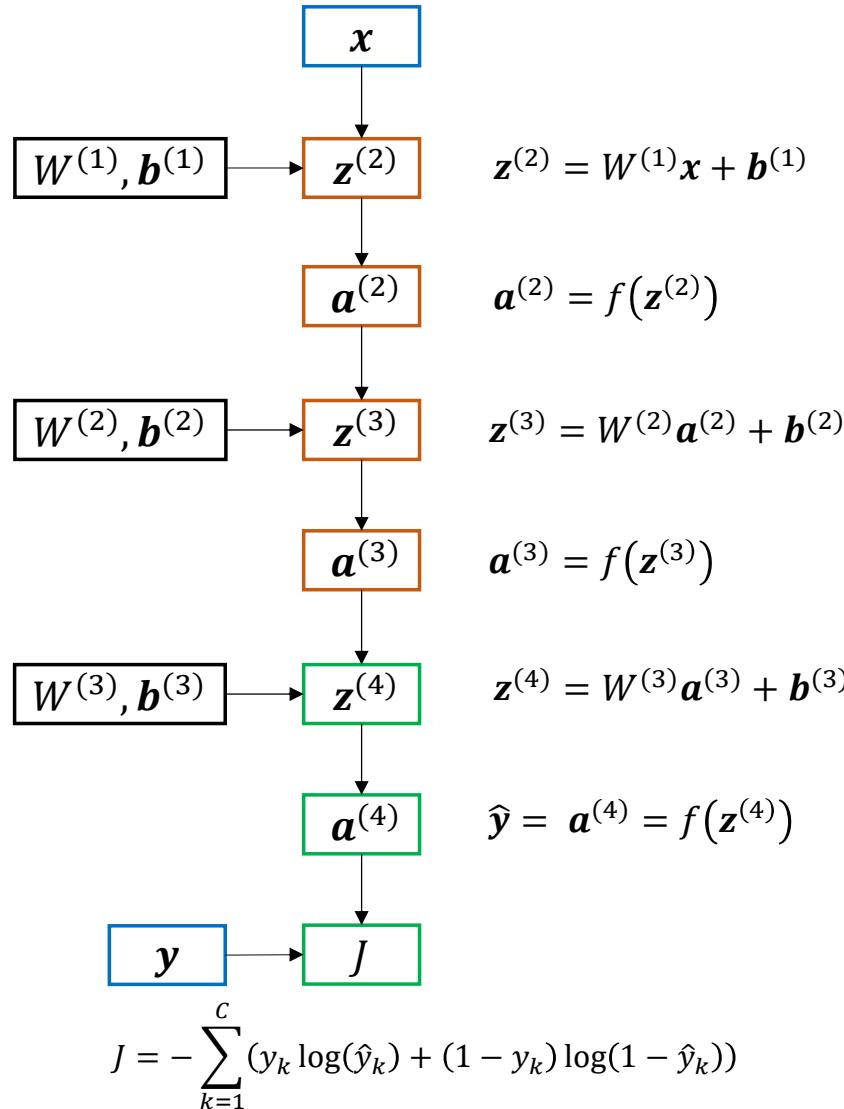
$$\mathbf{z}^{(3)} = W^{(2)}\mathbf{a}^{(2)} + \mathbf{b}^{(2)} \text{ and } \mathbf{a}^{(3)} = f(\mathbf{z}^{(3)})$$

$$\mathbf{z}^{(4)} = W^{(3)}\mathbf{a}^{(3)} + \mathbf{b}^{(3)} \text{ and } \mathbf{a}^{(4)} = f(\mathbf{z}^{(4)})$$

The predicted output, $\hat{\mathbf{y}} = \mathbf{a}^{(4)}$



Back Propagation Learning Algorithm



Cost function J is a function of $a^{(4)}$, $a^{(4)}$ is in turn a function of $z^{(4)}$. $z^{(4)}$ is function of $W^{(3)}, b^{(3)}$ and $a^{(3)}$.

So, to calculate the partial derivatives of the cost function J w.r.t. $W^{(3)}, b^{(3)}$ we will have to use the **chain rule**.

$$\frac{\partial J}{\partial W^{(3)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial W^{(3)}} \quad \text{and} \quad \frac{\partial J}{\partial b^{(3)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial b^{(3)}}$$

Similarly, to calculate the partial derivatives of the cost function J w.r.t. $W^{(2)}, b^{(2)}$ we will use the following formula:

$$\frac{\partial J}{\partial W^{(2)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial W^{(2)}}$$

$$\text{and} \quad \frac{\partial J}{\partial b^{(2)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial b^{(2)}}$$

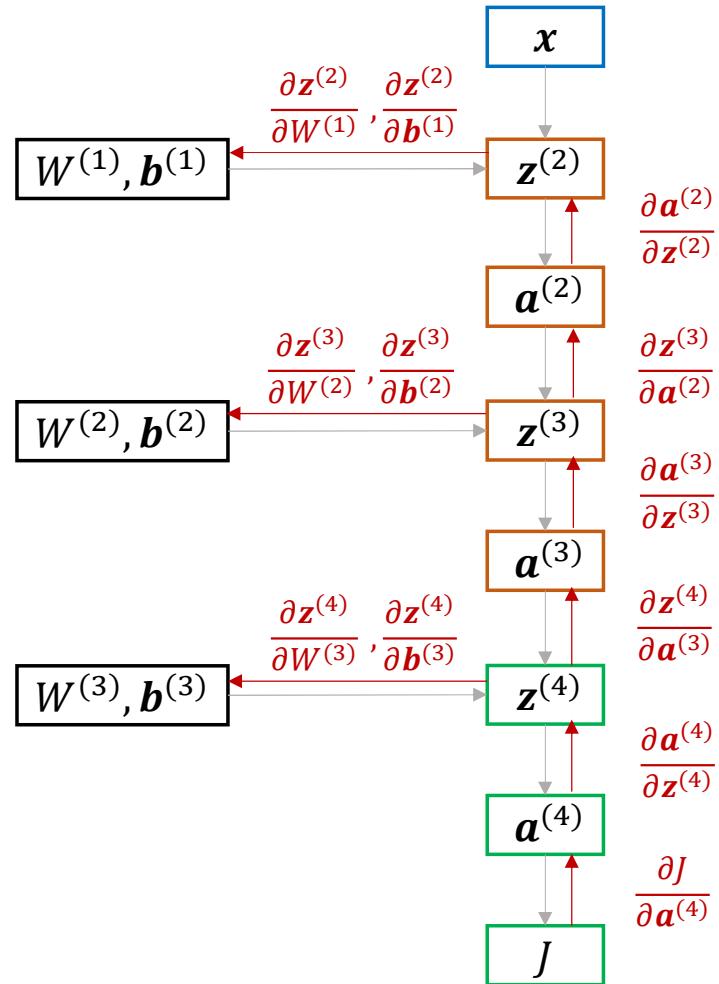
Likewise, for calculating the partial derivatives of the cost function J w.r.t. $W^{(1)}, b^{(1)}$ we will use:

$$\frac{\partial J}{\partial W^{(1)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial W^{(1)}}$$

$$\text{and} \quad \frac{\partial J}{\partial b^{(1)}} = \frac{\partial J}{\partial a^{(4)}} \cdot \frac{\partial a^{(4)}}{\partial z^{(4)}} \cdot \frac{\partial z^{(4)}}{\partial a^{(3)}} \cdot \frac{\partial a^{(3)}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial b^{(1)}}$$

Back Propagation Learning Algorithm

This can be better understood from the following figure:



$$z^{(2)} = W^{(1)}x + b^{(1)}$$

$$a^{(2)} = f(z^{(2)})$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)}$$

$$a^{(3)} = f(z^{(3)})$$

$$z^{(4)} = W^{(3)}a^{(3)} + b^{(3)}$$

$$\hat{y} = a^{(4)} = f(z^{(4)})$$

Forward propagation

Back propagation

- It's called **backpropagation** because the error (loss gradient) is **propagated backward** from the output layer to the input layer.
- This backward flow of gradients updates weights layer by layer using the chain rule.

Back Propagation Learning Algorithm

Input to the algorithm: Labelled dataset: $\langle X, y \rangle$; **Output of the algorithm:** Trained weights and biases

1. For each epoch : (Number of epochs should be specified by user)
2. Perform a feedforward pass, computing the activations for layers L_2, L_3 , and so on up to the output layer L_{n_l} .
3. For each output unit i in layer n_l (the final/output layer), set

$$\delta_i^{(n_l)} = \frac{\partial J}{\partial a_i^{(n_l)}} \cdot f'(z_i^{(n_l)})$$

4. For $l = (n_l - 1), (n_l - 2), (n_l - 3), \dots, 2$

For each node i in layer l , set

$$\delta_i^{(l)} = \left(\sum_{j=1}^{s_{l+1}} W_{ji}^{(l)} \delta_j^{(l+1)} \right) \cdot f'(z_i^{(l)})$$

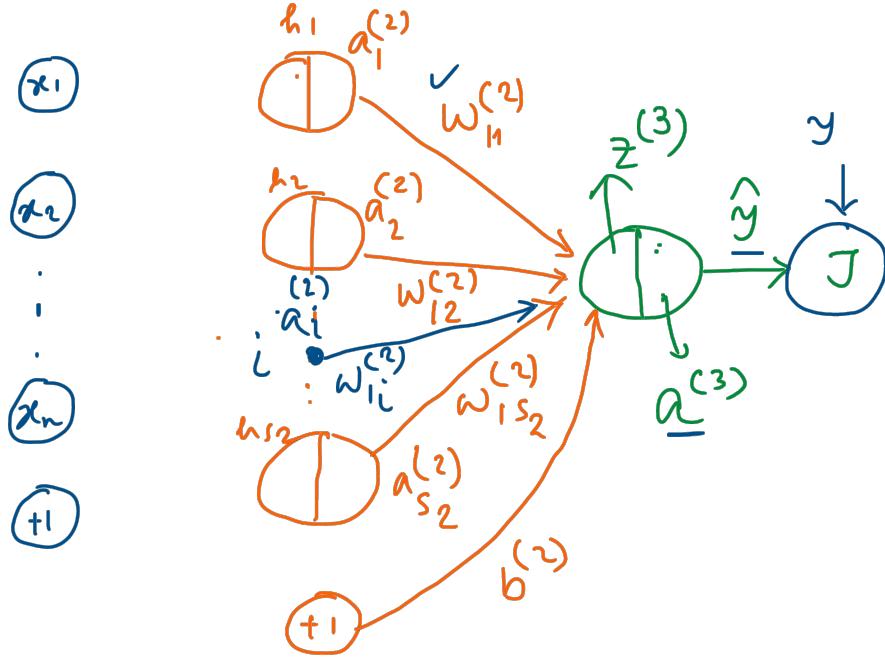
5. Compute the desired partial derivatives which are given as:

$$\frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) = a_j^{(l)} \cdot \delta_i^{(l+1)} \quad \text{and} \quad \frac{\partial}{\partial b_i^{(l)}} J(W, b) = \delta_i^{(l+1)}$$

6. Finally update the parameters:

$$W_{ij}^{(l)} \leftarrow W_{ij}^{(l)} - \eta \frac{\partial}{\partial W_{ij}^{(l)}} J(W, b) \quad \text{and} \quad b_i^{(l)} \leftarrow b_i^{(l)} - \eta \frac{\partial}{\partial b_i^{(l)}} J(W, b); \quad \eta \text{ is the learning rate}$$

Back-Propagation:



$$z^{(3)} = w_{11}^{(2)} a_1^{(2)} + \dots + w_{1s_2}^{(2)} a_{s_2}^{(2)} + b^{(2)} = \sum_{i=1}^{s_2} w_{1i}^{(2)} a_i^{(2)} + b^{(2)}$$

$$a^{(3)} = \hat{y} = f(z^{(3)}) \quad [f(\cdot) \rightarrow \text{activation function}]$$

$$J = - \left(y \log \hat{y} + (1-y) \log (1-\hat{y}) \right)$$

$$\frac{\partial J}{\partial w_{1i}^{(2)}} = \frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial w_{1i}^{(2)}}$$

$$= - \left(\frac{y}{\hat{y}} - \frac{1-y}{1-\hat{y}} \right) \cdot \frac{\partial f(z^{(3)})}{\partial z^{(3)}} \cdot \boxed{\frac{\partial z^{(3)}}{\partial w_{1i}^{(2)}}} \rightarrow a_i^{(2)}$$

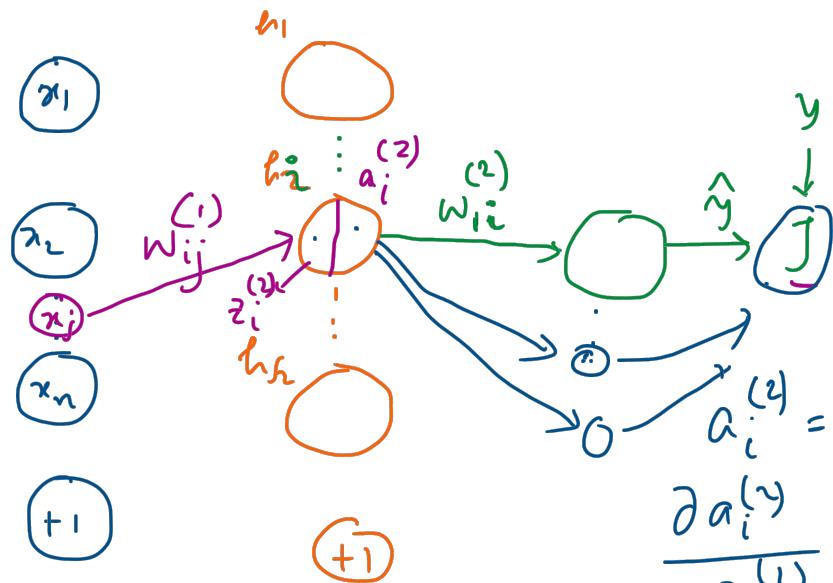
$$\sigma'(x) = - \frac{y(1-\hat{y}) - \hat{y}(1-y)}{\hat{y}(1-\hat{y})} \cdot f'(z^{(3)}) \cdot a_i^{(2)}$$

$$\frac{\partial J}{\partial b^{(2)}} = - \delta_i^{(3)} \quad (\text{Step } 5)$$

$$= - \frac{(y - \hat{y})}{\hat{y}(1-\hat{y})} \cdot f'(z^{(3)}) \cdot a_i^{(2)}$$

$$a_i^{(2)} = \boxed{- (y - \hat{y}) \cdot f'(z^{(3)}) \cdot a_i^{(2)}}$$

$$\frac{\partial J}{\partial w_{1i}^{(2)}} = - \underline{\delta_i^{(3)}} \cdot a_i^{(2)} \rightarrow (\text{Step } 5)$$



$$\frac{\partial J}{\partial w_{ij}^{(1)}} = \boxed{\frac{\partial J}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial z^{(3)}} \cdot \frac{\partial z^{(3)}}{\partial a_i^{(2)}} \cdot \frac{\partial a_i^{(2)}}{\partial z_i^{(2)}} \cdot \frac{\partial z_i^{(2)}}{\partial w_{ij}^{(1)}}}$$

$\delta_i^{(3)} \downarrow$ $w_{li}^{(2)}$ $f'(z_i^{(2)}) \cdot a_j^{(1)}$

$$\sum \delta_i^{(3)} w_{li}^{(2)} \cdot f'(z_i^{(2)}) \cdot a_j^{(1)}$$

Step - 1

$$\bar{z}^{(3)} = w_{11}^{(2)} a_1^{(2)} + w_{12}^{(2)} a_2^{(2)} + \dots + w_{ii}^{(2)} a_i^{(2)} + \dots + b^{(2)}$$

$$\frac{\partial \bar{z}^{(3)}}{\partial a_i^{(2)}} = w_{ii}^{(2)}$$

$$z_i^{(2)} = w_{11}^{(1)} x_1 + w_{12}^{(1)} x_2 + \dots + w_{ij}^{(1)} \cdot \underline{x_j} + \dots + b_i^{(1)}$$

$$\frac{\partial z_i^{(2)}}{\partial w_{ij}^{(1)}} = x_j = a_j^{(1)}$$

Back-Prop has a sequential dependency layerwise, but within a layer the weight updates can be performed parallelly.

Thank You