



Deep Learning CS60010

Abir Das

Assistant Professor

Computer Science and Engineering Department
Indian Institute of Technology Kharagpur

<http://cse.iitkgp.ac.in/~adas/>

Biological Neural Network

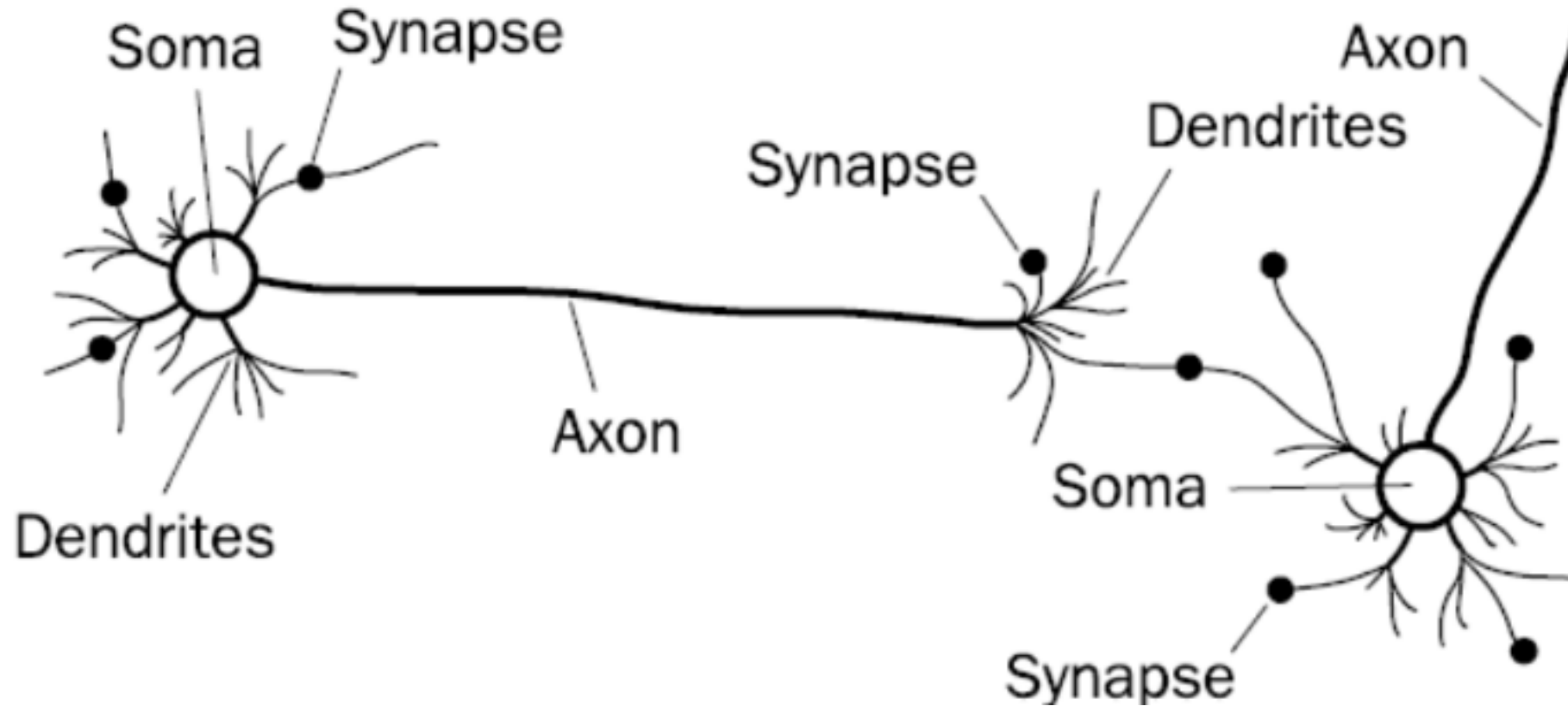


Image courtesy: F. A. Makinde et. al., "Prediction of crude oil viscosity using feed-forward back-propagation neural network (FFBPNN)". Petroleum & Coal 2012

deep CNN features (~2013) vs. IT neural "features"

Dear Computer Vision Researchers,

Thank you !

We are working toward the same goal....


See great related work by:
Kriegeskorte,
Oliva,
Gallant & Malik,
Gardner

Ability to predict IT neurons
(% variance explained)

Feature Type	Ability to predict IT neurons (% variance explained)
HMO features	~45
AlexNet features	~55
Zeiler & Fergus features	~58

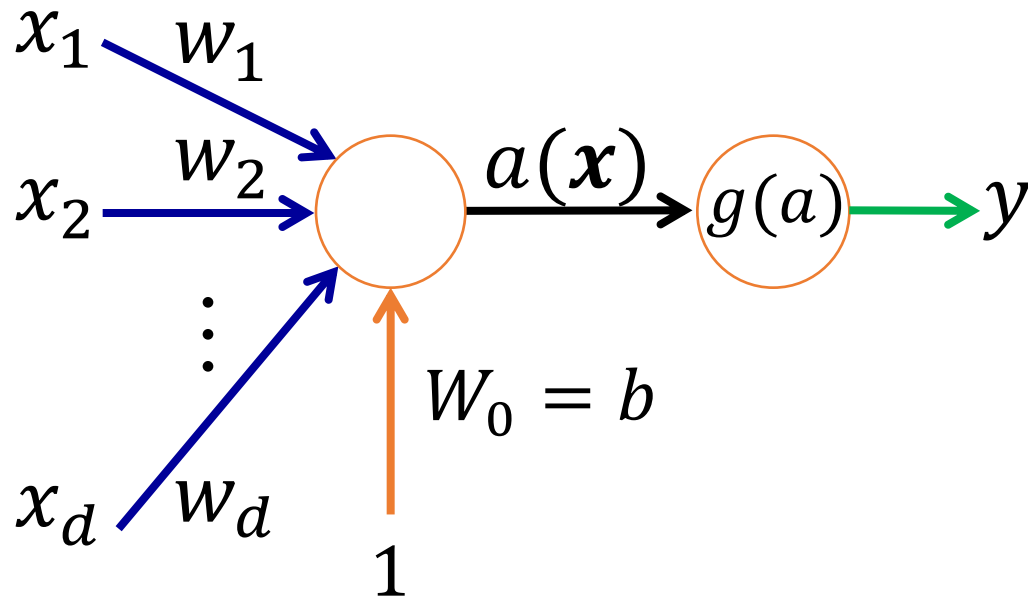
Zeiler & Fergus features
AlexNet features

N. Majaj N. and DiCarlo J.J. *PCLR* (2013);
N. Majaj N. and DiCarlo J.J. *PLoS Comp Bio* (2014)



CVPR 2017
July 21-26
HONOLULU

Artificial Neuron



$$\mathbf{w} = [w_1 \ w_2 \ \dots \ w_d]^T \quad \text{and} \quad \mathbf{x} = [x_1 \ x_2 \ \dots \ x_d]^T$$

$$a(\mathbf{x}) = b + \sum_{i=1}^d w_i x_i = [\mathbf{w}^T \ b] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$y = g(a(\mathbf{x}))$$

Terminologies:-

\mathbf{x} : input, \mathbf{w} : weights, \mathbf{b} : bias

a : pre-activation (input activation)

g : activation function

y : activation (output activation)

Perceptron

The New York Times

Electronic 'Brain' Teaches Itself

JULY 13, 1958

The Navy last week demonstrated the embryo of an electronic computer named the Perceptron which, when completed in about a year, is expected to be the first non-living mechanism able to "perceive, recognize and identify its surroundings without human training or control." Navy officers demonstrating a preliminary form of the device in Washington said they hesitated to call it a machine because it is so much like a "human being without life."

Dr. Frank Rosenblatt, research psychologist at the Cornell Aeronautical Laboratory, Inc., Buffalo, N. Y., designer of the Perceptron, conducted the demonstration. The machine, he said, would be the first electronic device to think as the

recognize the difference between right and left, almost the way a child learns.

When fully developed, the Perceptron will be designed to remember images and information it has perceived itself, whereas ordinary computers remember only what is fed into them on punch cards or magnetic tape.

Later Perceptrons, Dr. Rosenblatt said, will be able to recognize people and call out their names. Printed pages, longhand letters and even speech commands are within its reach. Only one more step of development, a difficult step, he said, is needed for the device to hear speech in one language and instantly translate it to speech or writing in another language.

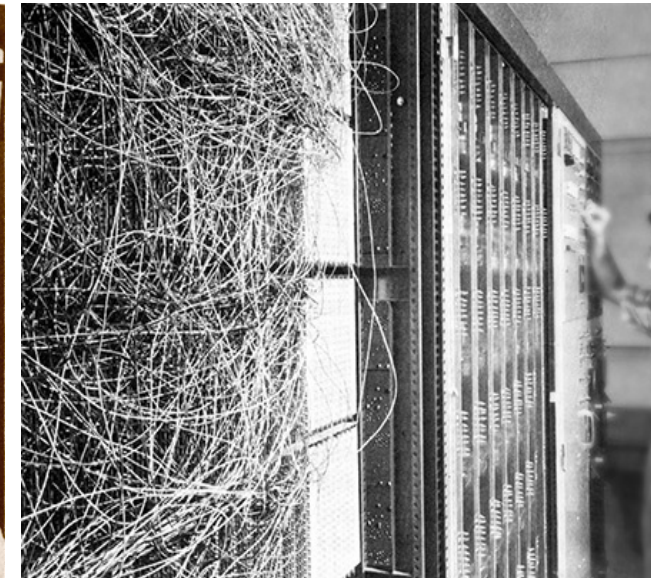


Image courtesy: <https://blogs.umass.edu/comphon/2017/06/15/did-frank-roosenblatt-invent-deep-learning-in-1962/>

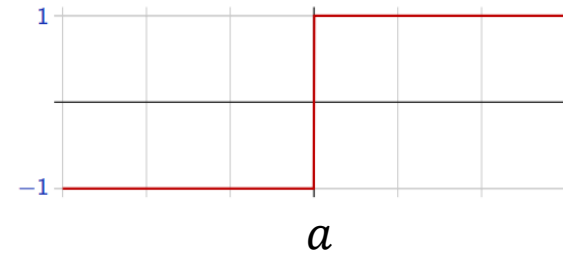
Perceptron

$\mathbf{x} \in \mathbb{R}^d$ and $y \in \{0, 1\}$ – Binary Classification

$$g(a) = \begin{cases} 1, & a \geq 0 \\ 0, & a < 0 \end{cases} \quad (\text{Rosenblatt, 1957})$$

To make things simpler, the response is taken as $y \in \{-1, 1\}$

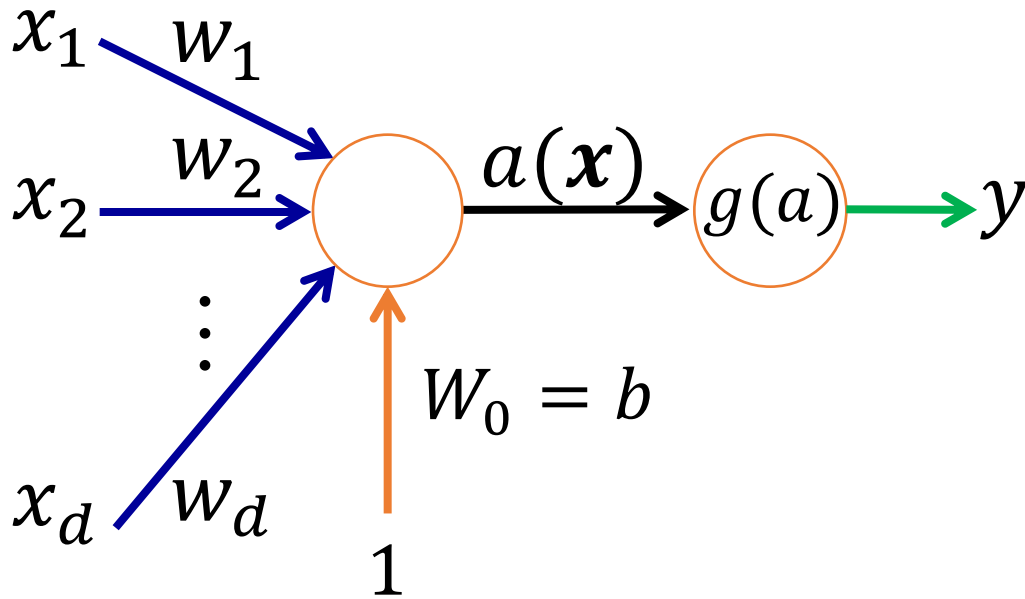
$$g(a) = \begin{cases} 1, & a \geq 0 \\ -1, & a < 0 \end{cases}$$



The perceptron classification rule, thus, translates to

$$y = \begin{cases} 1, & \mathbf{w}^T \mathbf{x} + b \geq 0 \\ -1, & \mathbf{w}^T \mathbf{x} + b < 0 \end{cases}$$

Remember: $\mathbf{w}^T \mathbf{x} + b = 0$ represents a hyperplane.



Perceptron Learning Algorithm

Given a training set: $\{\mathbf{x}^{(i)}, y^{(i)} : \mathbf{x}^{(i)} \in \mathbb{R}^d, y^{(i)} \in \{-1, 1\}\}$ for $i = 1, 2, \dots, N$

1. Start with $k = 1$ and $\mathbf{w}^{(k)} = [0 \ 0 \ \dots \ 0]^T$
2. Loop until all examples are correctly classified:

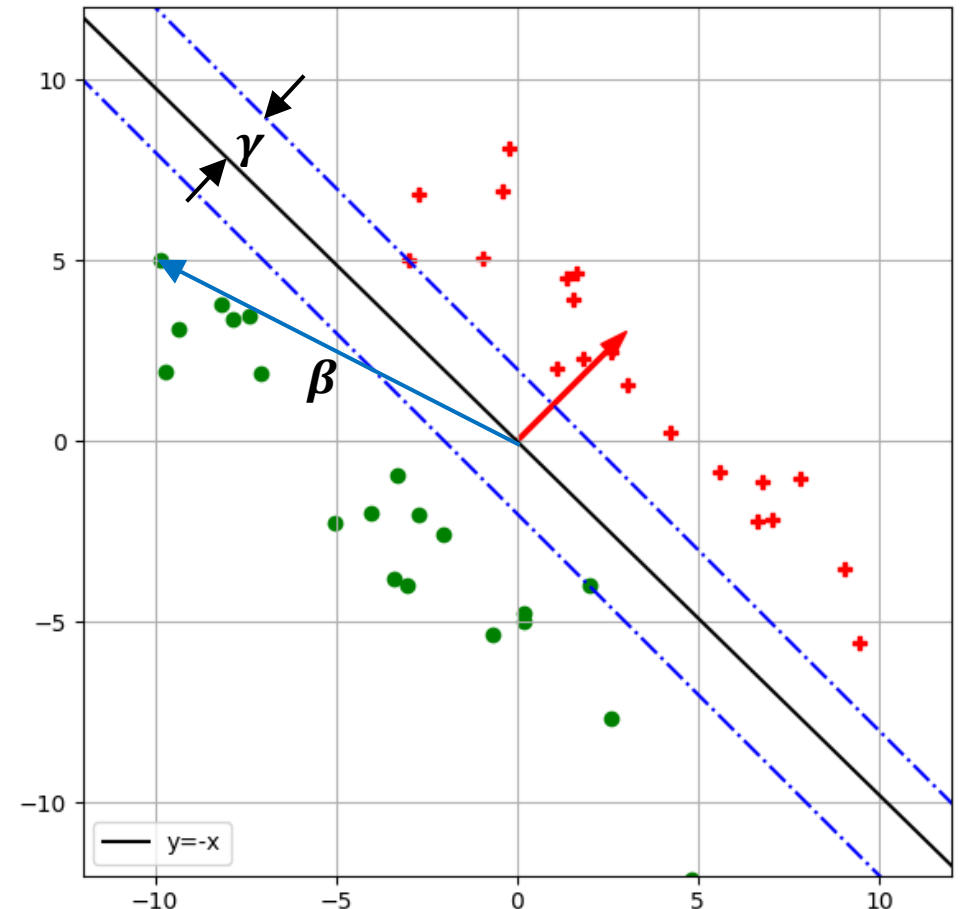
$\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} + y^{(k)} \mathbf{x}^{(k)}$ if $\mathbf{x}^{(k)}$ is misclassified.

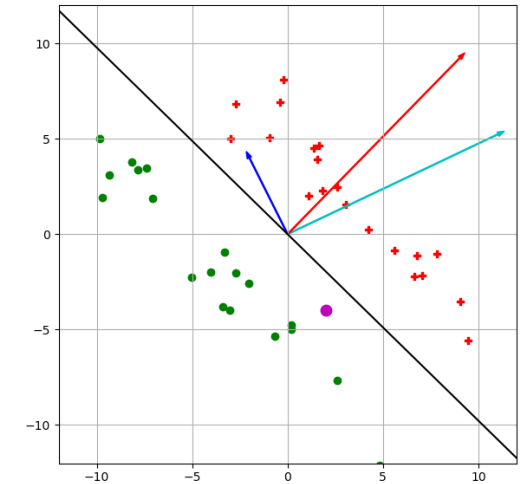
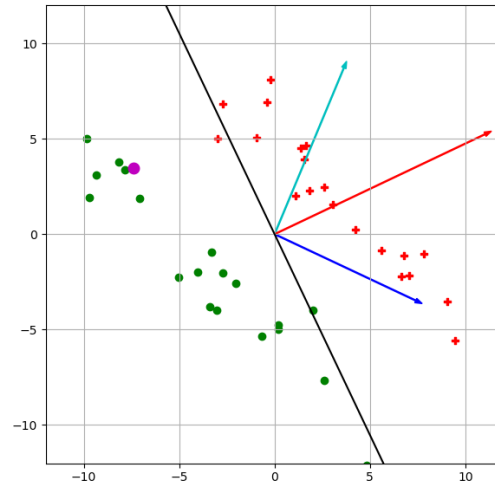
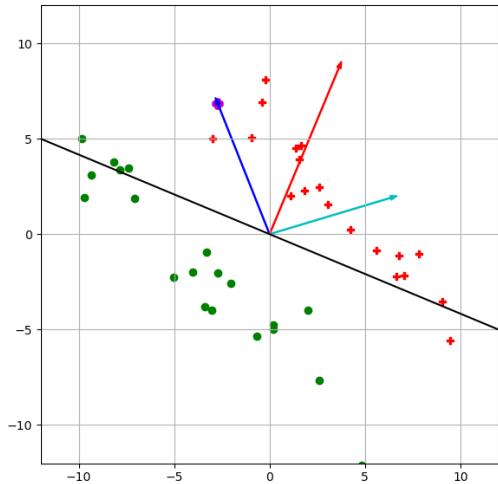
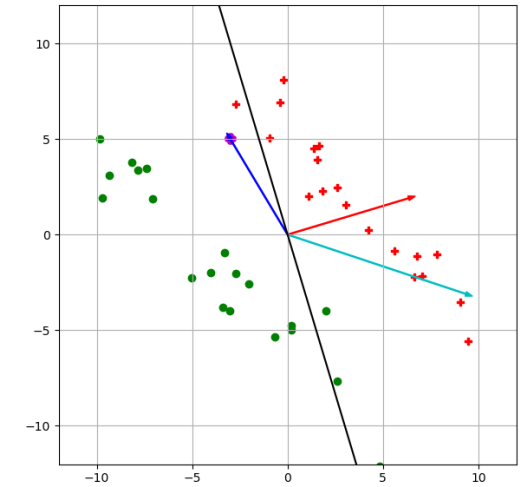
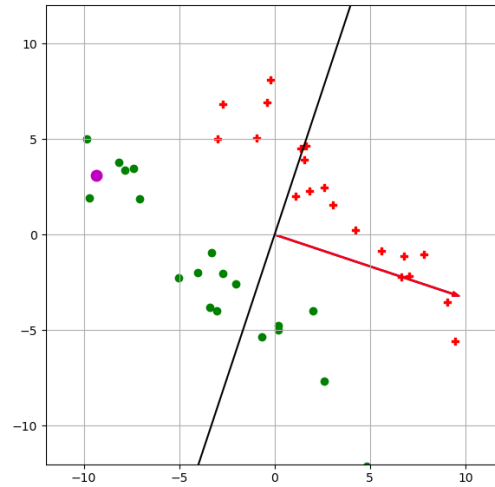
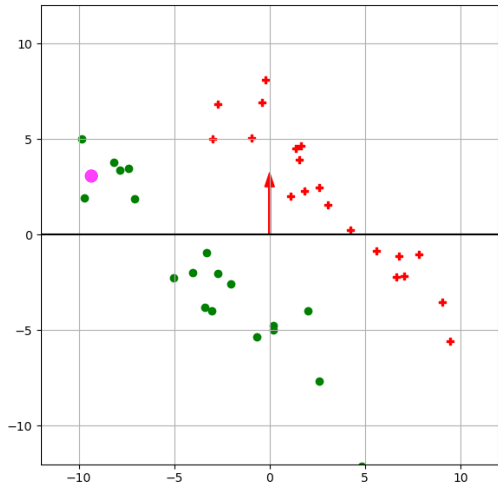
[This means $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} + \mathbf{x}^{(k)}$ or $\mathbf{w}^{(k+1)} \leftarrow \mathbf{w}^{(k)} - \mathbf{x}^{(k)}$]

$k \leftarrow k + 1$ or reset $k = 1$ if $k == N$

Convergence Theorem:-

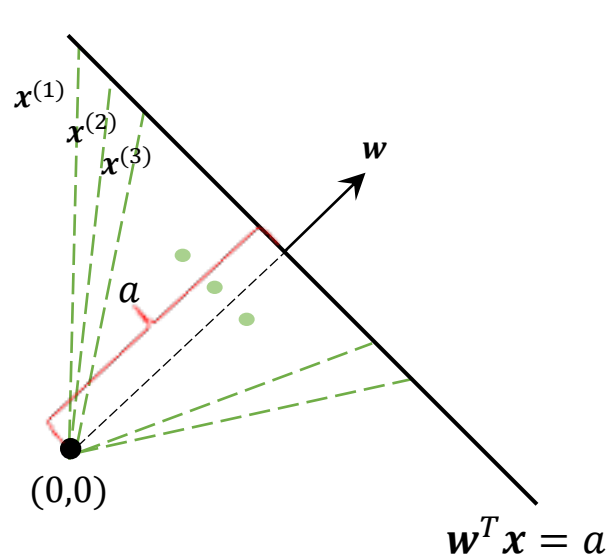
For a finite and linearly separable set of data, the perceptron learning algorithm will find a linear separator in at most $\frac{\beta^2}{\gamma^2}$ iterations where the maximum length of any data point is β and γ is the maximum margin of the linear separators.



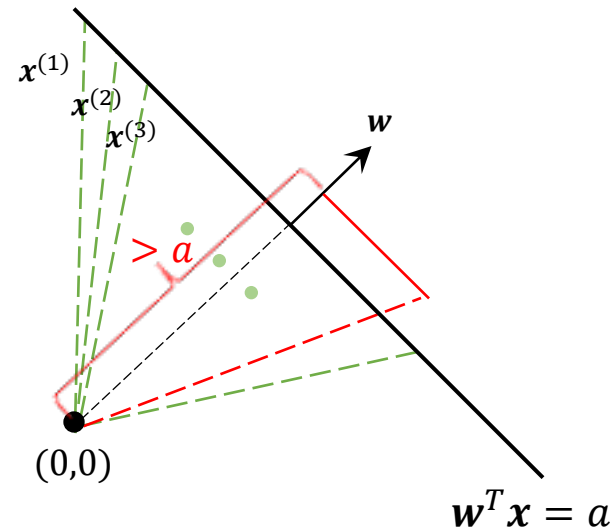


Perceptron Learning Algorithm – Convergence Proof

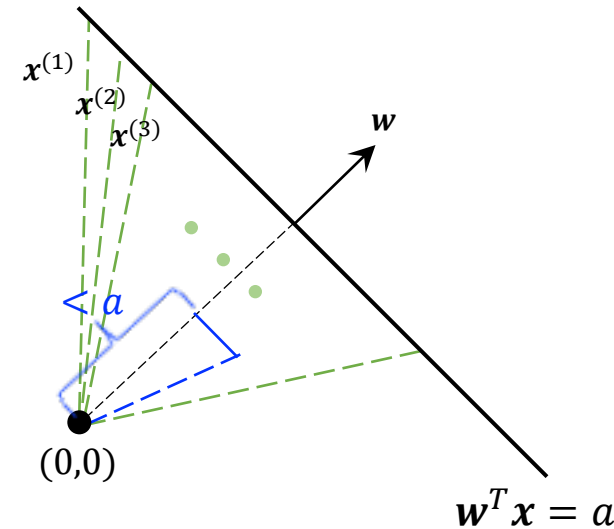
Some Important Brush-ups:-



For a hyperplane $w^T x = a$, w is the normal vector (often assumed to be unit normal vector) to the hyperplane and a is the distance from the origin.



Thus a hyperplane with normal vector w , can separate two sets of points by examining if $w^T x - a$ is > 0 or < 0 . (Remember our notation says $-a$ is nothing but b).



Perceptron Learning Algorithm – Convergence Proof

- Thus a simple classification rule can be $\mathbf{w}^T \mathbf{x} + b > 0$ or < 0 .
- For simplicity, let us assume that the hyperplane passes through origin – This does not affect the problem in our hand as it just shifts the origin. But that means $b = 0$. Thus the classification rule becomes $\text{sgn}(\mathbf{w}^T \mathbf{x}^{(i)}) = y^{(i)} \forall (\mathbf{x}^{(i)}, y^{(i)})$.
- There may not exist such hyperplane that can separate two sets of data. If it exists then the data is known to be linearly separable.
- A more formal definition is – Two sets P and N of d dimensional points are said to be linearly separable if there exists $d + 1$ real numbers $\mathbf{w} = (w_0, w_1, \dots, w_{d+1})$ such that every point $\mathbf{x} = (1, x_1, x_2, \dots, x_{d+1}) \in P$, satisfies $\text{sgn}(\mathbf{w}^T \mathbf{x}) = 1$ and every point $\mathbf{x} = (1, x_1, x_2, \dots, x_{d+1}) \in N$, satisfies $\text{sgn}(\mathbf{w}^T \mathbf{x}) = -1$,

Perceptron Learning Algorithm – Convergence Proof

Without loss of generalization we will also assume the following while proving the convergence theorem.

- We normalize all the points $\mathbf{x}^{(i)}$ so that $\max ||\mathbf{x}^{(i)}|| = 1$, i.e., maximum length of any data point is 1. Notice that this does not affect the algorithm or the solution as, $\text{sgn}(\mathbf{w}^T \mathbf{x}^{(i)}) =$

$$\text{sgn}\left(\mathbf{w}^T \frac{\mathbf{x}^{(i)}}{||\mathbf{x}^{(i)}||}\right)$$

- Margin of separation: This is the minimum distance between any data point and the hyperplane $\mathbf{w}^T \mathbf{x} = 0$. Such a hyperplane is defined to be the *maximum margin separator*.

$$\gamma = \min_x |\mathbf{w}^T \mathbf{x}|$$

Perceptron Learning Algorithm – Convergence Proof

- Let \mathbf{w}^* be the normal to the maximum margin linear separator. We want our \mathbf{w} in the perceptron algorithm in each step getting closer and closer to this (sort of) ideal \mathbf{w}^* .
- This means that the inner product between \mathbf{w} and \mathbf{w}^* should increase with each update.
- Remember the update rule: $\mathbf{w}' = \mathbf{w} + \mathbf{x}^{(i)} y^{(i)}$
- So, $(\mathbf{w}')^T \mathbf{w}^* = (\mathbf{w} + \mathbf{x}^{(i)})^T \mathbf{w}^*$ (assuming $y^{(i)} = +1$)

$$= \mathbf{w}^T \mathbf{w}^* + (\mathbf{x}^{(i)})^T \mathbf{w}^*$$

Previous $\mathbf{w}^T \mathbf{w}^*$

This is at least γ

- Similarly for a negative training example,

$$(\mathbf{w}')^T \mathbf{w}^* = (\mathbf{w} - \mathbf{x}^{(i)})^T \mathbf{w}^* = \mathbf{w}^T \mathbf{w}^* - (\mathbf{x}^{(i)})^T \mathbf{w}^*$$

$(\mathbf{x}^{(i)})^T \mathbf{w}^*$ is negative as the example is negative and thus $-(\mathbf{x}^{(i)})^T \mathbf{w}^* \geq \gamma$.

- So, in both cases $(\mathbf{w}')^T \mathbf{w}^*$ is increasing by at least γ . Which means length of \mathbf{w}' increases with each update → **Observation 1**

Perceptron Learning Algorithm – Convergence Proof

- length of \mathbf{w}' is given by $||\mathbf{w}'||^2$. Let us see, how this changes with each update.
- $||\mathbf{w}'||^2 = (\mathbf{w}')^T \mathbf{w}' = (\mathbf{w} + \mathbf{x}^{(i)} y^{(i)})^T (\mathbf{w} + \mathbf{x}^{(i)} y^{(i)}) = ||\mathbf{w}||^2 + 2y^{(i)} (\mathbf{x}^{(i)})^T \mathbf{w} + (y^{(i)})^2 ||\mathbf{x}^{(i)}||^2$
- Let us consider the case $y^{(i)} = +1$ first.

- $||\mathbf{w}'||^2 = ||\mathbf{w}||^2 + 2(\mathbf{x}^{(i)})^T \mathbf{w} + ||\mathbf{x}^{(i)}||^2$

This is negative as update occurs only on misclassification and for $y^{(i)} = +1$, misclassification means $(\mathbf{x}^{(i)})^T \mathbf{w} \leq 0$.

This is ≤ 1 , as we assumed every point is rescaled to unit ball.

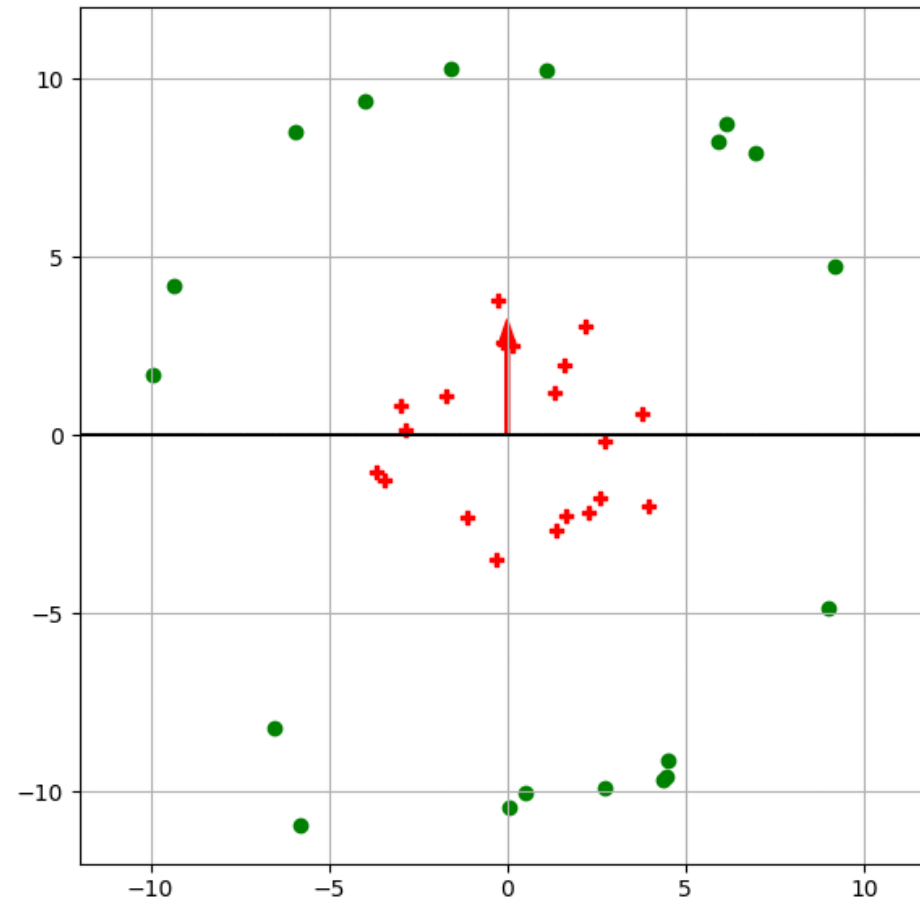
- So, $||\mathbf{w}'||^2 \leq ||\mathbf{w}||^2 + 1$

Perceptron Learning Algorithm – Convergence Proof

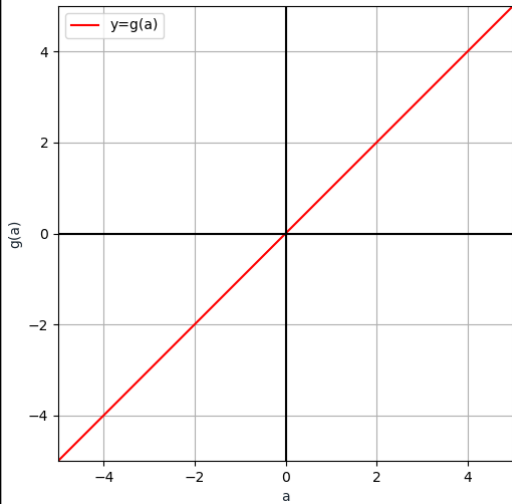
- Let us now consider the case $y^{(i)} = -1$.
- $||\mathbf{w}'||^2 = ||\mathbf{w}||^2 - 2(\mathbf{x}^{(i)})^T \mathbf{w} + ||\mathbf{x}^{(i)}||^2$
- Using similar arguments, $(\mathbf{x}^{(i)})^T \mathbf{w} \geq 0$, so, $-2(\mathbf{x}^{(i)})^T \mathbf{w} \leq 0$ and also $||\mathbf{x}^{(i)}||^2 \leq 1$
- So, here also $||\mathbf{w}'||^2 \leq ||\mathbf{w}||^2 + 1$
- So, the length of the updated vector is increasing but it is increasing in a controlled way. →
Observation 2

Perceptron Learning Algorithm – Convergence Proof

- Say after k iterations, \mathbf{w} is denoted as \mathbf{w}_k .
- So **observation 1** implies, $\mathbf{w}_k^T \mathbf{w}^* \geq k\gamma$ and **observation 2** implies $||\mathbf{w}_k||^2 \leq k$ [remember the initial value of \mathbf{w} is $\mathbf{0}$]
- Now we have to apply the Cauchy-Swartz inequality.
- $||\mathbf{w}_k|| ||\mathbf{w}^*|| \geq \mathbf{w}_k^T \mathbf{w}^*$, but $||\mathbf{w}^*|| = 1$, so,
- $||\mathbf{w}_k|| \geq \mathbf{w}_k^T \mathbf{w}^* \geq k\gamma$ [**observation 1**]
- This implies $k^2\gamma^2 \leq ||\mathbf{w}_k||^2 \leq k$ [**observation 2**]
- So, $k \leq \frac{1}{\gamma^2}$

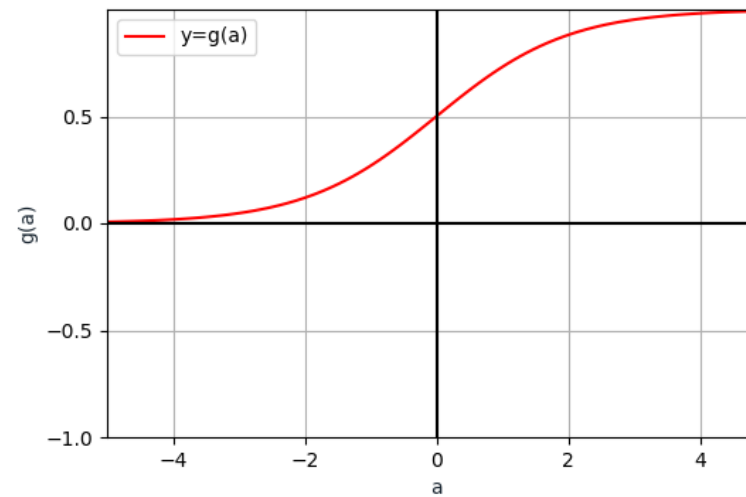


Popular Choices of Activation Functions



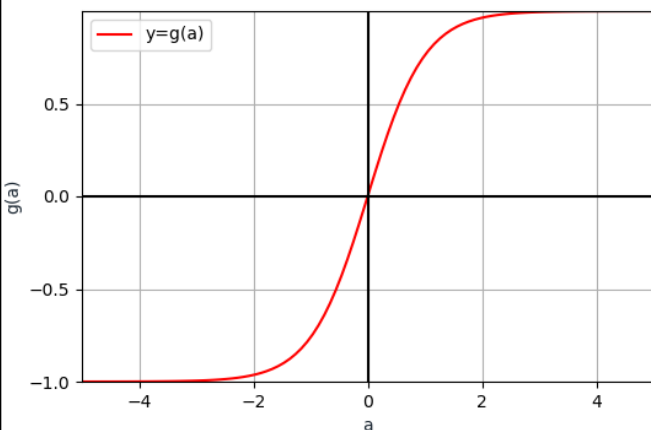
Linear activation function

- $g(a) = a$
- Unbounded
- $g'(a) = 1$



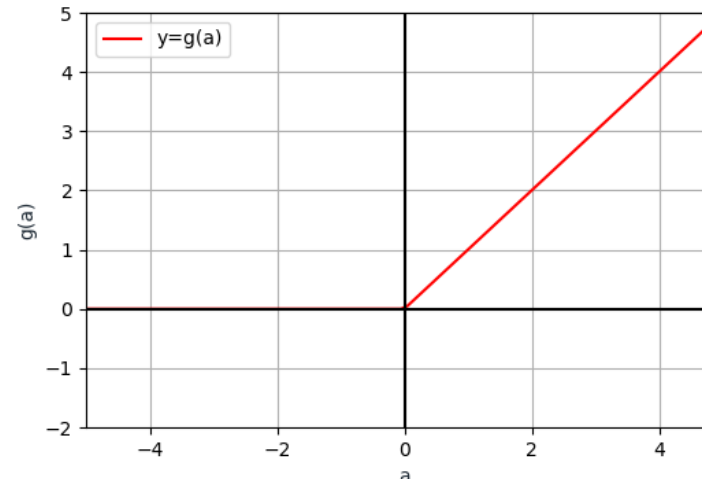
Sigmoid activation function

- $g(a) = \sigma(a) = \frac{1}{1 + \exp(-a)}$
- Bounded (0, 1)
- Always positive
- $g'(a) = g(a)(1 - g(a))$



tanh activation function

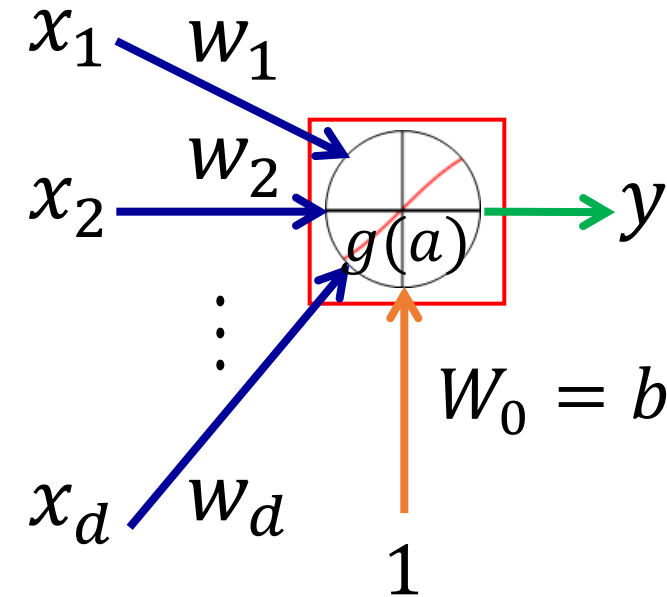
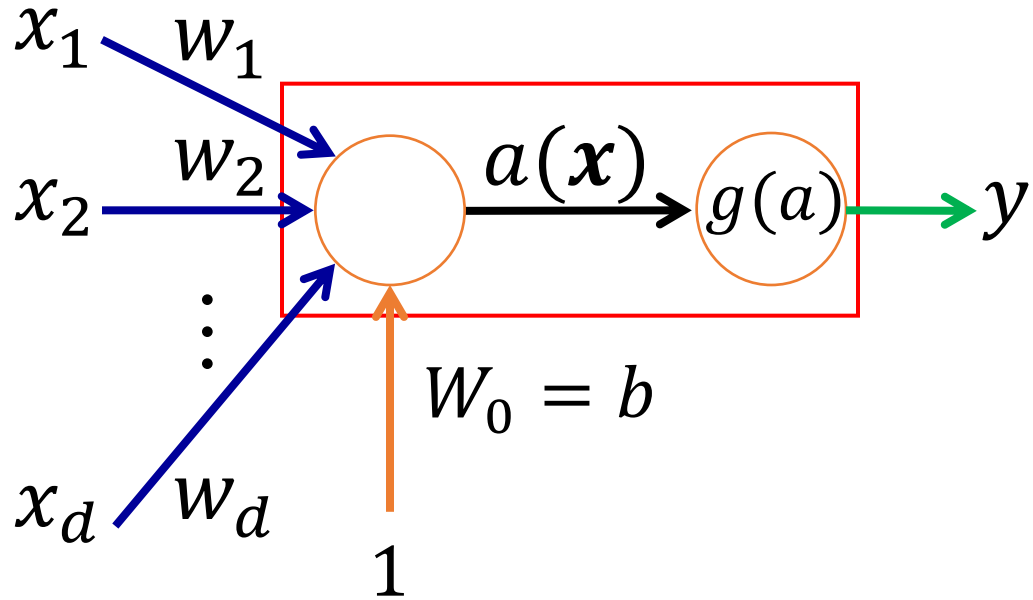
- $g(a) = \tanh(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)}$
- Bounded (-1, 1)
- Can be positive or negative
- $g'(a) = 1 - g^2(a)$



ReLU activation function

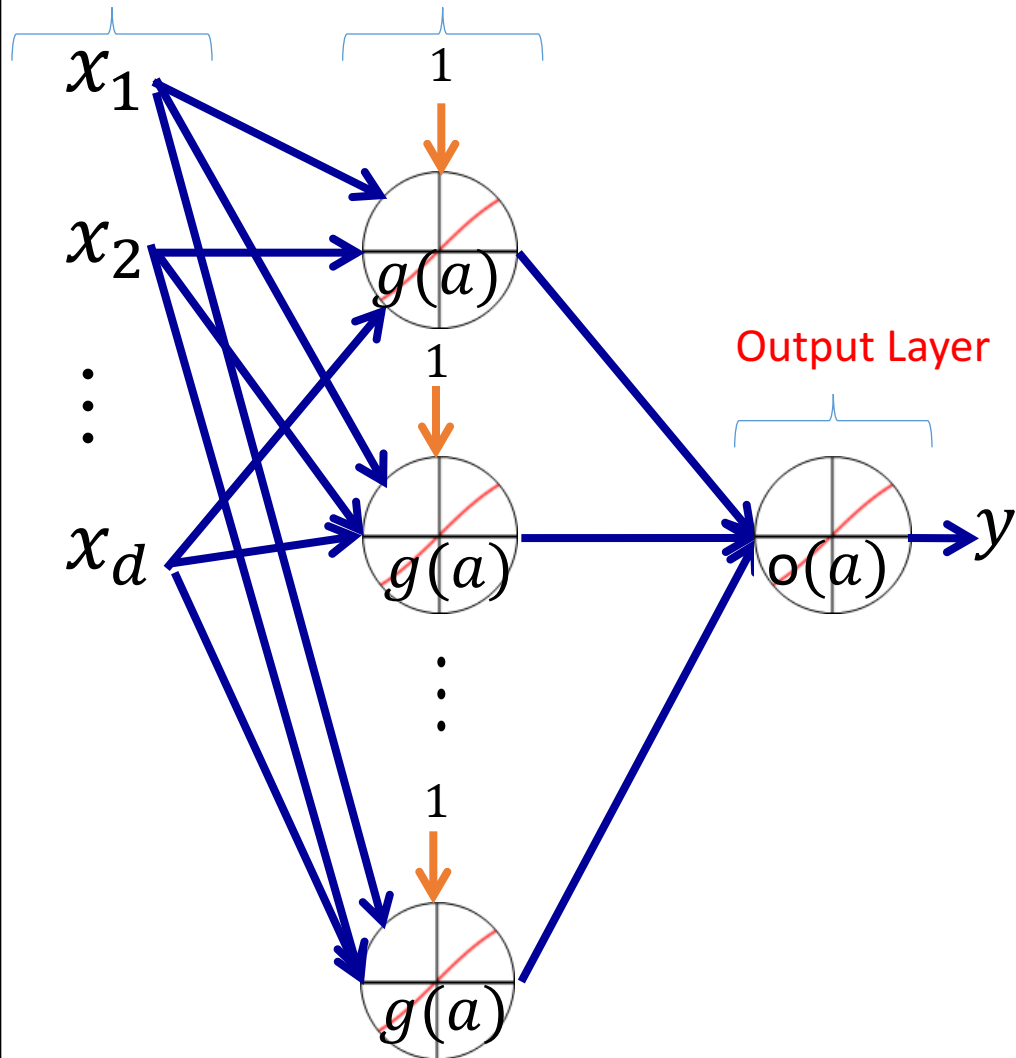
- $g(a) = \max(0, a)$
- Bounded below by 0
- But not upper-bounded
- $g'(a) = \begin{cases} 1, & a \geq 0 \\ 0, & a < 0 \end{cases}$

Multilayer Neural Network



A bit of notation change

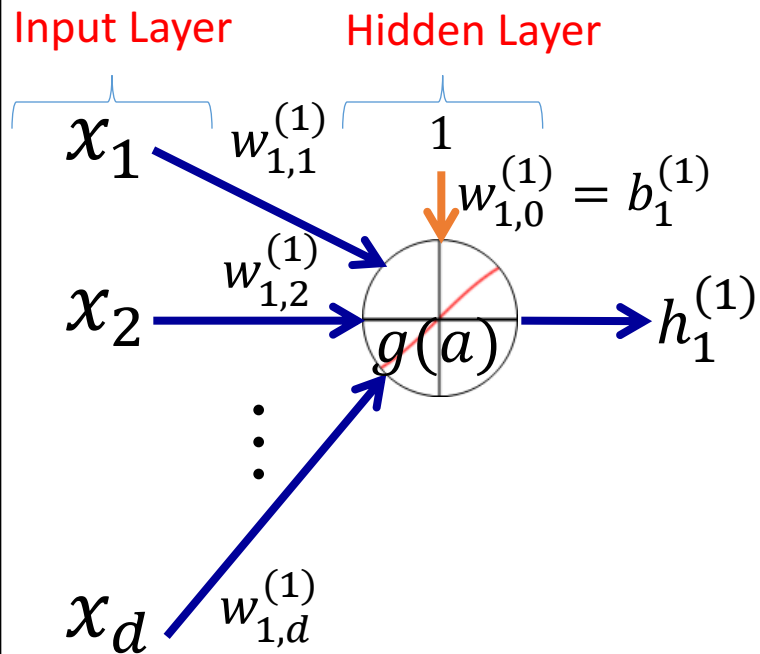
Multilayer Neural Network



Intuition:-

Hidden Layer: Extracts better representation of the input data

Output layer: Does the classification



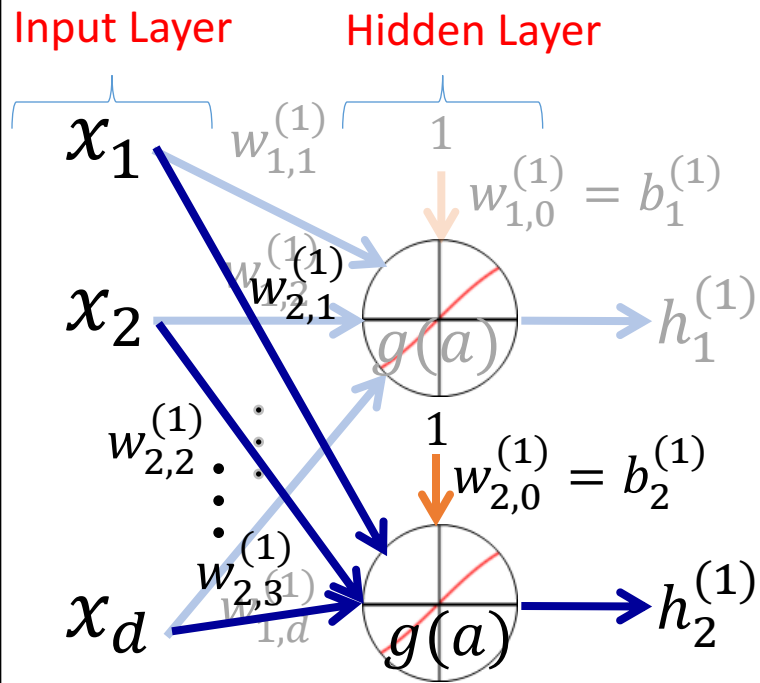
Multilayer Neural Network

$$a_1^{(1)}(\mathbf{x}) = b_1^{(1)} + \sum_{i=1}^d w_{1,i}^{(1)} x_i = \begin{bmatrix} \mathbf{w}_1^{(1)} & b_1^{(1)} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

Row vector $\begin{bmatrix} w_{1,1}^{(1)}, w_{1,2}^{(1)}, \dots, w_{1,d}^{(1)} \end{bmatrix}$

Column vector $\begin{bmatrix} x_1 & x_2 & \dots & x_d \end{bmatrix}^T$

$$h_1^{(1)} = g \left(a_1^{(1)}(\mathbf{x}) \right)$$



Multilayer Neural Network

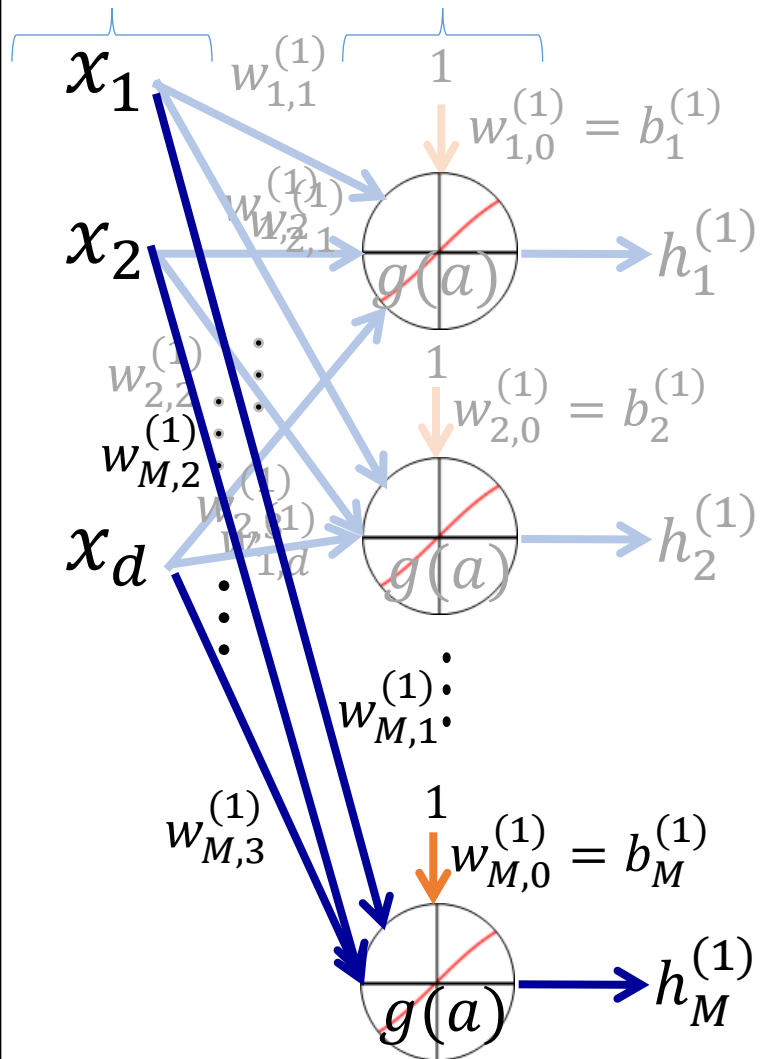
$$a_1^{(1)}(\mathbf{x}) = b_1^{(1)} + \sum_{i=1}^d w_{1,i}^{(1)} x_i = [\mathbf{w}_1^{(1)} b_1^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$h_1^{(1)} = g\left(a_1^{(1)}(\mathbf{x})\right)$$

$$a_2^{(1)}(\mathbf{x}) = b_2^{(1)} + \sum_{i=1}^d w_{2,i}^{(1)} x_i = [\mathbf{w}_2^{(1)} b_2^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$h_2^{(1)} = g\left(a_2^{(1)}(\mathbf{x})\right)$$

Input Layer Hidden Layer



Multilayer Neural Network

$$a_1^{(1)}(\mathbf{x}) = b_1^{(1)} + \sum_{i=1}^d w_{1,i}^{(1)} x_i = [\mathbf{w}_1^{(1)} b_1^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

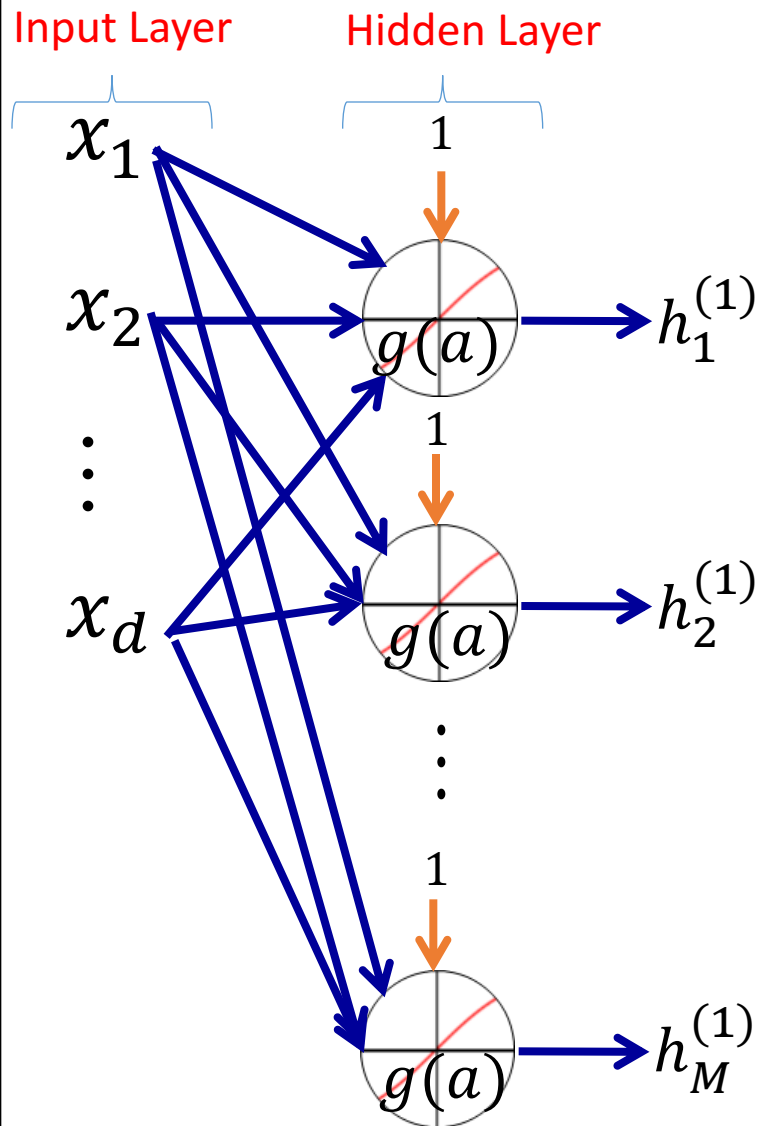
$$h_1^{(1)} = g\left(a_1^{(1)}(\mathbf{x})\right)$$

$$a_2^{(1)}(\mathbf{x}) = b_2^{(1)} + \sum_{i=1}^d w_{2,i}^{(1)} x_i = [\mathbf{w}_2^{(1)} b_2^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$h_2^{(1)} = g\left(a_2^{(1)}(\mathbf{x})\right)$$

$$a_M^{(1)}(\mathbf{x}) = b_M^{(1)} + \sum_{i=1}^d w_{M,i}^{(1)} x_i = [\mathbf{w}_M^{(1)} b_M^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

$$h_M^{(1)} = g\left(a_M^{(1)}(\mathbf{x})\right)$$



Multilayer Neural Network

$$\begin{aligned} a_1^{(1)}(x) &= [\mathbf{w}_1^{(1)} \ b_1^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \\ a_2^{(1)}(x) &= [\mathbf{w}_2^{(1)} \ b_2^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \\ &\vdots \\ a_M^{(1)}(x) &= [\mathbf{w}_M^{(1)} \ b_M^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix} \end{aligned}$$

$$\begin{bmatrix} a_1^{(1)}(x) \\ a_2^{(1)}(x) \\ \vdots \\ a_M^{(1)}(x) \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^{(1)} & b_1^{(1)} \\ \mathbf{w}_2^{(1)} & b_2^{(1)} \\ \vdots & \vdots \\ \mathbf{w}_M^{(1)} & b_M^{(1)} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

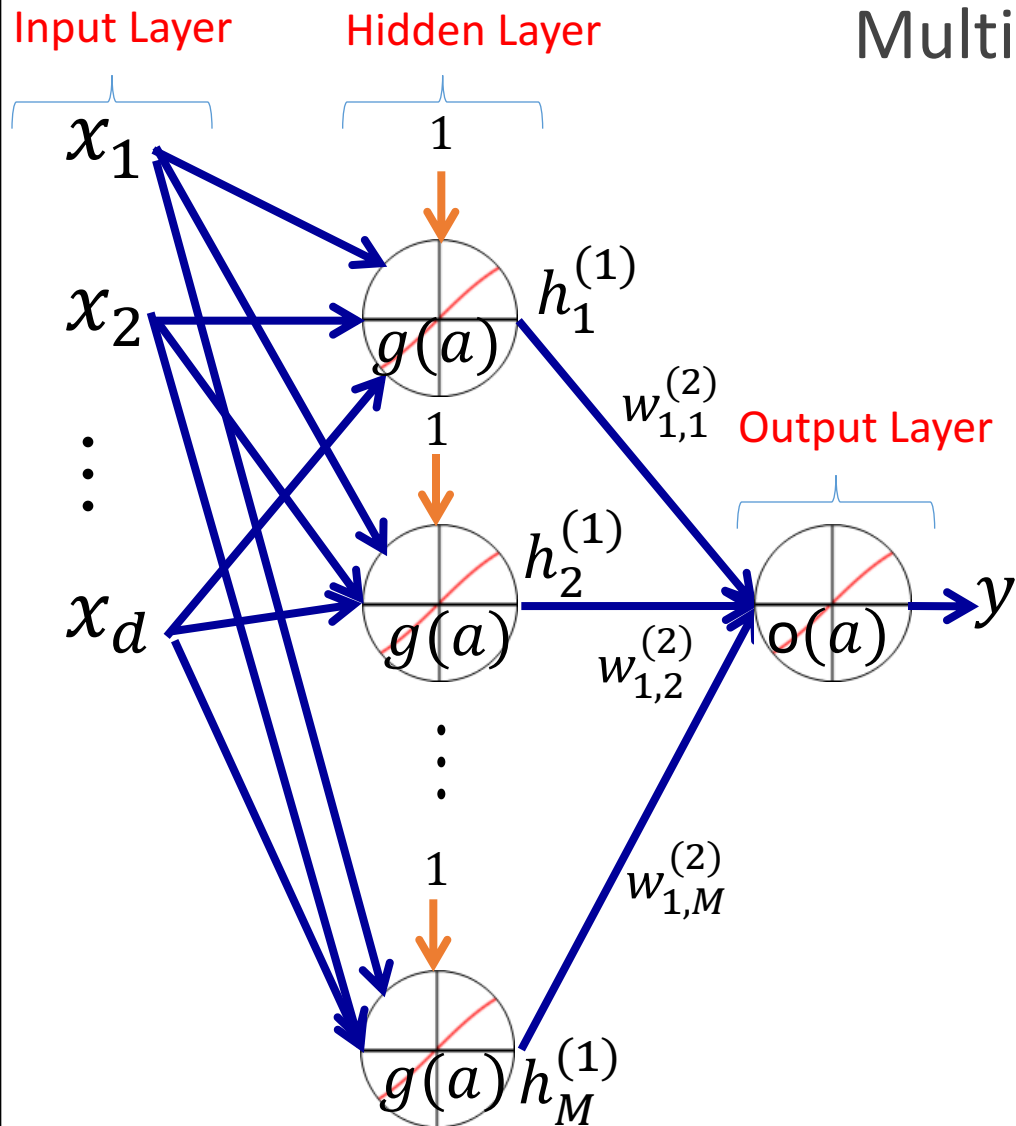
$$\mathbf{a}^{(1)}(x) = [\mathbf{W}^{(1)} \ \mathbf{b}^{(1)}] \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

Column vector
Matrix
Column vector

Column vector

Hidden Layer Activation

$$\begin{bmatrix} h_1^{(1)} \\ h_2^{(1)} \\ \vdots \\ h_M^{(1)} \end{bmatrix} = \begin{bmatrix} g(a_1^{(1)}) \\ g(a_2^{(1)}) \\ \vdots \\ g(a_M^{(1)}) \end{bmatrix} \quad \mathbf{h}^{(1)} = \mathbf{g}(\mathbf{a}^{(1)})$$



Output Layer Pre-activation

$$a_1^{(2)}(\mathbf{h}^{(1)}) = \left[\mathbf{w}_1^{(2)} b_1^{(2)} \right] \begin{bmatrix} \mathbf{h}^{(1)} \\ 1 \end{bmatrix}$$

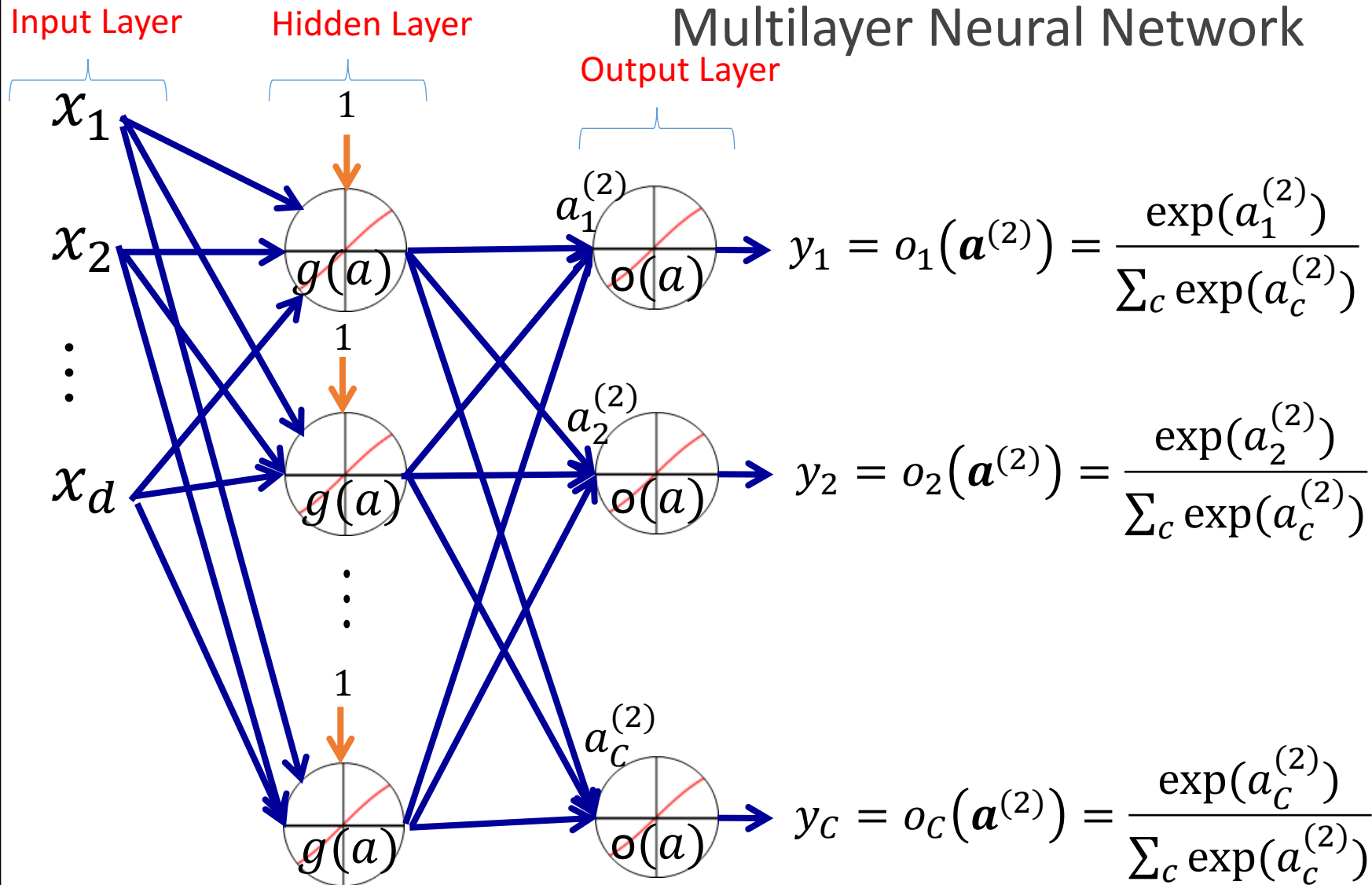
Output Layer Activation

$$y = o \left(a_1^{(2)}(\mathbf{h}^{(1)}) \right)$$

o can have many options

- sigmoid
- tanh
- Relu etc.

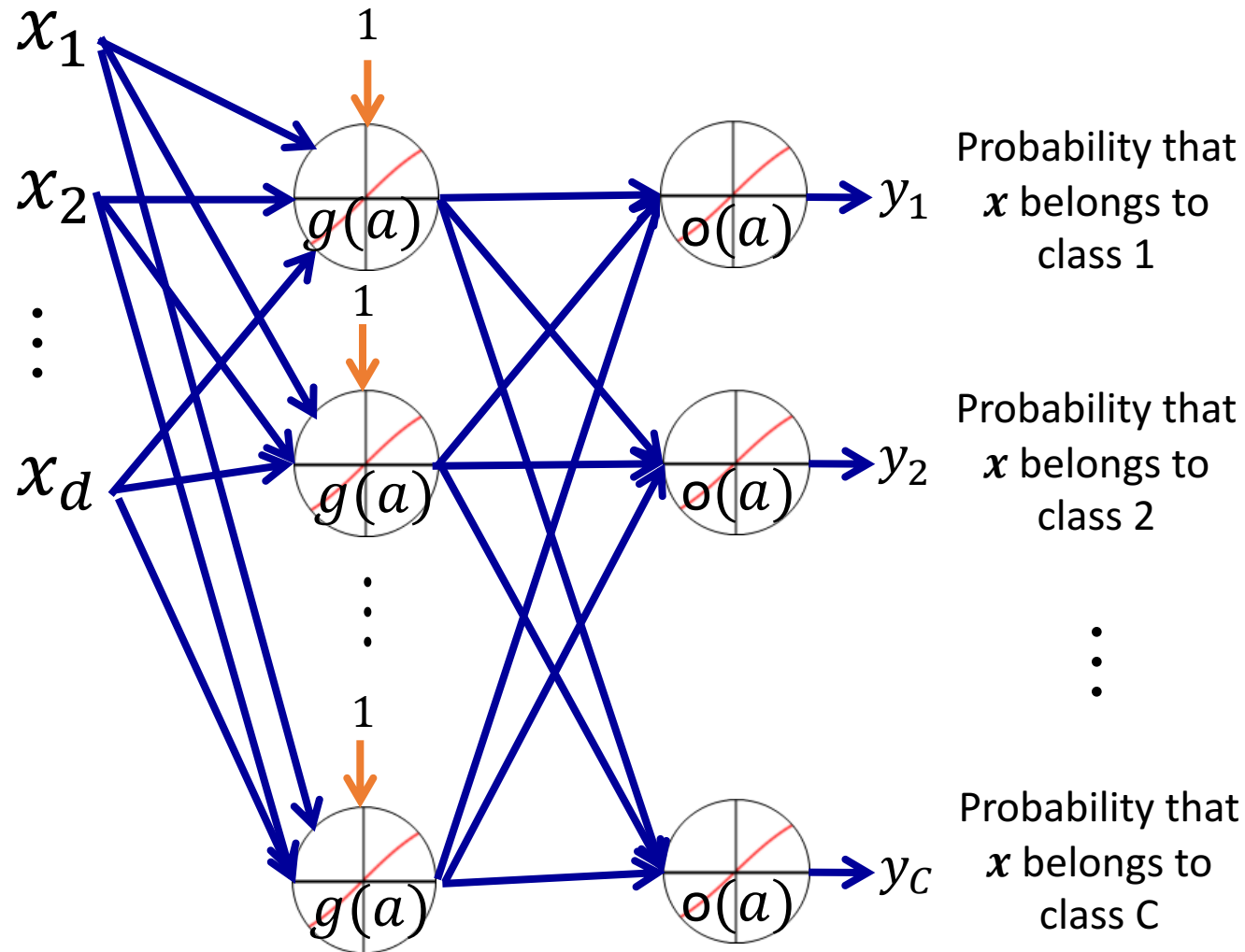
One special is softmax function



Softmax function:

Exponentiate each component of output activation vector and elementwise divide by the sum of the exponentials.

Training a Neural Network – Loss Function

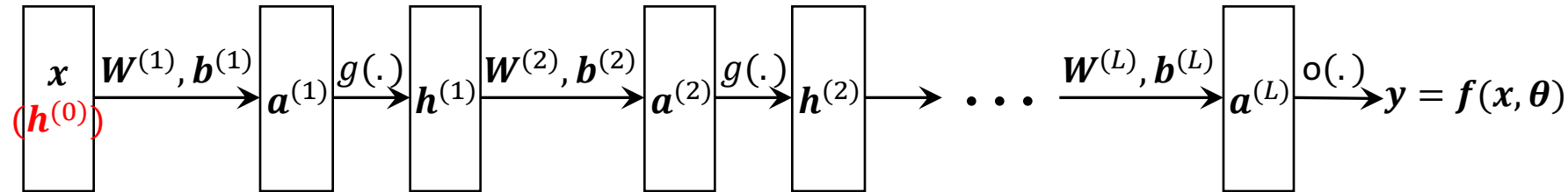


So we can aim to maximize the probability corresponding to the correct class for any example x

$$\begin{aligned} \max y_c \\ \equiv \max (\log y_c) \\ \equiv \min (-\log y_c) \end{aligned}$$

Can be equivalently expressed as
 $-\sum_i \mathbb{I}(i = c) \log(y_i)$ known as cross-entropy loss

Forward Pass in a Nutshell



θ is the collection of all learnable parameters
i.e., all W and b

Hidden layer pre-activation:

For $l = 1, \dots, L$; $\mathbf{a}^{(l)} = \mathbf{W}^{(l)} \mathbf{h}^{(l-1)} + \mathbf{b}^{(l)}$

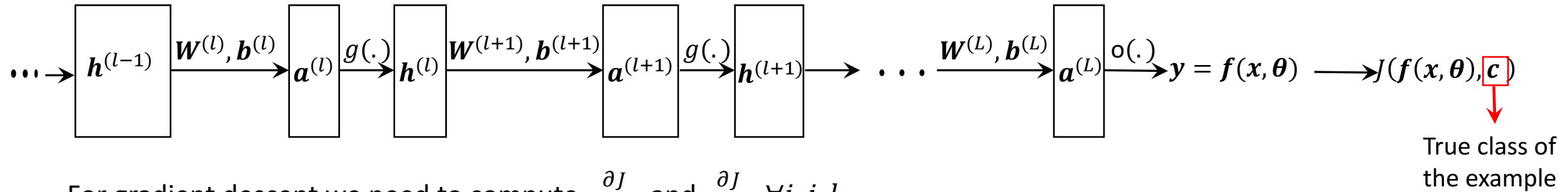
Hidden layer activation:

For $l = 1, \dots, L - 1$; $\mathbf{h}^{(l)} = g(\mathbf{a}^{(l)})$

Output layer activation:

For $l = L$; $\mathbf{h}^{(L)} = o(\mathbf{a}^{(L)}) = \mathbf{f}(\mathbf{x}, \theta)$

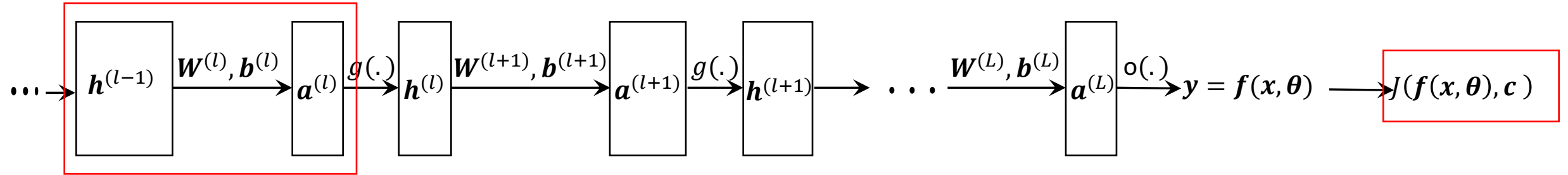
Backpropagation



For gradient descent we need to compute $\frac{\partial J}{\partial w_{i,j}^{(l)}}$ and $\frac{\partial J}{\partial b_i^{(l)}} \forall i, j, l$

Backpropagation consists of intelligent use of the chain rule of differentiation for computation of these gradients

Backpropagation



$W_{i,j}^{(l)}$ influences J through $a_i^{(l)}$. That means J is a function of $a_i^{(l)}$ and $a_i^{(l)}$, in turn, is a function of $W_{i,j}^{(l)}$

$$\frac{\partial J}{\partial W_{i,j}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \frac{\partial a_i^{(l)}}{\partial W_{i,j}^{(l)}}$$

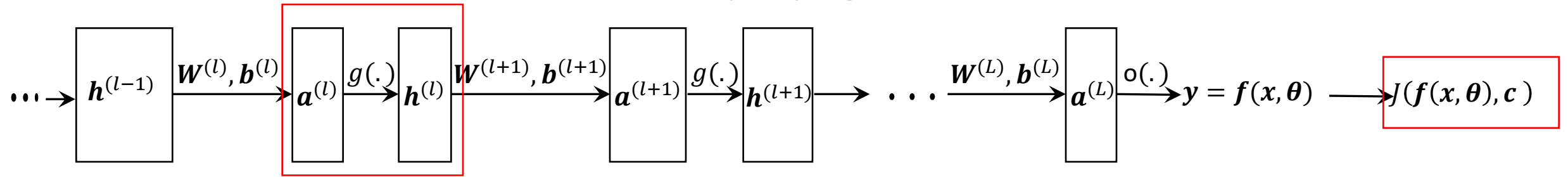
$$a_i^{(l)} = \sum_j W_{i,j}^{(l)} h_j^{(l-1)} + b_i^{(l)} = W_{i,1}^{(l)} h_1^{(l-1)} + W_{i,2}^{(l)} h_2^{(l-1)} + \dots + W_{i,j}^{(l)} h_j^{(l-1)} + \dots + b_i^{(l)}$$

$$\Rightarrow \frac{\partial J}{\partial W_{i,j}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \frac{\partial a_i^{(l)}}{\partial W_{i,j}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} h_j^{(l-1)} \quad \dots (1)$$

$$\text{And similarly } \frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \frac{\partial a_i^{(l)}}{\partial b_i^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \quad \dots (2)$$

We know $h_j^{(l-1)}$ from forward pass,
so we need $\frac{\partial J}{\partial a_i^{(l)}}$

Backpropagation

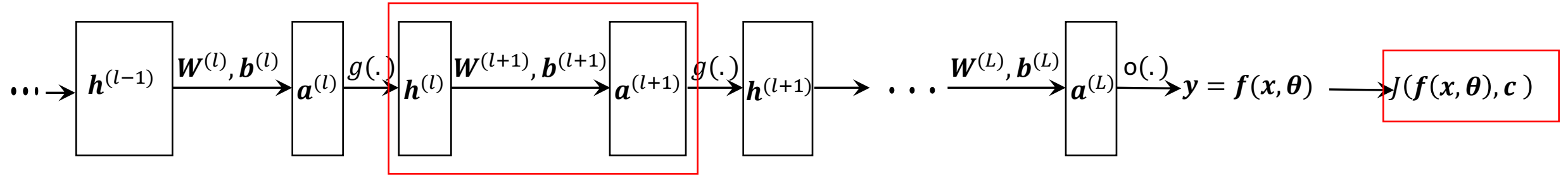


$a_i^{(l)}$ influences J through $h_i^{(l)}$. That means J is a function of $h_i^{(l)}$ and $h_i^{(l)}$, in turn, is a function of $a_i^{(l)}$

$$\frac{\partial J}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} \frac{\partial h_i^{(l)}}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} g' \left(a_i^{(l)} \right) \quad \dots (3)$$

We know $g' \left(a_i^{(l)} \right)$ from forward pass, so we need $\frac{\partial J}{\partial h_i^{(l)}}$

Backpropagation



$h_i^{(l)}$ influences J through $a_{1:M_{l+1}}^{(l+1)}$. That means J is a function of $a_{1:M_{l+1}}^{(l+1)}$ and each of $a_{1:M_{l+1}}^{(l+1)}$, in turn, is a function of $h_i^{(l)}$

Number of neurons in hidden layer $(l + 1)$

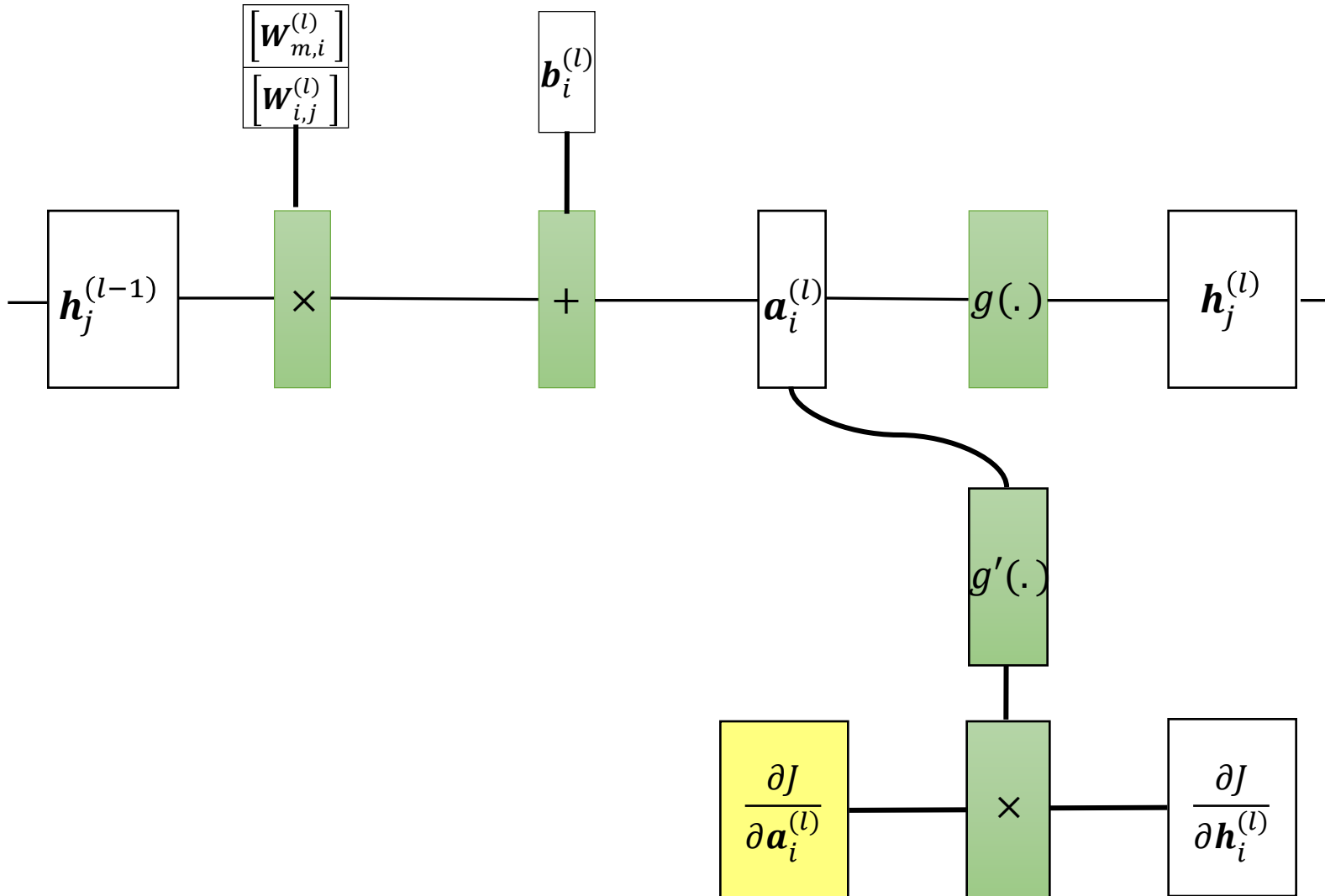
$$\frac{\partial J}{\partial h_i^{(l)}} = \sum_{m=1}^{M_{l+1}} \frac{\partial J}{\partial a_m^{(l+1)}} \frac{\partial a_m^{(l+1)}}{\partial h_i^{(l)}}$$

Remember If, $u = f(x, y)$, where $x = \phi(t)$, $y = \psi(t)$, then $\frac{\partial u}{\partial t} = \frac{\partial u}{\partial x} \frac{\partial x}{\partial t} + \frac{\partial u}{\partial y} \frac{\partial y}{\partial t}$

$$a_m^{(l+1)} = \sum_j w_{m,j}^{(l+1)} h_j^{(l)} + b_m^{(l+1)} = w_{m,1}^{(l+1)} h_1^{(l)} + w_{m,2}^{(l+1)} h_2^{(l)} + \dots + w_{m,i}^{(l+1)} h_i^{(l)} + \dots + b_m^{(l+1)}$$

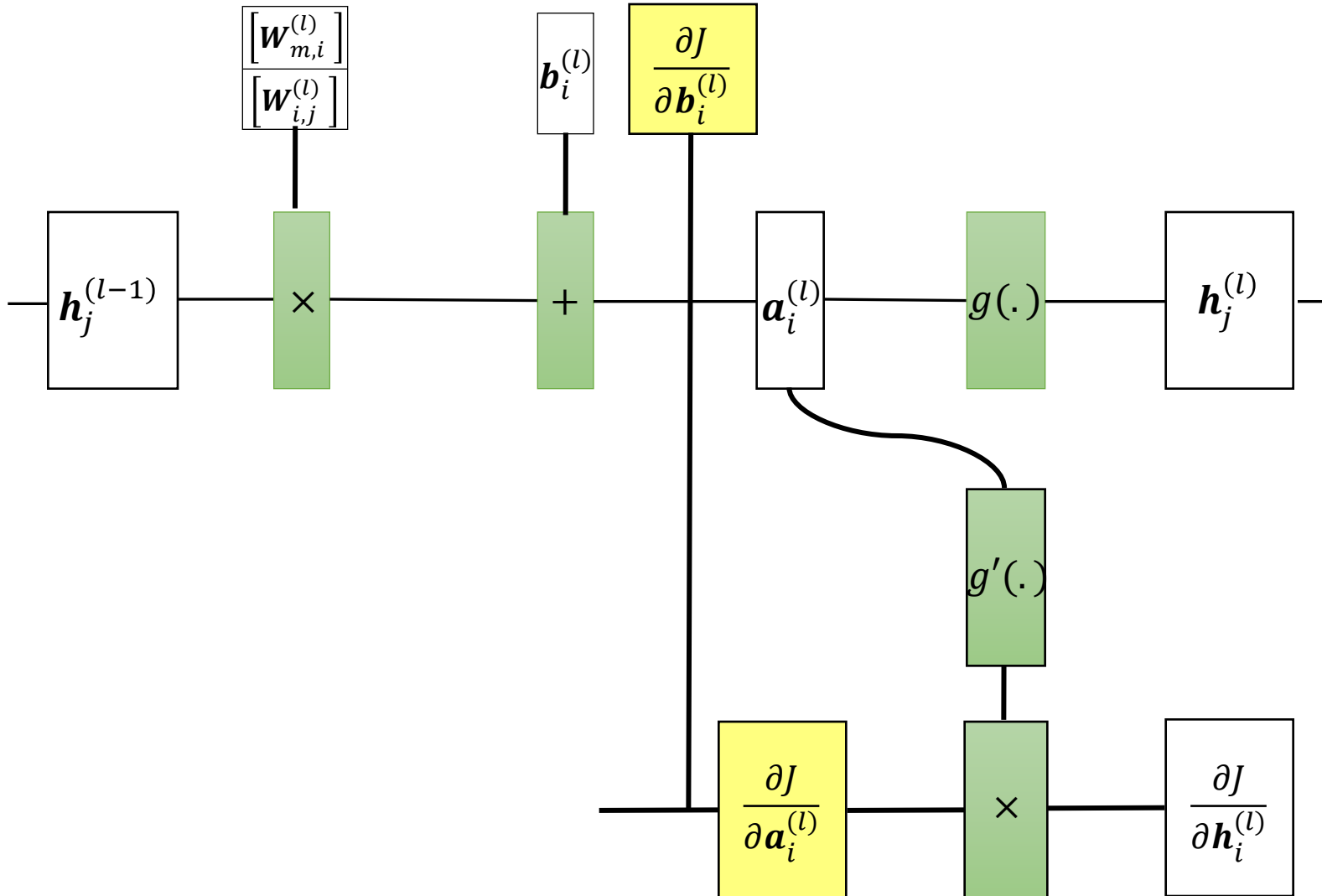
$$\Rightarrow \frac{\partial J}{\partial h_i^{(l)}} = \sum_{m=1}^{M_{l+1}} \frac{\partial J}{\partial a_m^{(l+1)}} w_{m,i}^{(l+1)} \quad \dots (4)$$

Backpropagation



$$\frac{\partial J}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} g' \left(a_i^{(l)} \right) \quad \dots (3)$$

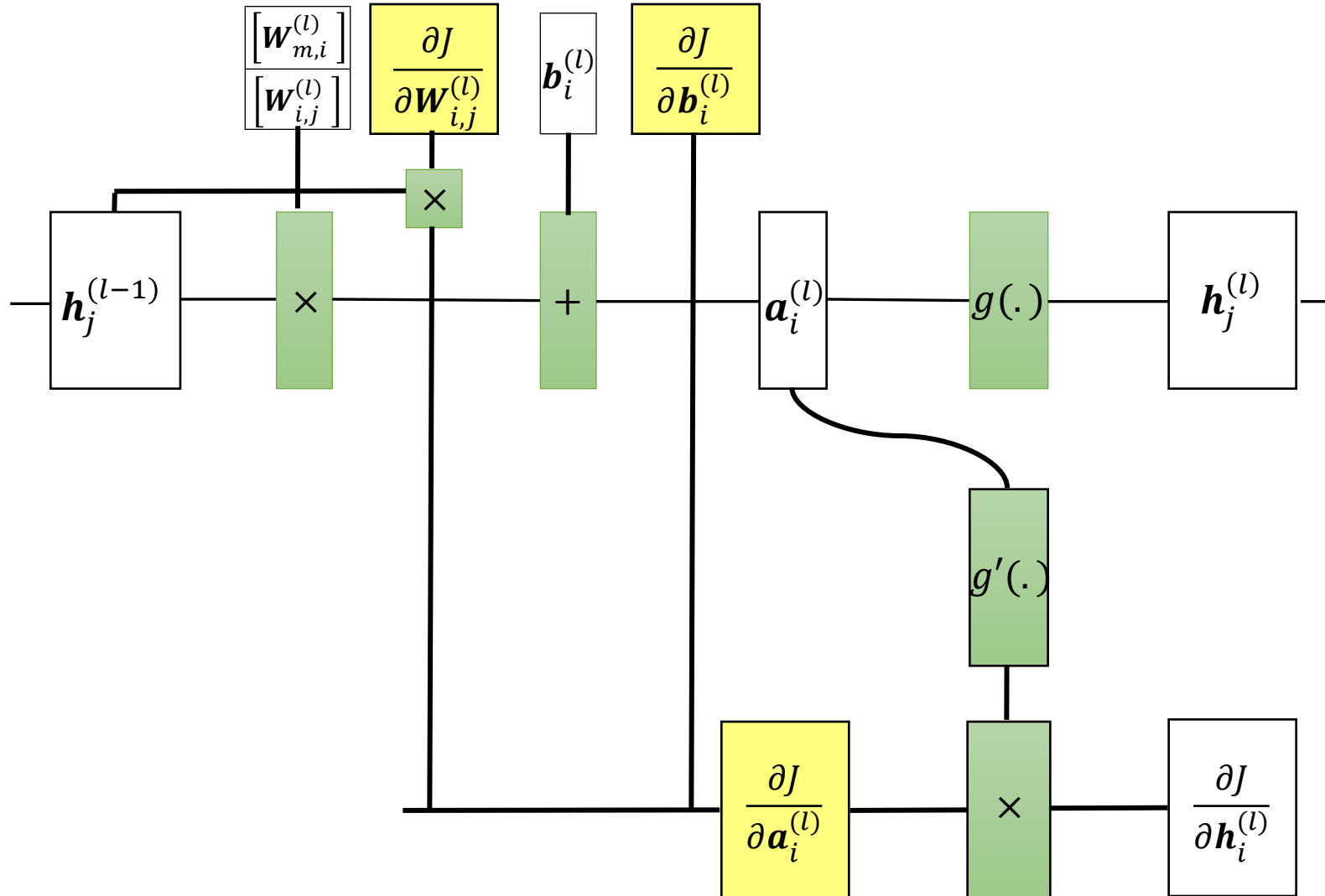
Backpropagation



$$\frac{\partial J}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} g' \left(a_i^{(l)} \right) \quad \dots (3)$$

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \quad \dots (2)$$

Backpropagation

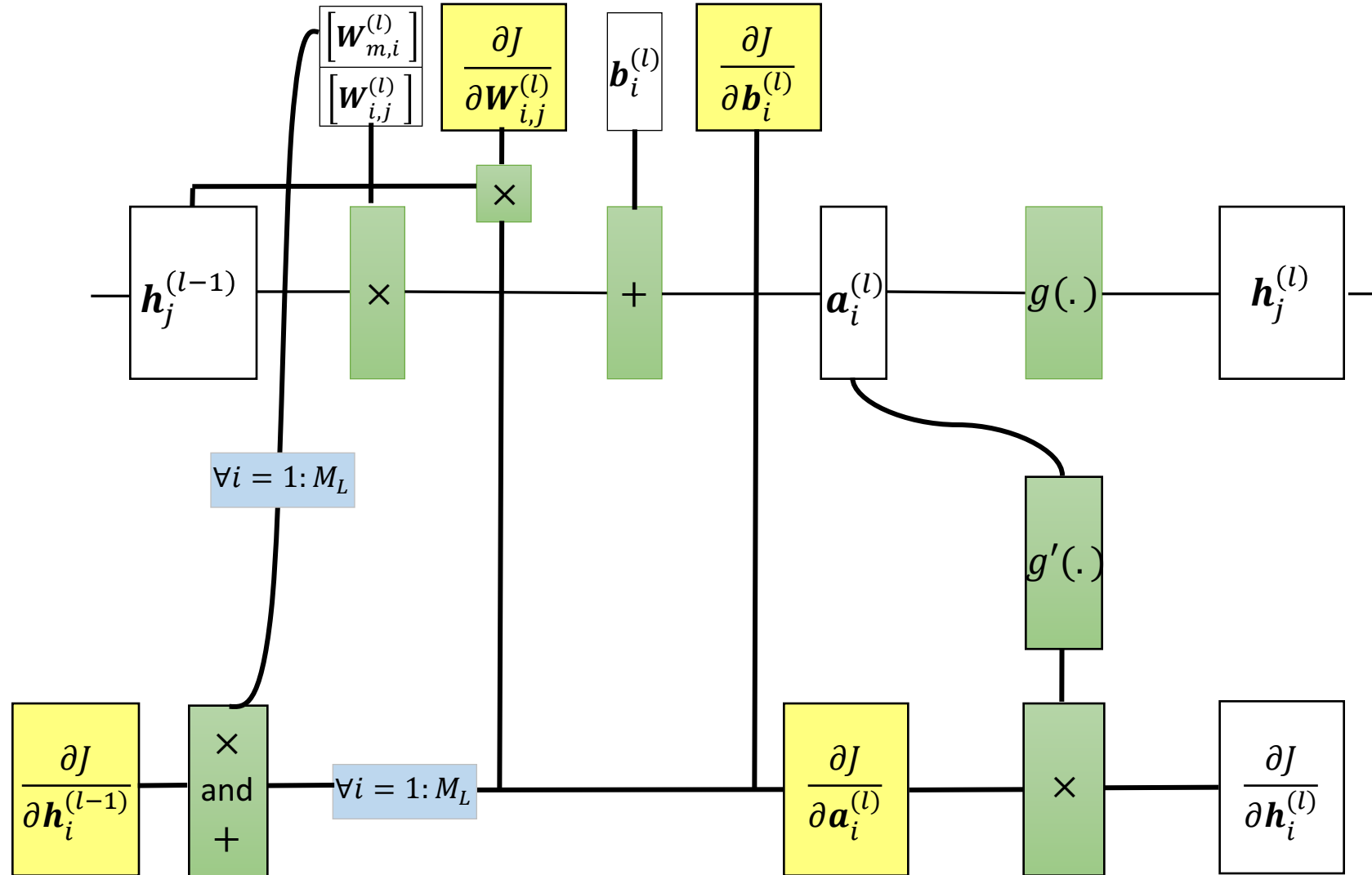


$$\frac{\partial J}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} g' \left(a_i^{(l)} \right) \quad \dots (3)$$

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \quad \dots (2)$$

$$\frac{\partial J}{\partial w_{i,j}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} h_j^{(l-1)} \quad \dots (1)$$

Backpropagation



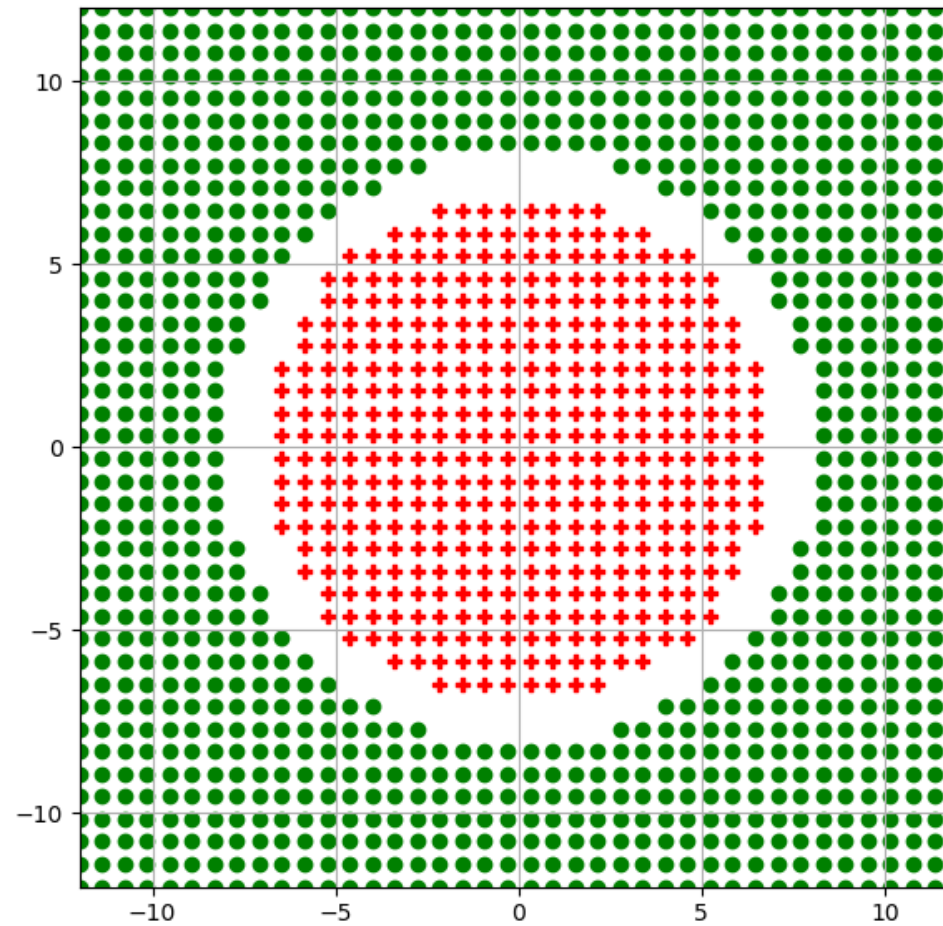
$$\frac{\partial J}{\partial a_i^{(l)}} = \frac{\partial J}{\partial h_i^{(l)}} g' (a_i^{(l)}) \quad \dots (3)$$

$$\frac{\partial J}{\partial b_i^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} \quad \dots (2)$$

$$\frac{\partial J}{\partial w_{i,j}^{(l)}} = \frac{\partial J}{\partial a_i^{(l)}} h_j^{(l-1)} \quad \dots (1)$$

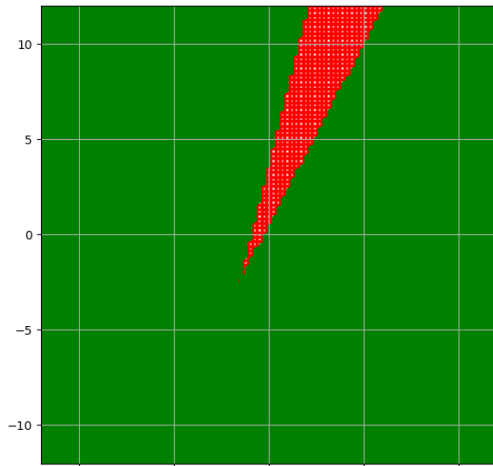
$$\frac{\partial J}{\partial h_i^{(l-1)}} = \sum_{m=1}^{M_l} \frac{\partial J}{\partial a_m^{(l)}} w_{m,i}^{(l)} \quad \dots (4)$$

Training Multilayer Neural Network for non-linearly Separable Data

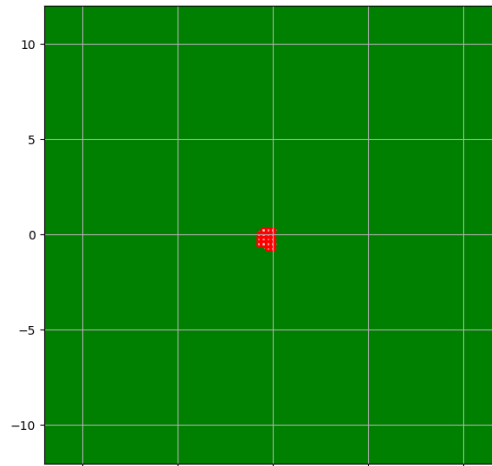


Training Data

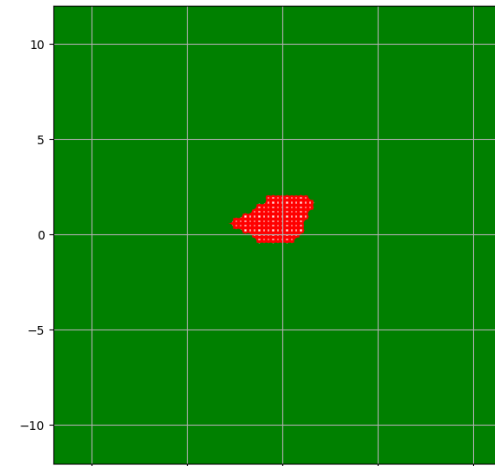
Learned Decision Boundary with Single Hidden Layer



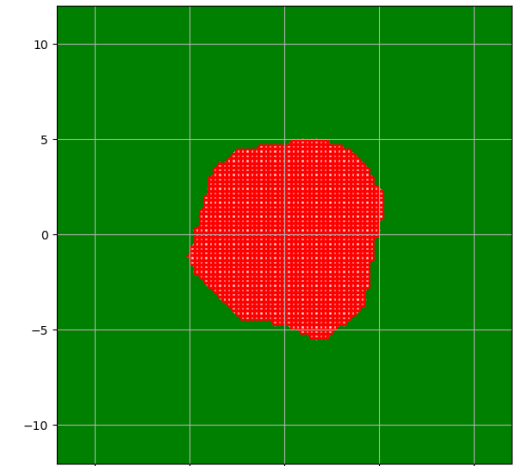
of hidden neurons = 2



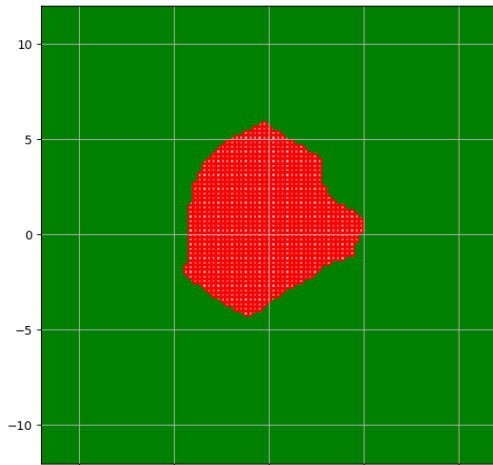
of hidden neurons = 4



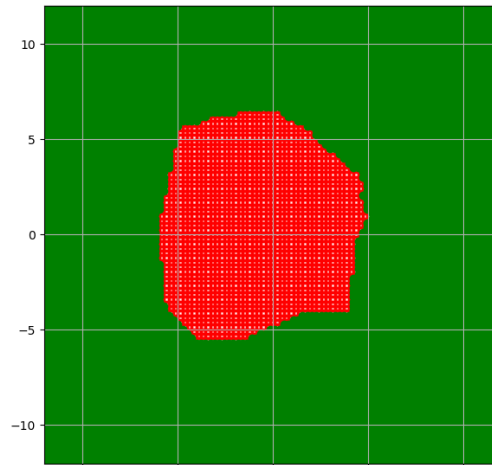
of hidden neurons = 8



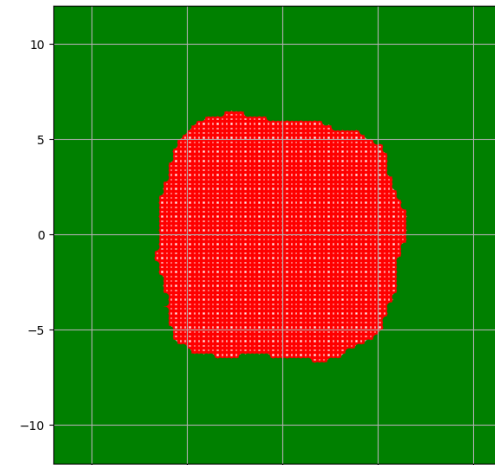
of hidden neurons = 16



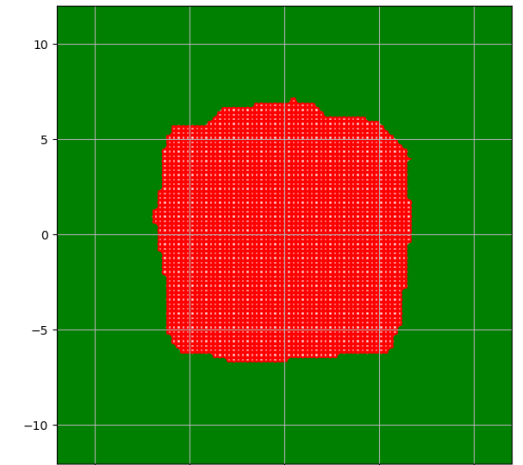
of hidden neurons = 32



of hidden neurons = 64

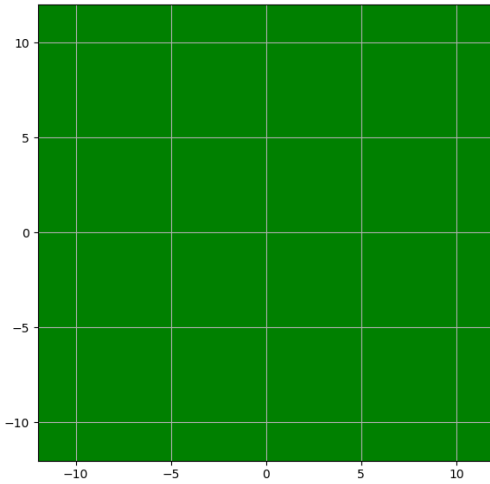


of hidden neurons = 128

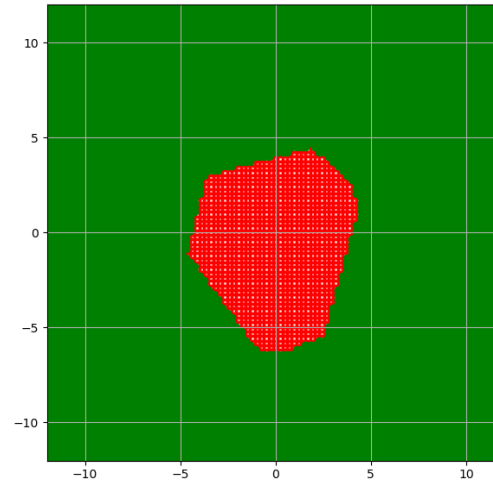


of hidden neurons = 256

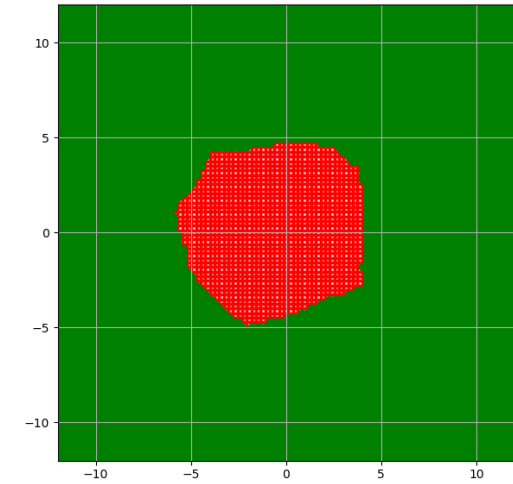
Learned Decision Boundary with Two Hidden Layers



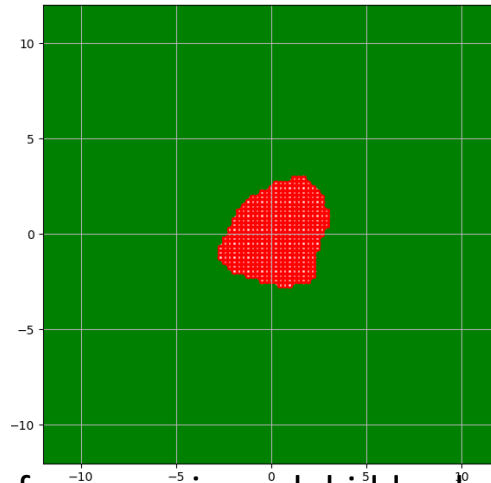
of neurons in each hidden layer = 2



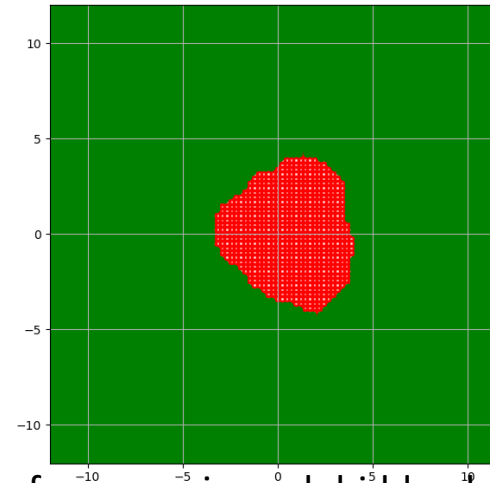
of neurons in each hidden layer = 4



of neurons in each hidden layer = 8



of neurons in each hidden layer = 16



of neurons in each hidden layer = 32