

Universidad de Guadalajara



Thompson Rivers University Seminar

Introduction to Machine Learning

Dr. Carlos Villaseñor

Content

Day 3 (Dr. Carlos Villaseñor)

- Introduction to neural networks
 - Biological neuron
 - Artificial neuron
 - Perceptron algorithm

Rest

- Gradient Descent
- Linear Neurons and Logistic neurons (MSE and BCE)
- One-layer Network (Softmax and CCE)
- How to model the last layer of a net

Rest

- Multi-layer Perceptron
- Keras/TensorFlow framework

Course repository



<https://github.com/Dr-Carlos-Villasenor/TRSeminar.git>

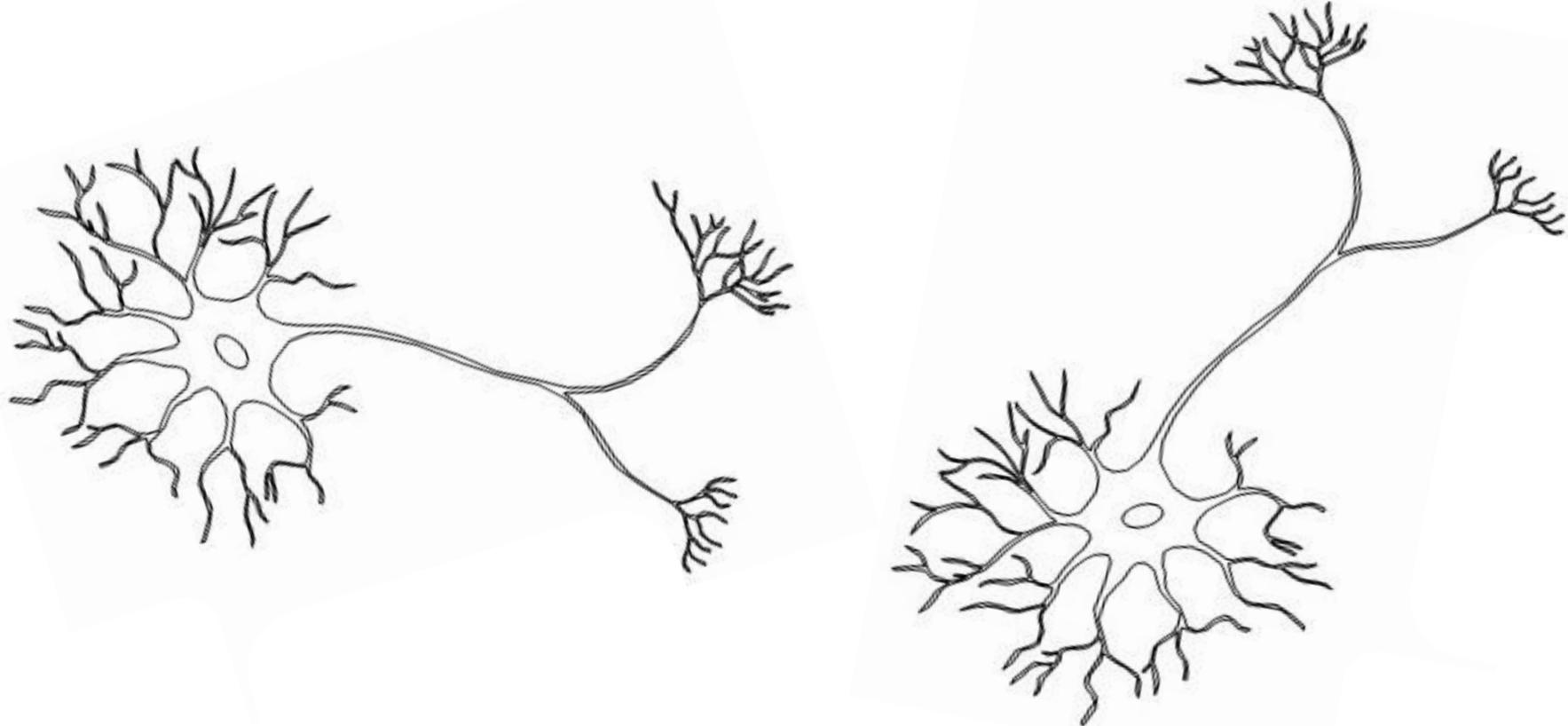
The Biological Neuron

- In 1906, Santiago Ramón y Cajal won the Nobel Prize in Medicine in recognition of his work on the structure of the nervous system.
- He described neurons as information-processing units that connect and form dynamic networks to fulfill all the necessary functions.



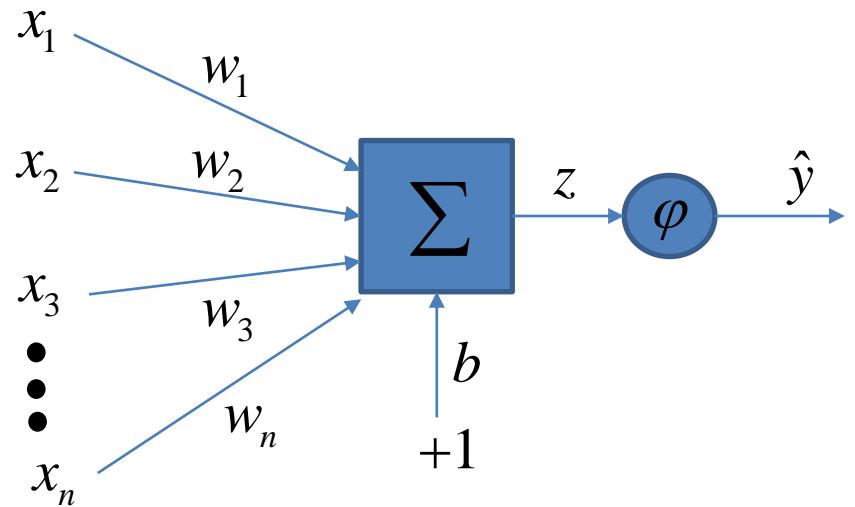
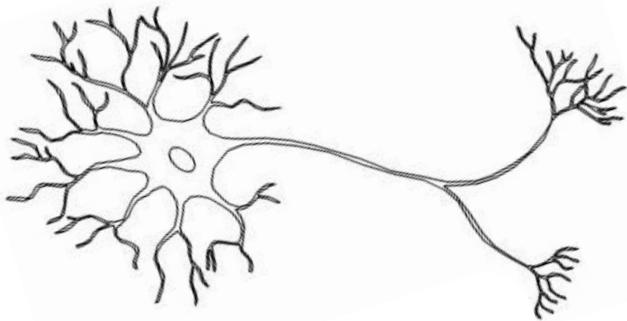
Ramón, S. (1906). Cajal. *The Nobel Prize in Physiology or Medicine.*

The Biological Neuron



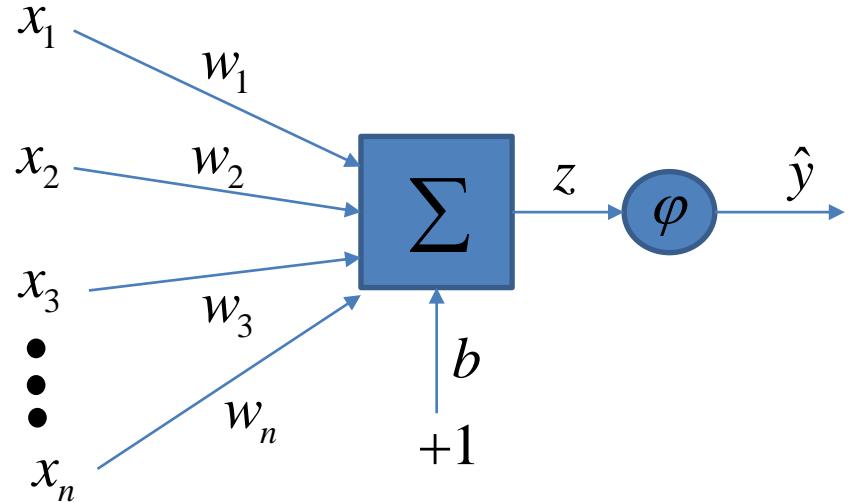
The artificial neuron

In 1943, McCulluch and Pitts developed the first mathematical model of a neuron



McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4), 115-133.

The artificial neuron



$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$\hat{y} = \varphi(z)$$

$$\varphi(z) = \begin{cases} 1 & z \geq 0 \\ 0 & z < 0 \end{cases}$$

McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4), 115-133.

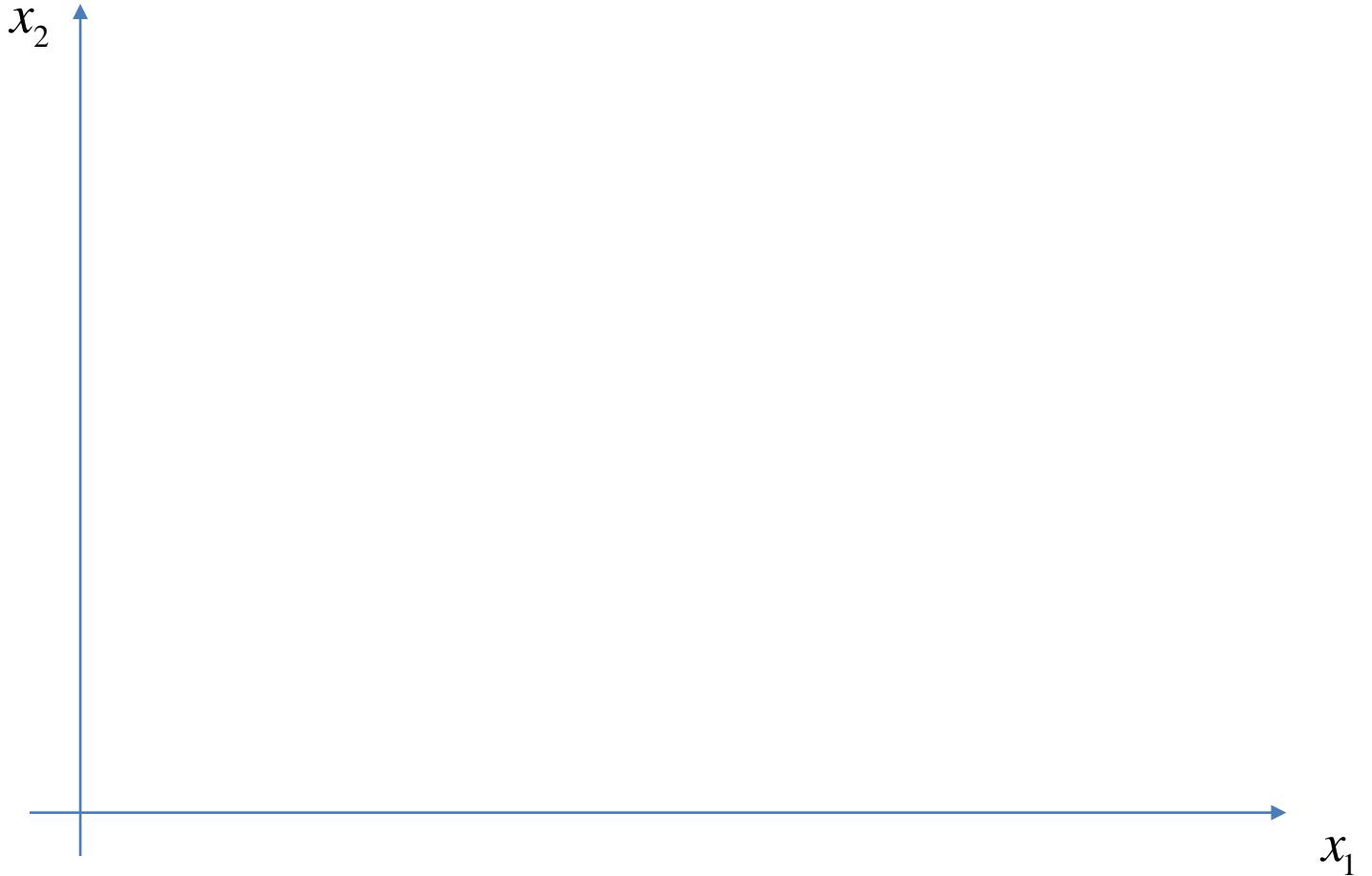
The artificial neuron

$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(p)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(p)} \end{bmatrix}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}^{(1)} & \hat{y}^{(2)} & \dots & \hat{y}^{(p)} \end{bmatrix}$$

McCulloch, W. S., & Pitts, W. (1943). A logical calculus of the ideas immanent in nervous activity. *The bulletin of mathematical biophysics*, 5(4), 115-133.

The artificial neuron



Realizar actividad L01.2

The Perceptron algorithm



Data: $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(p)}, y^{(p)})\}$ $x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \{0,1\}$

For $e \in \{1, 2, \dots, \text{epochs}\}$

For $i \in \{1, 2, \dots, p\}$

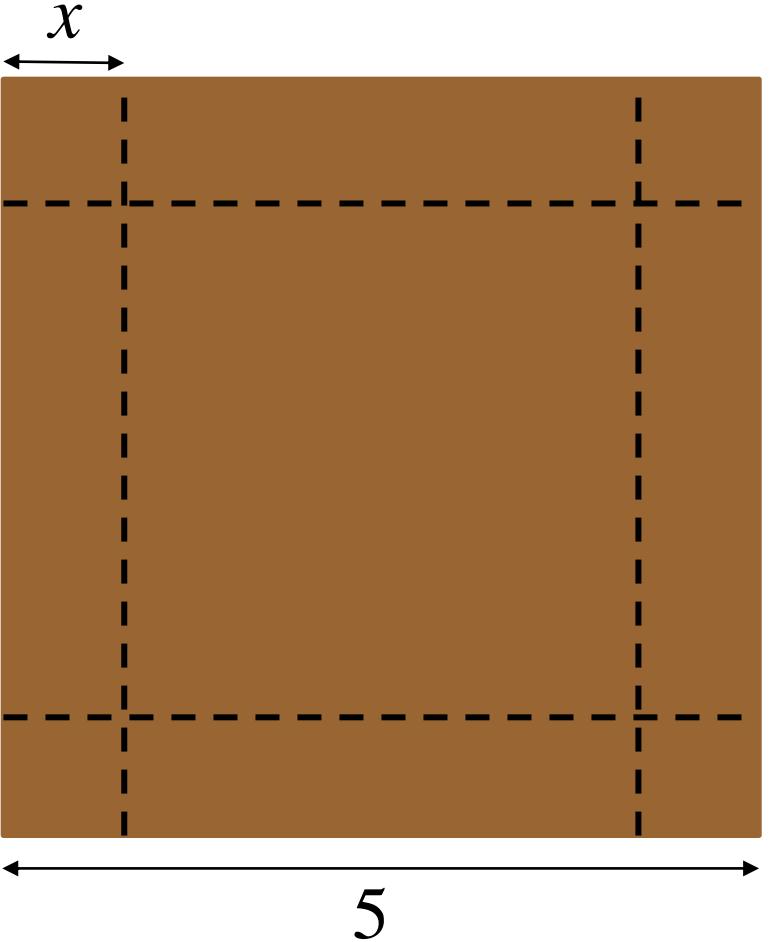
$$\hat{y} = \varphi(w^T x^{(i)} + b)$$

$$w \leftarrow w + \eta(y^{(i)} - \hat{y})x^{(i)}$$

$$b \leftarrow b + \eta(y^{(i)} - \hat{y})$$

Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. Psychological review, 65(6), 386.

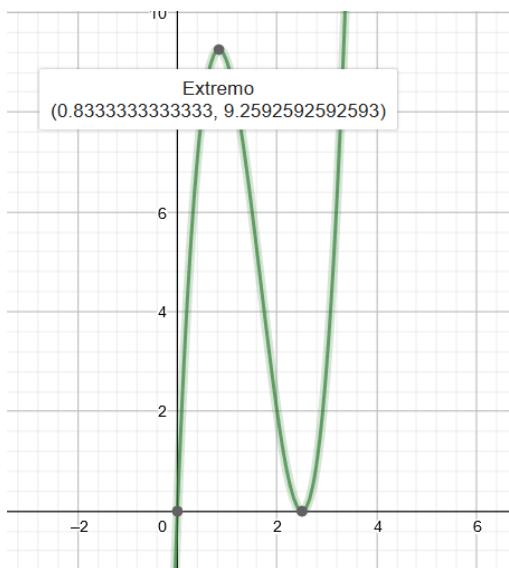
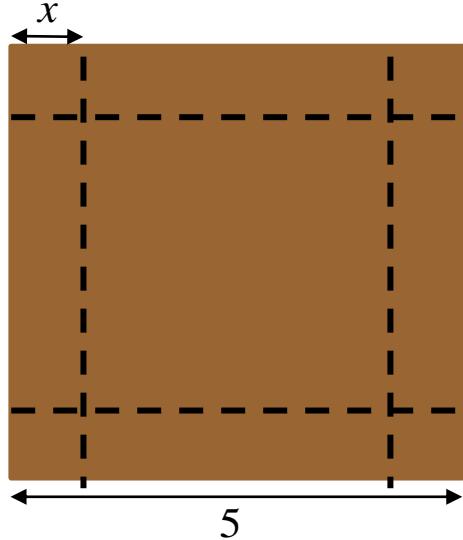
Optimization



$$\max_x f(x)$$

x

Optimization



$$f(x) = x(5 - 2x)^2$$

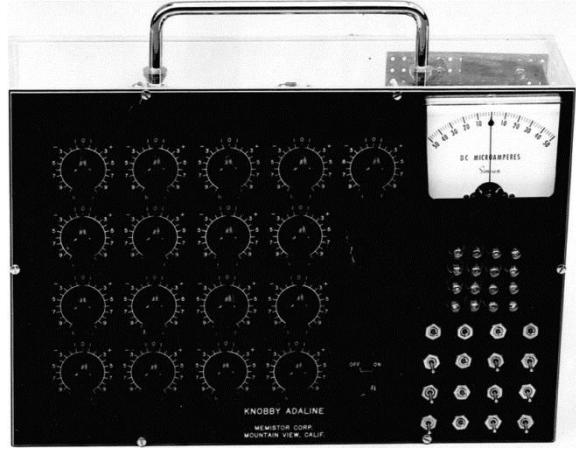
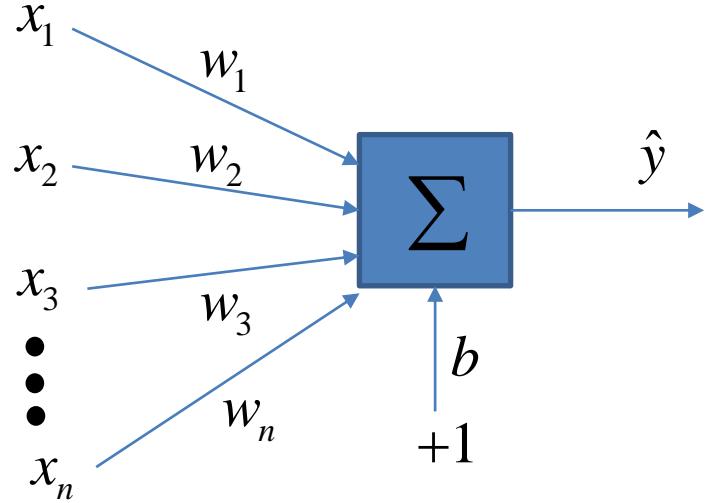
$$f(x) = 4x^3 - 20x^2 + 25x$$

$$\frac{d}{dx} f(x) = 0$$

$$\frac{d}{dx} f(x) = 12x^2 - 40x + 25 = 0$$

$$x_{1,2} = \frac{5}{2}, \frac{5}{6}$$

Linear Neuron

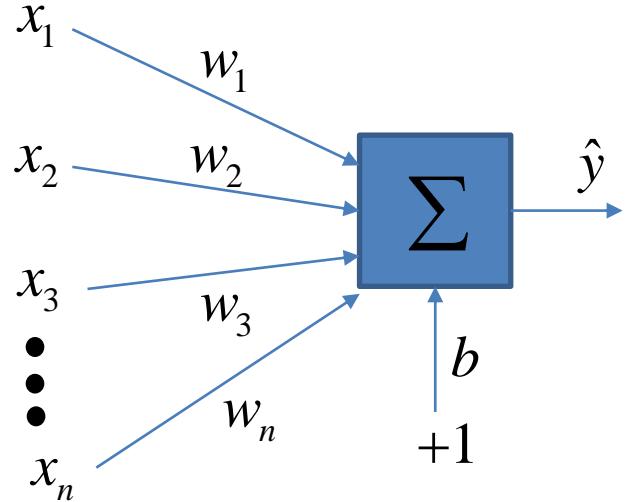


$$\hat{y} = w^T x + b \quad x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \mathbb{R}$$

Widrow, B. (1960). *Adaptive "adaline" Neuron Using Chemical "memistors."*

Widrow, B. (1962). Generalization and information storage in network of ADALINE 'neurons'. *Self-organizing systems-1962*, 435-462.

Linear Neuron



$$D = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(p)}, y^{(p)})\}$$

$$x^{(i)} \in \mathbb{R}^n, y^{(i)} \in \mathbb{R}$$

$$\hat{y} = w^T x + b$$

Online Optimization Problem:

$$\min_{w,b} \frac{1}{2} (y - \hat{y})^2$$

Global optimization problem (for the entire batch of data):

$$\min_{w,b} \frac{1}{2p} \sum_{i=1}^p (y^{(i)} - \hat{y}^{(i)})^2$$

Gradient Descent (SGD)

$$\frac{\delta}{\delta w_i} L(y, \hat{y}) = -(y - \hat{y})x_i$$

$$\nabla_w L(y, \hat{y}) = \begin{bmatrix} \frac{\delta}{\delta w_1} \\ \frac{\delta}{\delta w_2} \\ \vdots \\ \frac{\delta}{\delta w_n} \end{bmatrix} L(y, \hat{y}) = \begin{bmatrix} -(y - \hat{y})x_1 \\ -(y - \hat{y})x_2 \\ \vdots \\ -(y - \hat{y})x_n \end{bmatrix} = -(y - \hat{y}) \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} = -(y - \hat{y})x$$

Propagate Multiple Entries

$$D = \left\{ \left(x^{(1)}, y^{(1)} \right), \left(x^{(2)}, y^{(2)} \right), \dots, \left(x^{(p)}, y^{(p)} \right) \right\}$$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(p)} \\ | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times p} \quad Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(p)} \end{bmatrix} \in \mathbb{R}^{1 \times p}$$

$$\hat{y} = \begin{bmatrix} w_1 & w_2 & \dots & w_n \end{bmatrix} \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} + b$$

Broadcasting

$$\hat{y} = [w_1 \quad w_2 \quad \cdots \quad w_n] \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(p)} \\ | & | & | & | \end{bmatrix} + b$$

For greater comfort, let's redefine w

$$w = [w_1 \quad w_2 \quad \cdots \quad w_n] \in \mathbb{R}^{1 \times n}$$

$$\hat{Y}_{1 \times p} = w_{1 \times n} X_{n \times p} + b_{1 \times 1}$$

Batch Descending Gradient (BGD)

$$E_2 = \frac{1}{2p} \sum_{i=1}^p (y^{(i)} - \hat{y}^{(i)})^2$$

$$X = \begin{bmatrix} | & | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(p)} \\ | & | & | & | \end{bmatrix} \in \mathbb{R}^{n \times p} \quad Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(p)} \end{bmatrix} \in \mathbb{R}^{1 \times p}$$

$$\nabla_w E_2 \quad \frac{\partial}{\partial b} E_2$$

Batch Descending Gradient (BGD)

$$E_2 = \frac{1}{2p} \sum_{i=1}^p (y^{(i)} - \hat{y}^{(i)})^2$$

$$\hat{y}^{(i)} = w x^{(i)} + b$$

$$\nabla_w E_2 = \nabla_w \left[\frac{1}{2p} \sum_{i=1}^p (y^{(i)} - \hat{y}^{(i)})^2 \right]$$

$$\nabla_w E_2 = \frac{1}{p} \sum_{i=1}^p \left[\nabla_w \frac{1}{2} (y^{(i)} - \hat{y}^{(i)})^2 \right]$$

$$\nabla_w E_2 = \frac{1}{p} \sum_{i=1}^p \left[-(y^{(i)} - \hat{y}^{(i)}) x^{(i)} \right]$$

$$\nabla_w E_2 = -\frac{1}{p} \sum_{i=1}^p \left[(y^{(i)} - \hat{y}^{(i)}) x^{(i)} \right]$$

Batch Descending Gradient (BGD)

For $e \in \{1, 2, \dots, \text{epochs}\}$

$$\hat{Y} = wX + b$$

$$w \in \mathbb{R}^{1 \times n}$$

$$w \leftarrow w + \frac{\eta}{p} (Y - \hat{Y}) X^T$$

$$X \in \mathbb{R}^{n \times p}$$

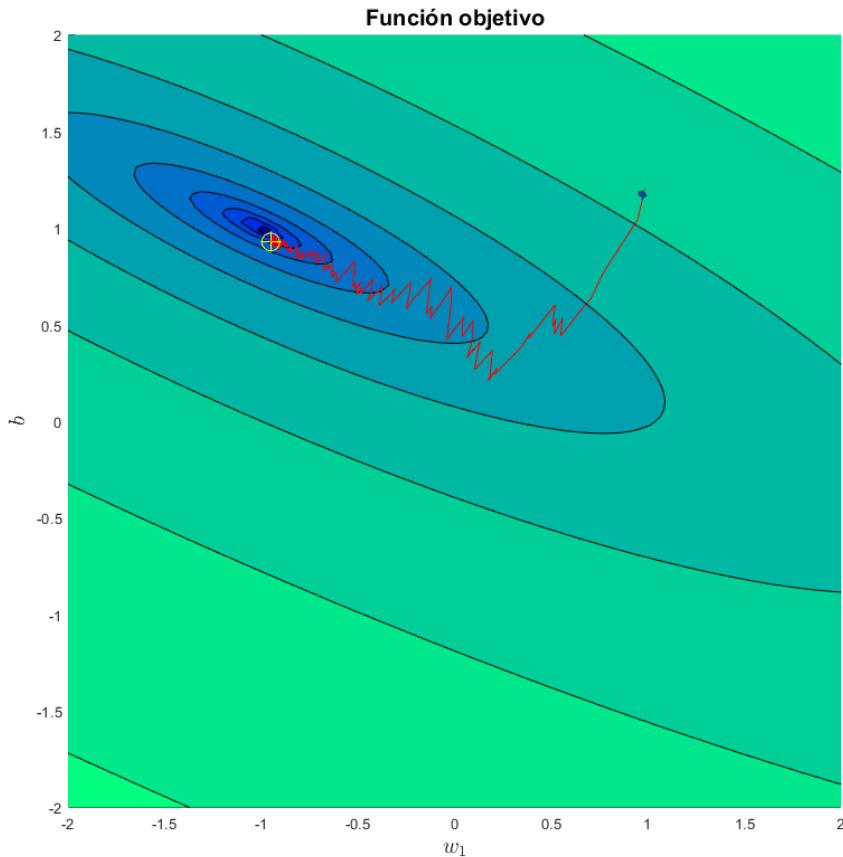
$$b \leftarrow b + \frac{\eta}{p} \sum_{i=1}^p (y^{(i)} - \hat{y}^{(i)})$$

$$Y, \hat{Y} \in \mathbb{R}^{1 \times p}$$

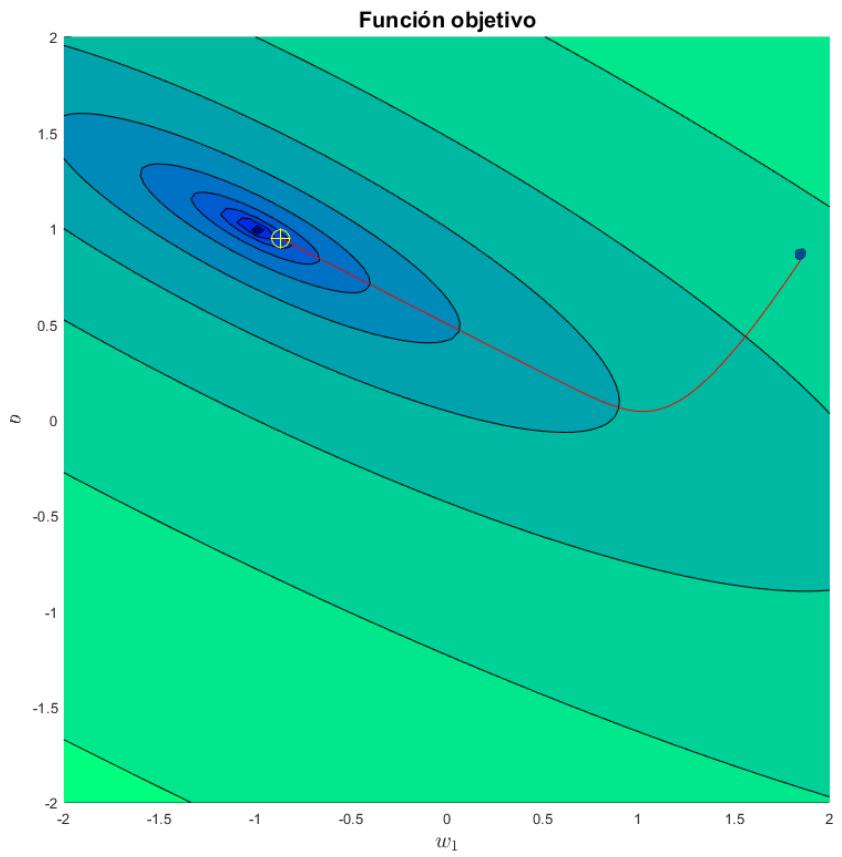
$$b, \eta \in \mathbb{R}$$

Mini-Batch Descending Gradient (mBGD)

SGD

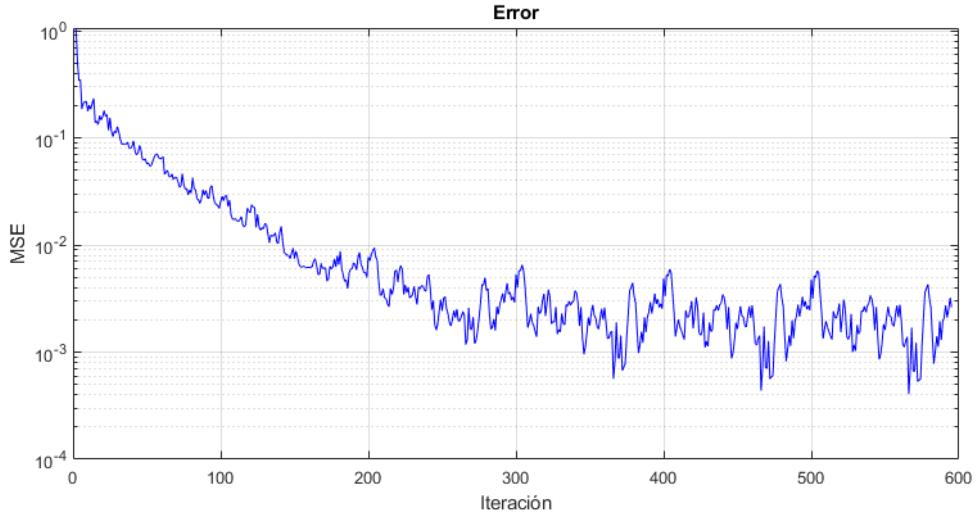


BGD

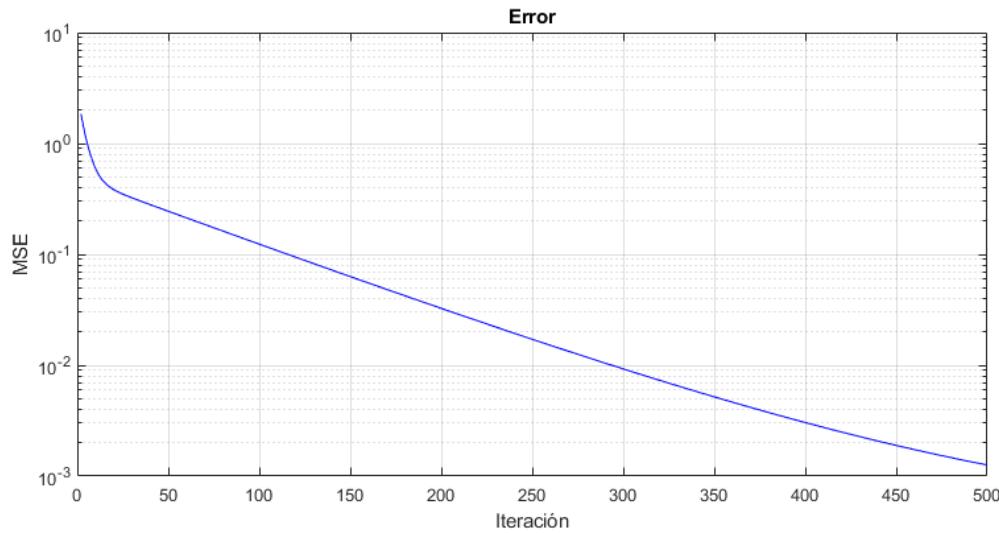


Mini-Batch Descending Gradient (mBGD)

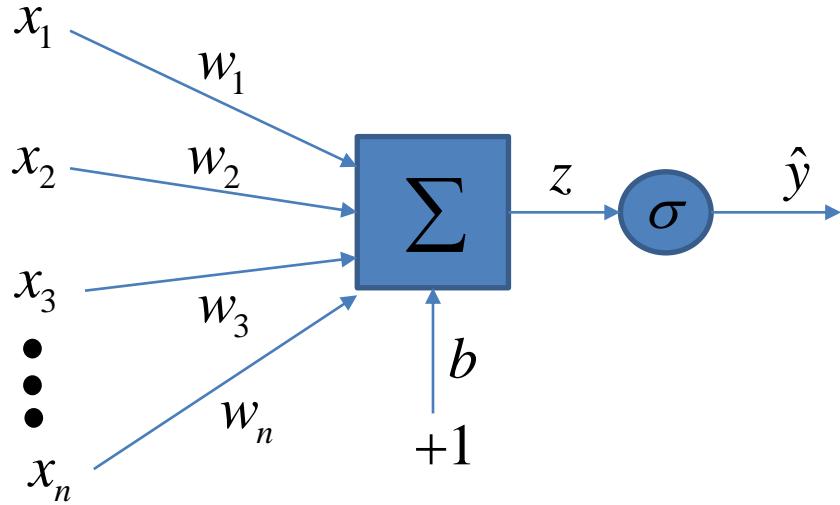
SGD



BGD

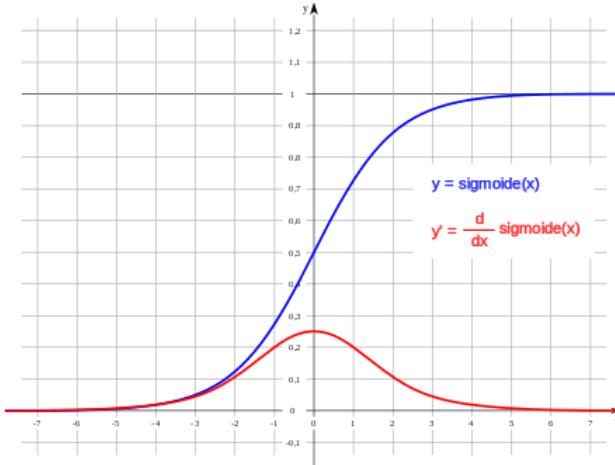


Sigmoid Neuron



$$z = w^T x + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Binary Cross Entropy

Optimization Problem

$$\min_{w,b} -\frac{1}{2} \left[y \log(\hat{y}) + (1-y) \log(1-\hat{y}) \right]$$

Learning algorithm

$$w \leftarrow w - \eta \nabla_w L(y, \hat{y})$$

$$b \leftarrow b - \eta \frac{\delta}{\delta b} L(y, \hat{y})$$

Update equations

$$w \leftarrow w + \eta (y - \hat{y}) x$$

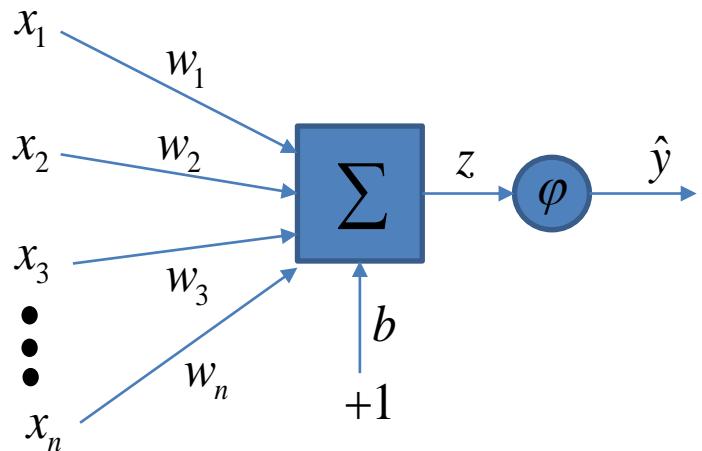
$$b \leftarrow b + \eta (y - \hat{y})$$

$$z = w^T x + b$$

$$\hat{y} = \sigma(z) = \frac{1}{1 + e^{-z}}$$

Delta Rule

- Suppose $\varphi(z)$ is a differentiable function then



$$\hat{y} = \varphi(w^T x + b)$$

$$\min_{w,b} \frac{1}{2} (y - \hat{y})^2$$

$$w \leftarrow w + \eta (y - \hat{y}) \frac{d\varphi(z)}{dz} x$$

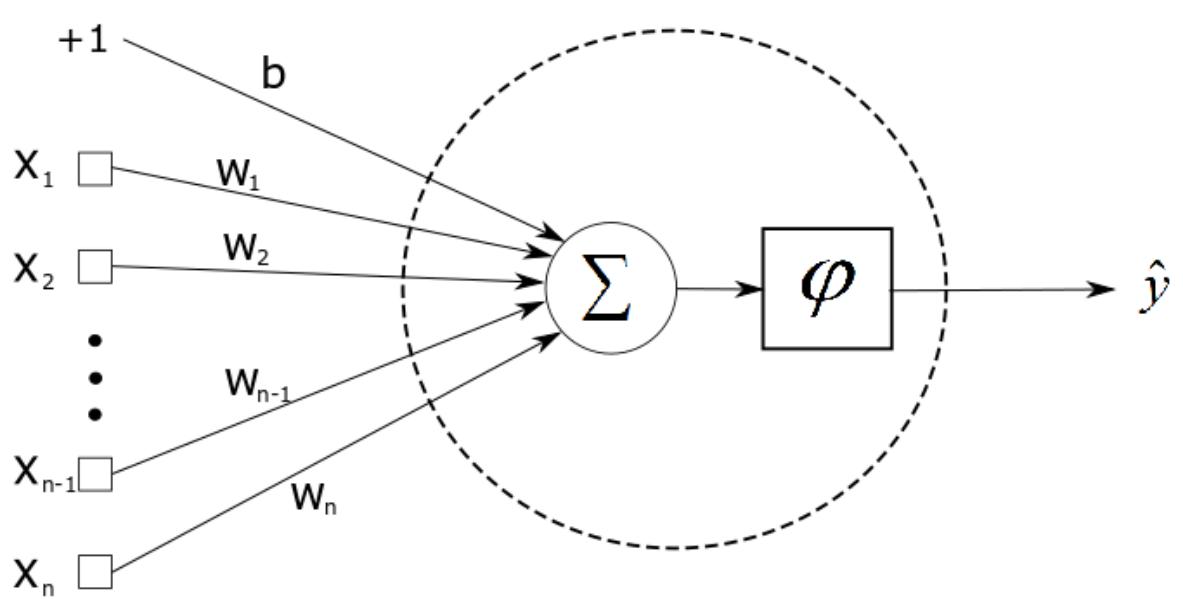
$$b \leftarrow b + \eta (y - \hat{y}) \frac{d\varphi(z)}{dz}$$

Delta Rule



Function name	Function	Derivative
Linear or identity	z	1
Sigmoid	$\frac{1}{1 + e^{-z}}$	$\hat{y}(1 - \hat{y})$
Hyperbolic tangent	$\tanh(z)$	$(1 + \hat{y})(1 - \hat{y})$
Rectified Linear Unit (ReLU)	$\max(0, z)$	$\begin{cases} 1 & z > 0 \\ 0 & z \leq 0 \end{cases}$

One-Layer Network

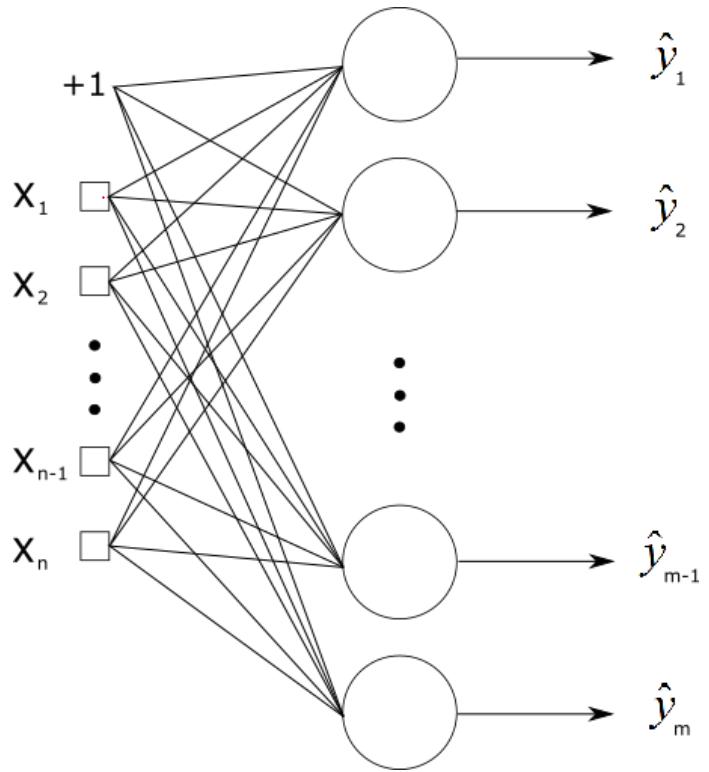


$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^n$$

$$b \in \mathbb{R}$$

$$\hat{y} = \varphi(w^T x + b) = \varphi \left(\begin{bmatrix} w_1 & w_2 & \cdots & w_n \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + b \right)$$

One-Layer Network



$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^n \quad \hat{y} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} \in \mathbb{R}^m \quad b = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \in \mathbb{R}^m$$

$$W = \begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

One-Layer Network

$$\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix} = \begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix}$$

$$\begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_m \end{bmatrix} = \varphi \left(\begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix} \right)$$

$$z = wx + b$$

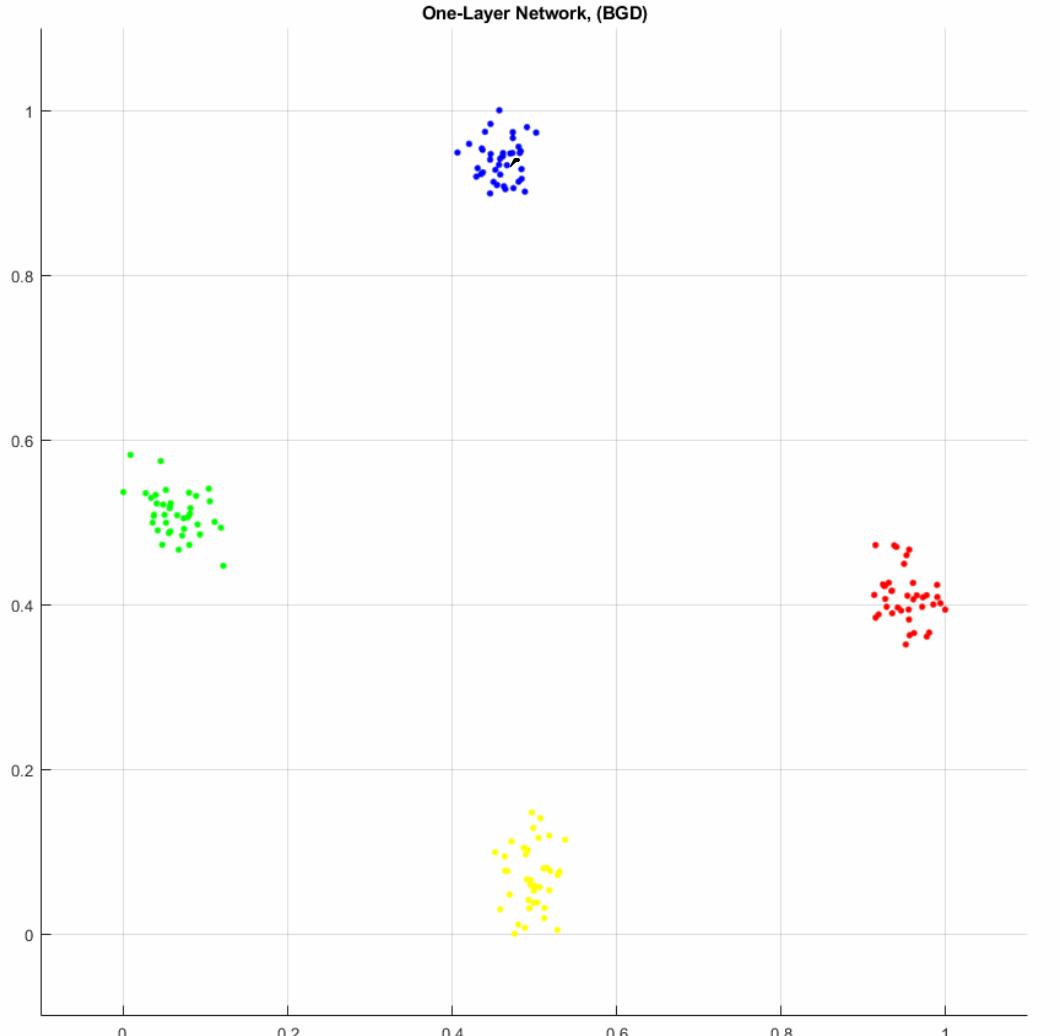
$$y = \varphi(z)$$

One-Layer Network

$$\hat{Y} = \varphi \left(\begin{bmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{m,1} & w_{m,2} & \cdots & w_{m,n} \end{bmatrix} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(p)} \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \cdots & x_n^{(p)} \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{bmatrix} \right)$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1^{(1)} & \hat{y}_1^{(2)} & \cdots & \hat{y}_1^{(p)} \\ \hat{y}_2^{(1)} & \hat{y}_2^{(2)} & \cdots & \hat{y}_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{y}_m^{(1)} & \hat{y}_m^{(2)} & \cdots & \hat{y}_m^{(p)} \end{bmatrix}$$

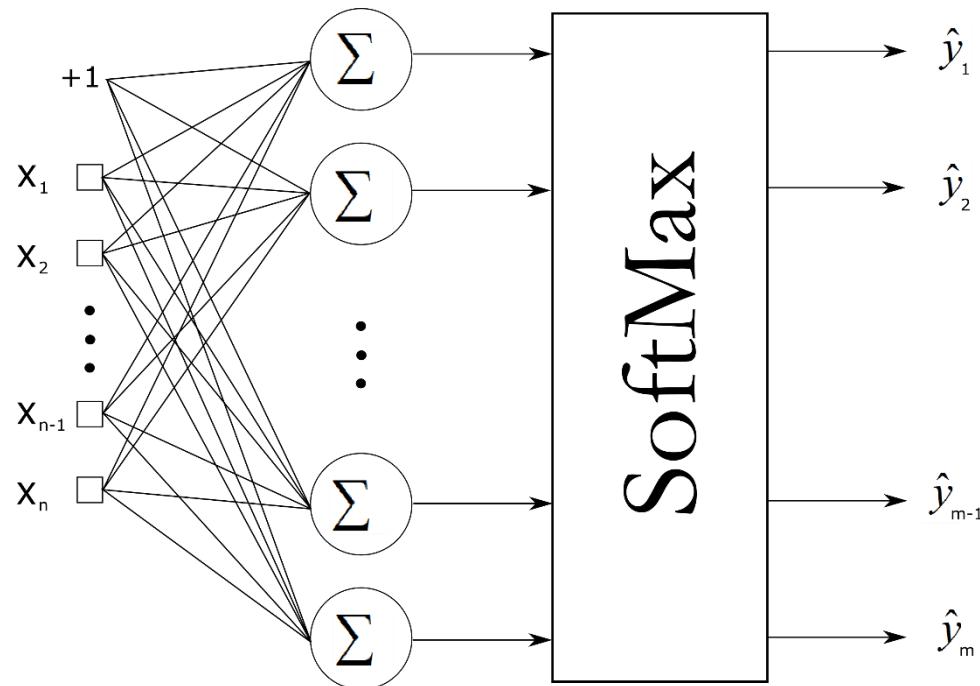
One vs All



SoftMax (Categorical Cross Entropy)

$$E = -\sum_{j=1}^m y_j \log(\hat{y}_j)$$

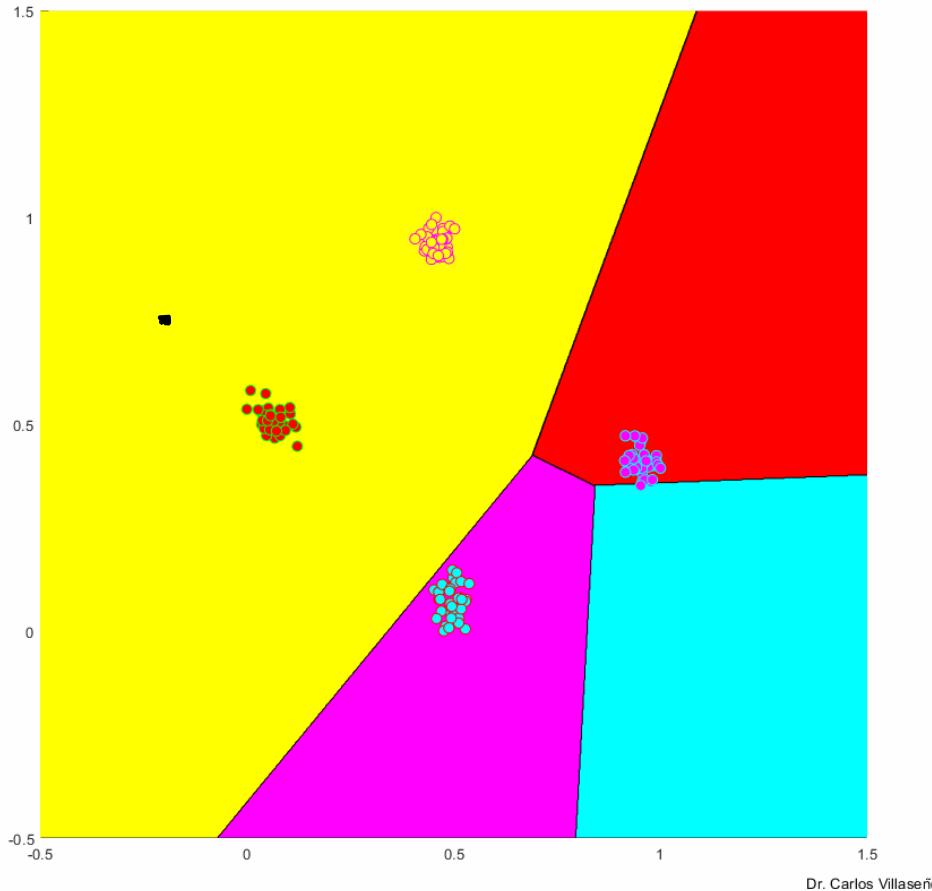
$$E = -\frac{1}{p} \sum_{i=1}^p \sum_{j=1}^m y_j^{(i)} \log(\hat{y}_j^{(i)})$$



$$S(z)_j = \frac{e^{z_j}}{\sum_{i=1}^m e^{z_i}}$$

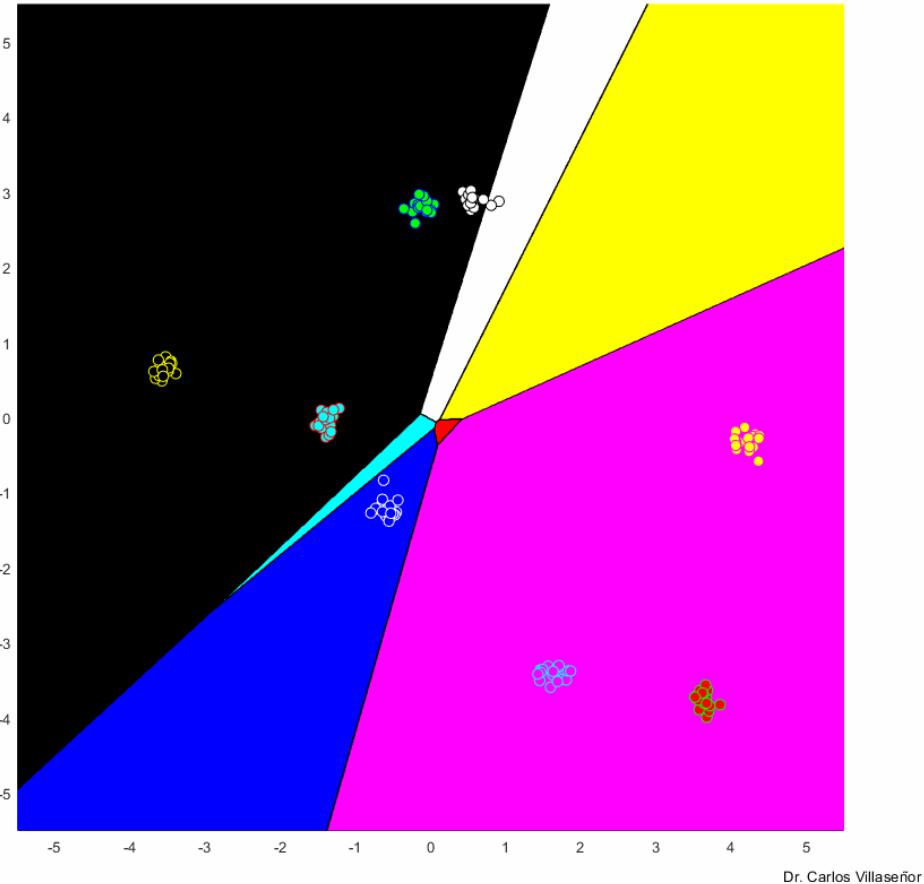
Softmax Layer

Softmax Regression



Softmax Layer

Softmax Regression

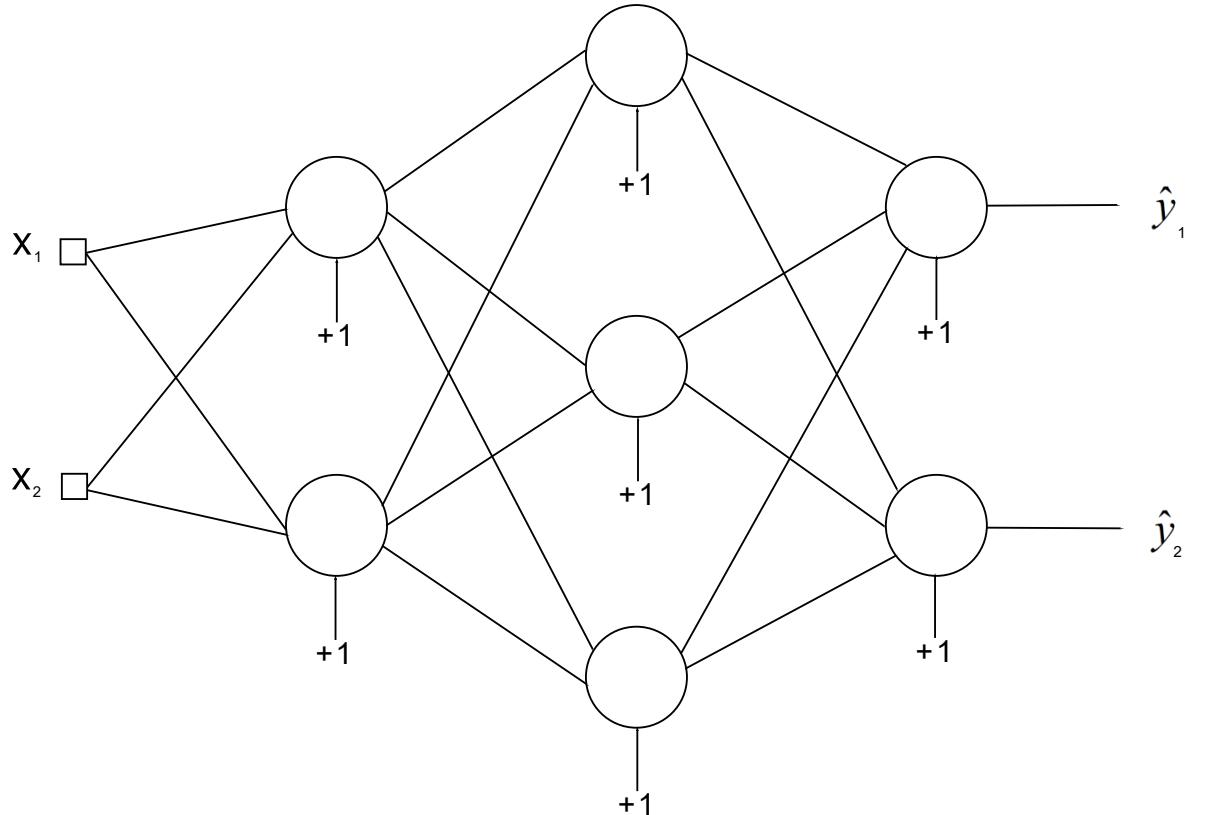


Summary

Problem	Activation function	Loss Function
Regression	Linear	Mean square error(MSE)
Multi-label classification	Sigmoid	Binary Cross Entropy Error(BCE)
One-winner qualification	SoftMax	Categorical cross-entropy error(CCE)

Dense Network

Enumeration of layers and neurons:



$$L =$$

$$n_1 =$$

$$n_0 =$$

$$n_L =$$

n_l : Number of neurons or inputs in the layer number “l”

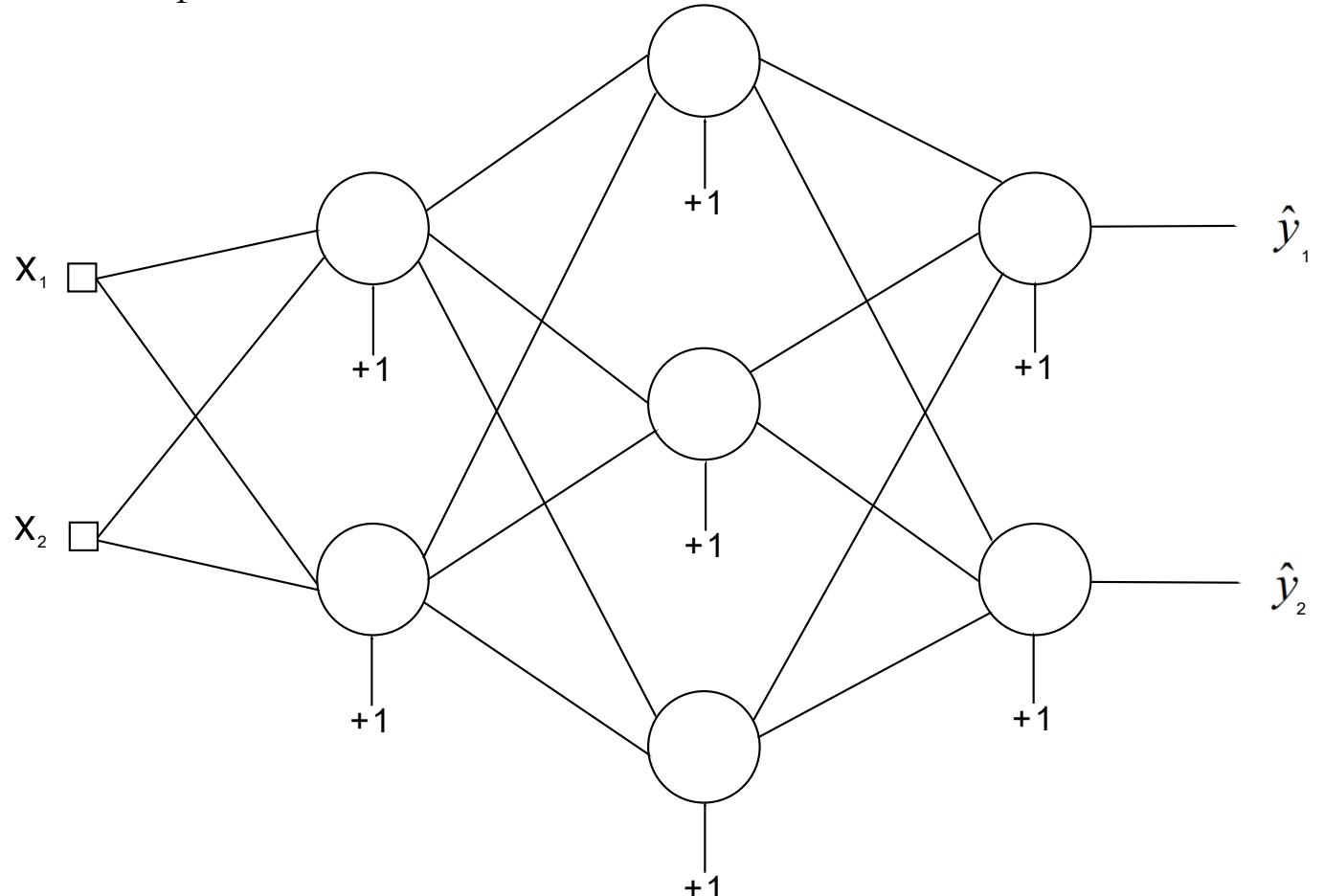
L : Last Layer

Dense Network

$$w^{[l]} = \begin{bmatrix} w_{1,1}^{[l]} & w_{1,2}^{[l]} & \cdots & w_{1,n_{l-1}}^{[l]} \\ w_{2,1}^{[l]} & w_{2,2}^{[l]} & \cdots & w_{2,n_{l-1}}^{[l]} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n_l,1}^{[l]} & w_{n_l,2}^{[l]} & \cdots & w_{n_l,n_{l-1}}^{[l]} \end{bmatrix} \in \mathbb{R}^{n_l \times n_{l-1}} \quad b^{[l]} = \begin{bmatrix} b_1^{[l]} \\ b_2^{[l]} \\ \vdots \\ b_{n_l}^{[l]} \end{bmatrix} \in \mathbb{R}^{n_l}$$

Dense Network

Inputs and outputs:



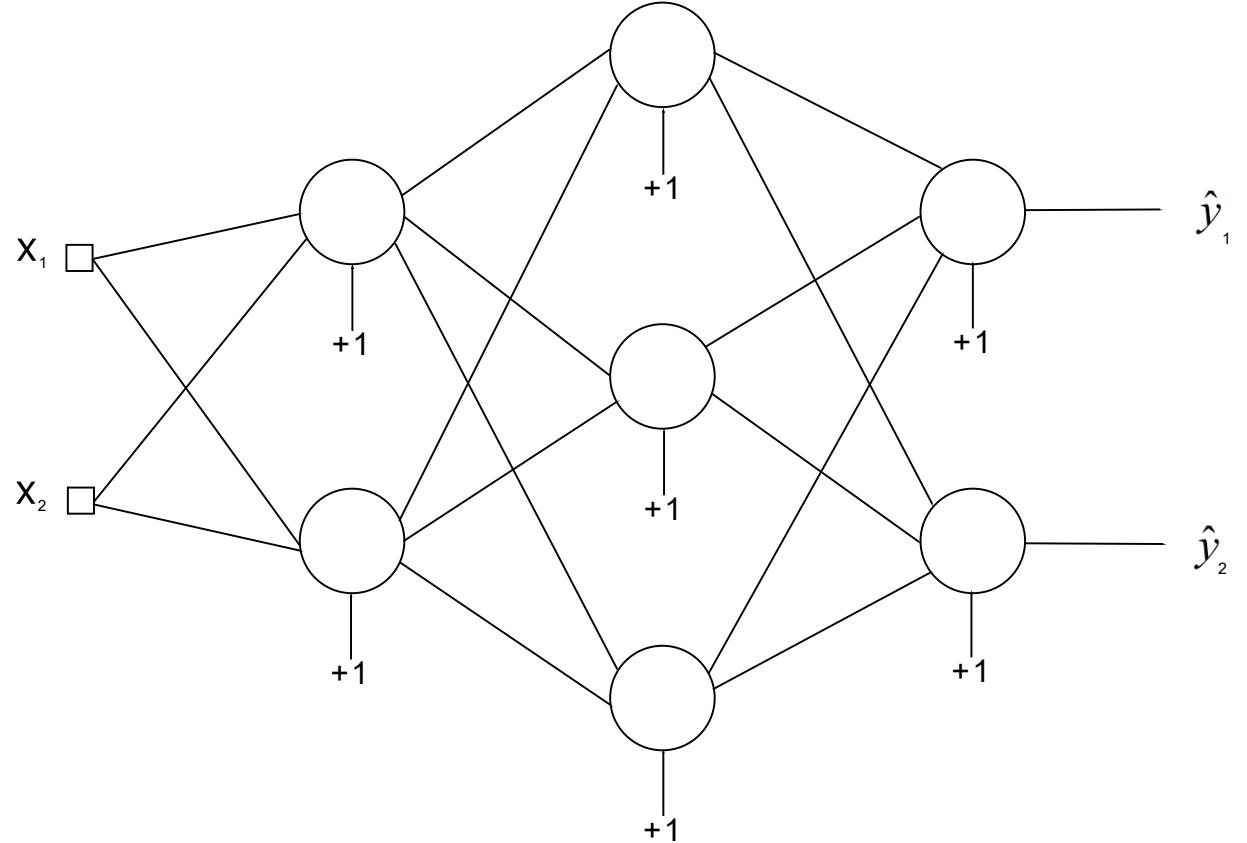
$a_i^{[l]}$: Output of the " l " layer of the " i " neuron

$$a_i^{[L]} = \hat{y}_i$$

$$a_i^{[0]} = x_i$$

$$a^{[l]} = \begin{bmatrix} a_1^{[l]} \\ a_2^{[l]} \\ \vdots \\ a_{n_l}^{[l]} \end{bmatrix} \in \mathbb{R}^{n_l}$$

Propagation



$$a^{[0]} = x$$

$$z^{[1]} = w^{[1]}a^{[0]} + b^{[1]}$$

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

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

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

$$z^{[3]} = w^{[3]}a^{[2]} + b^{[3]}$$

$$a^{[3]} = \varphi^{[3]}(z^{[3]})$$

$$\hat{y} = a^{[3]}$$

Propagation

$$a^{[0]} = x$$

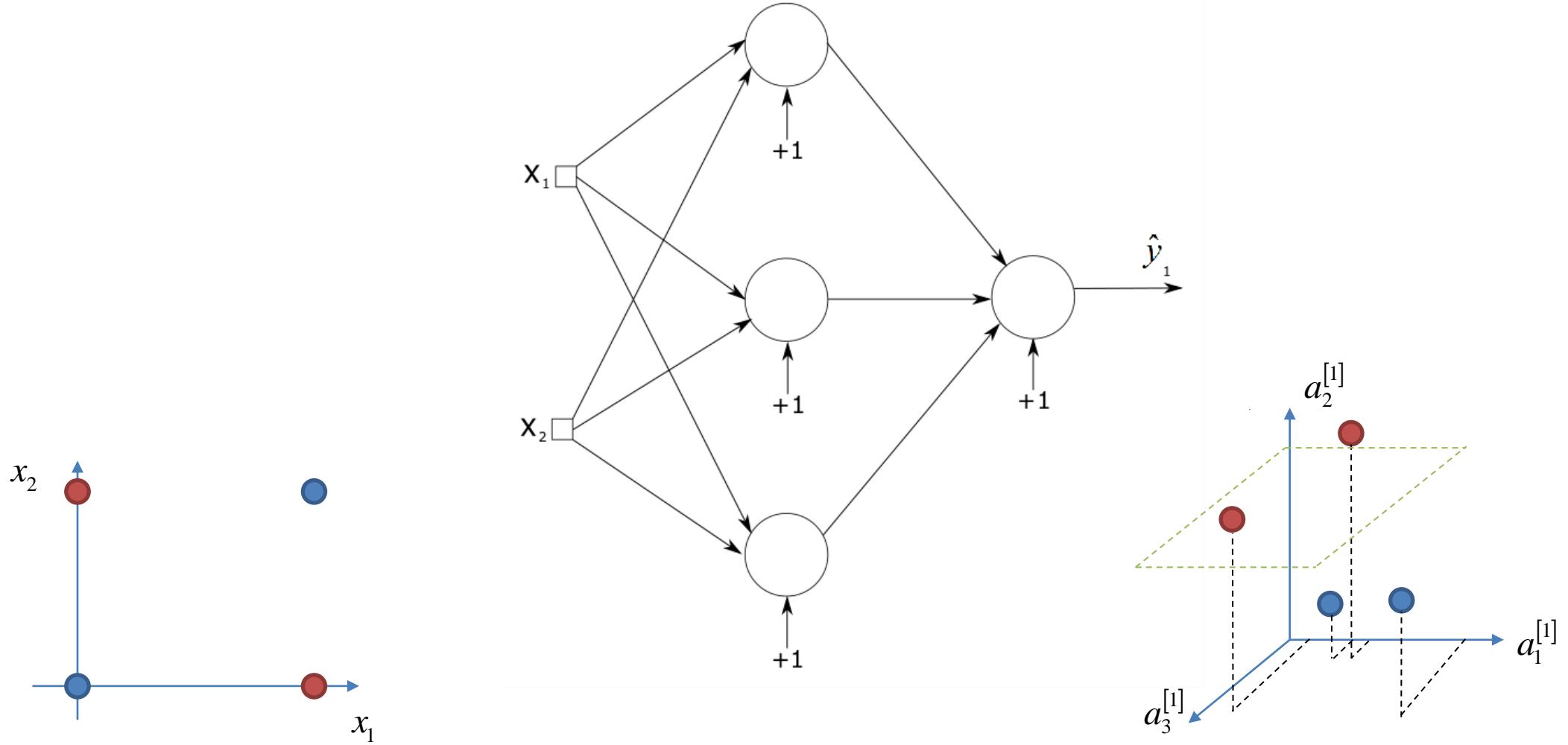
foreach $l \in \{1, \dots, L\}$

$$z^{[l]} = w^{[l]}a^{[l-1]} + b^{[l]}$$

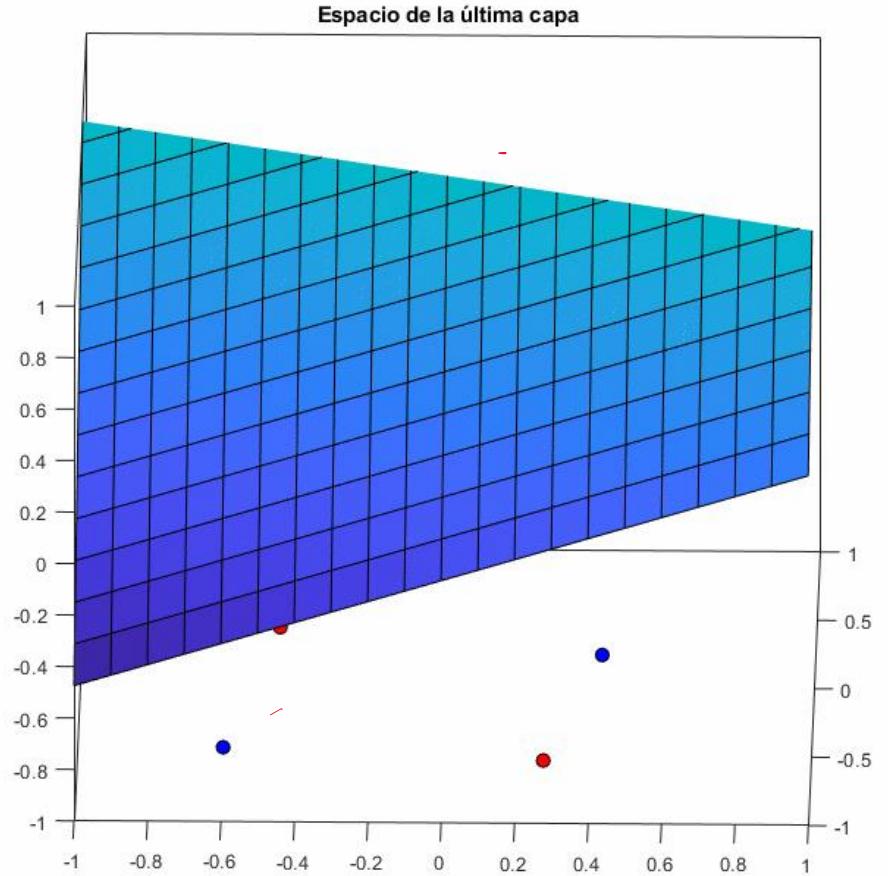
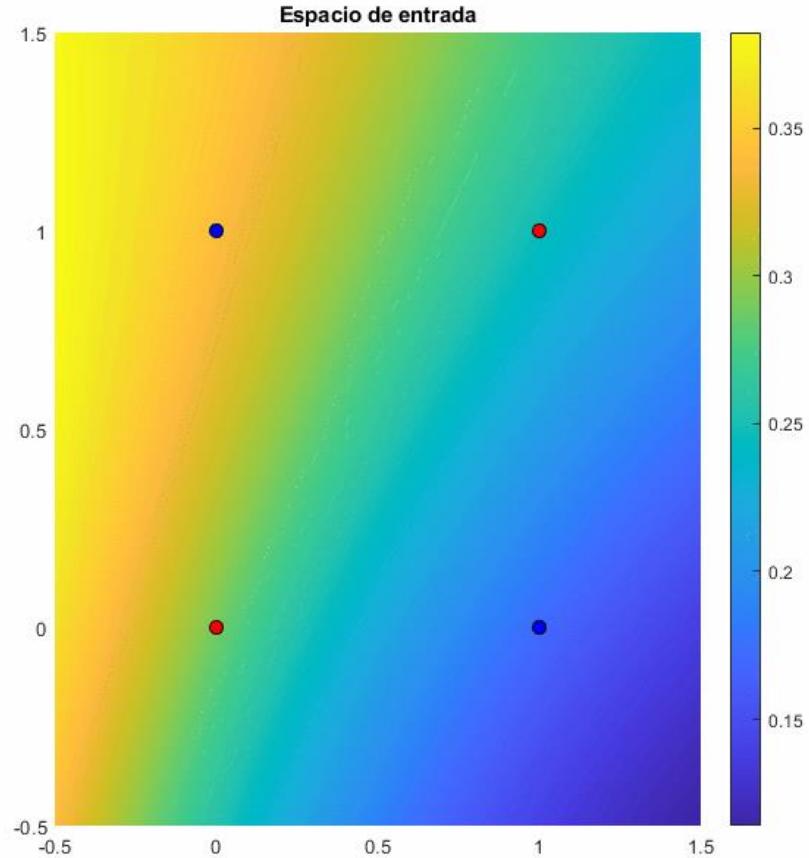
$$a^{[l]} = \varphi^{[l]}(z^{[l]})$$

$$\hat{y} = a^{[L]}$$

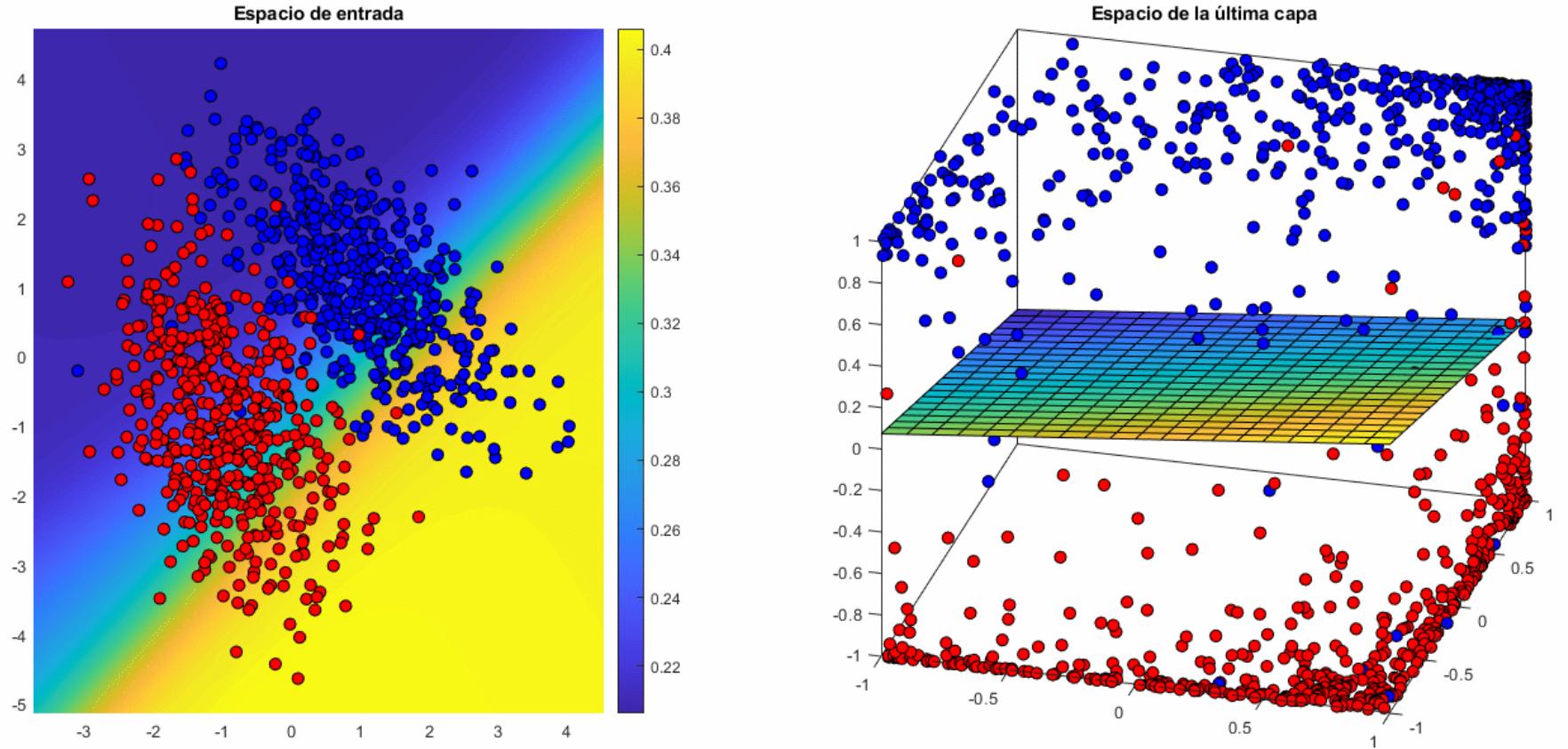
Geometric Representation



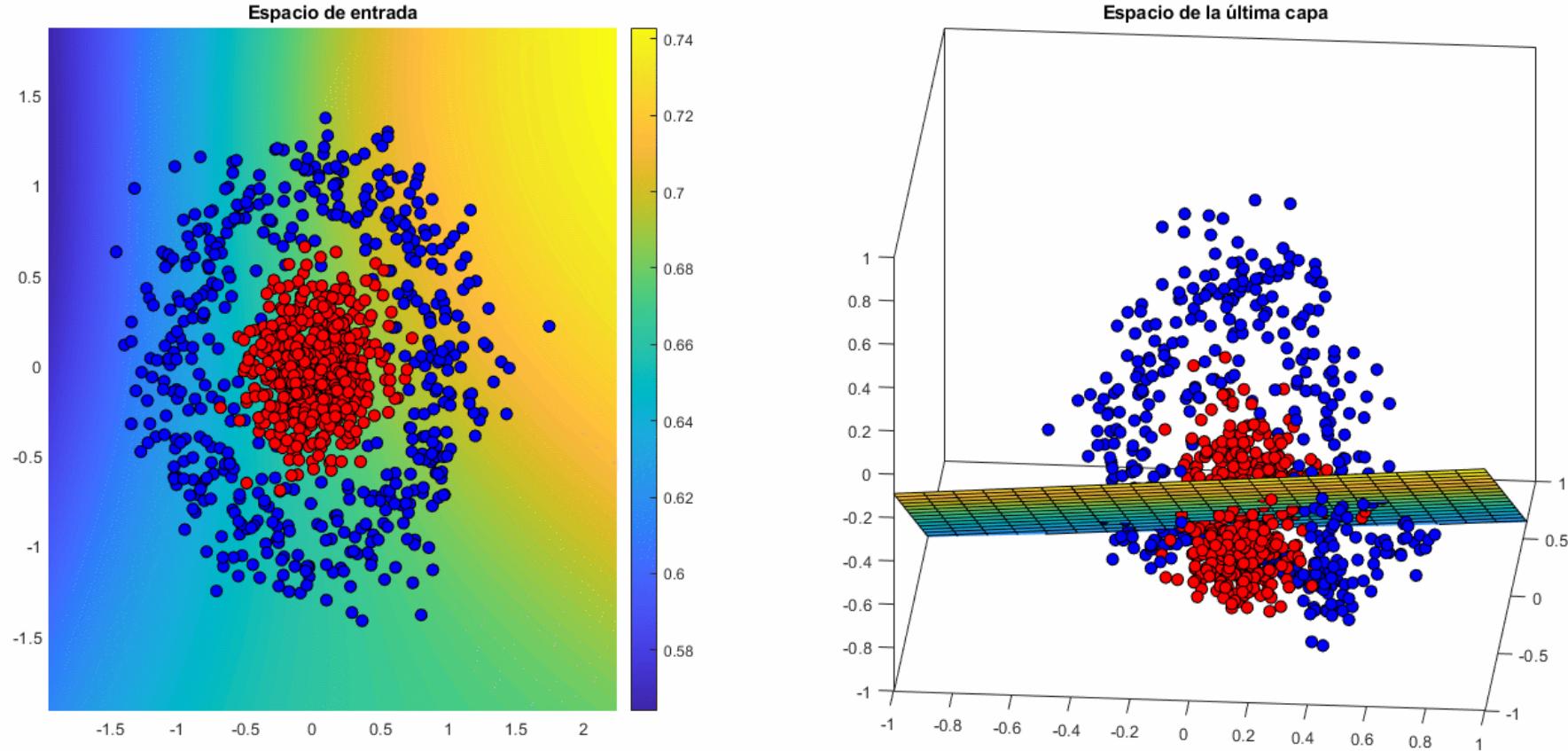
Geometric Representation



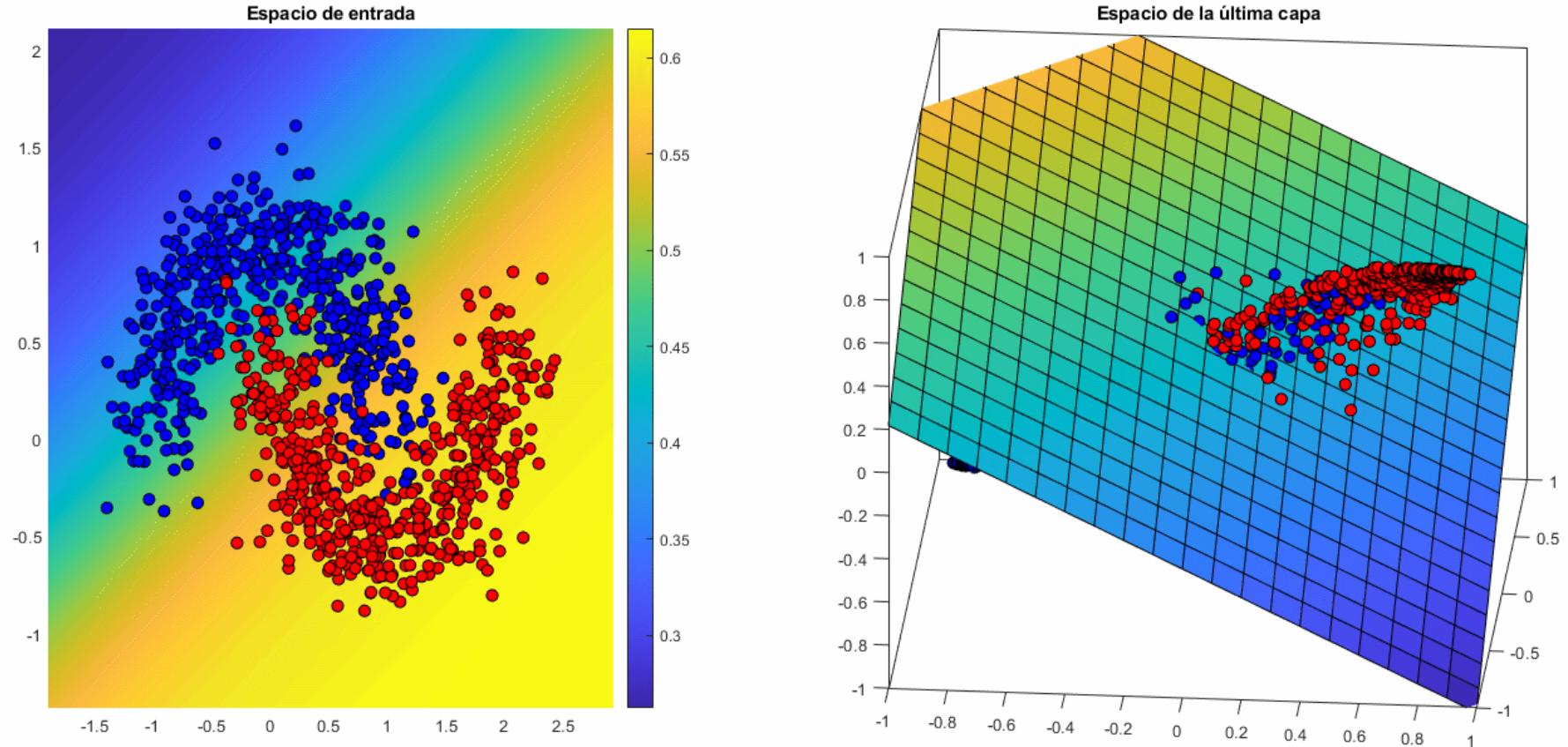
Geometric Representation



Geometric Representation

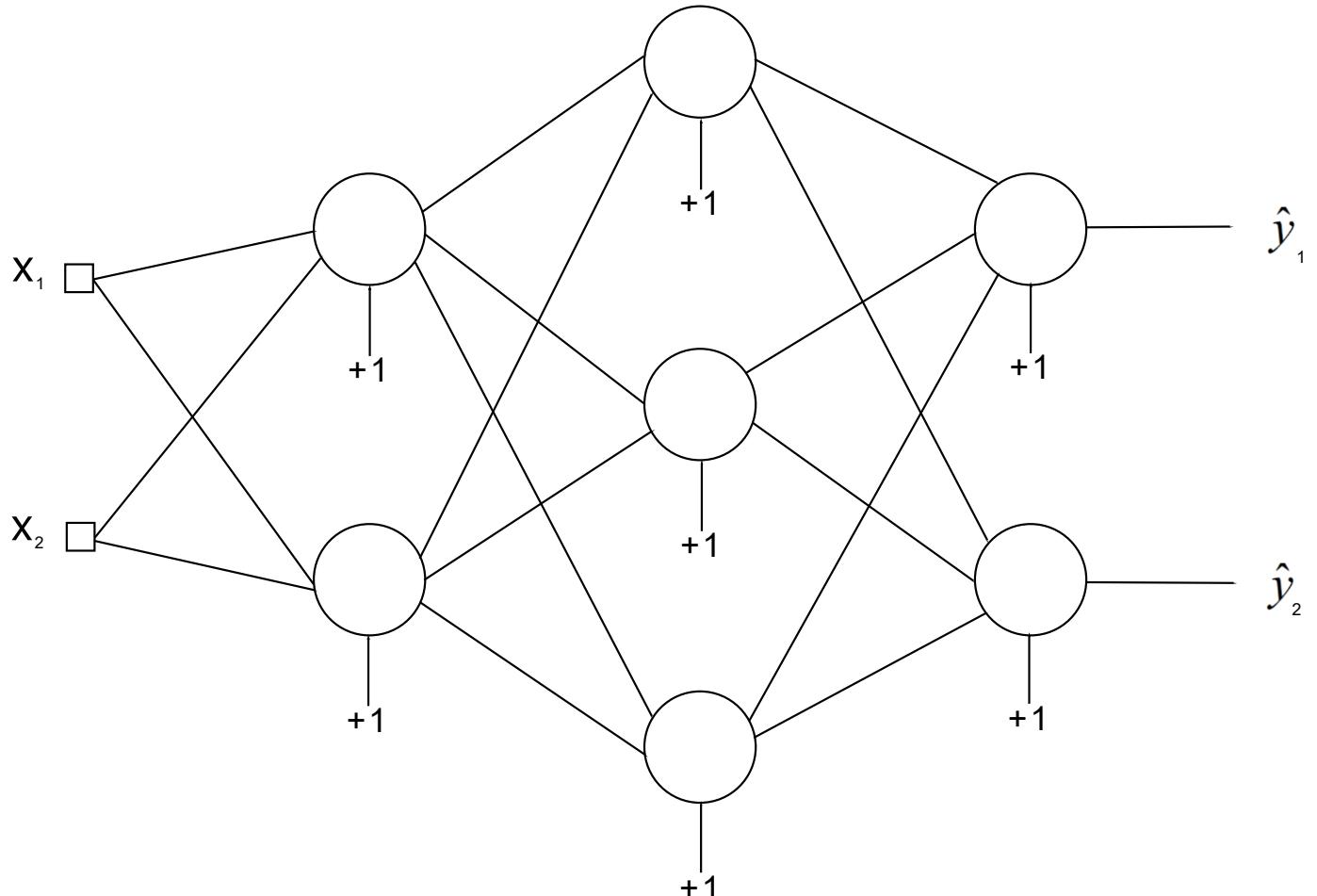


Geometric Representation



Backpropagation

$$da_i^{[L]} = \frac{d\phi_i^{[L]}}{dz}$$



$\delta_i^{[l]}$: The local gradient of the "i" neuron in the "l" layer

$$\delta_i^{[L]} = (y_i - a_i^{[L]}) da_i^{[L]}$$

Backpropagation

- Let's take a specific example:

$$\delta_1^{[2]} = \left[w_{1,1}^{[3]} \delta_1^{[3]} + w_{2,1}^{[3]} \delta_2^{[3]} \right] da_1^{[2]}$$

$$\delta_2^{[1]} = \left[w_{1,2}^{[2]} \delta_1^{[2]} + w_{2,2}^{[2]} \delta_2^{[2]} + w_{3,2}^{[2]} \delta_3^{[2]} \right] da_2^{[1]}$$

Backpropagation

- We finally come to these equations

$$\delta^{[l]} = \begin{cases} (y - a^{[l]}) \odot da^{[l]} & \text{Si } l = L \\ [w^{T[l+1]} \delta^{[l+1]}] \odot da^{[l]} & \text{Si } 0 < l < L \end{cases}$$

$$w^{[l]} \leftarrow w^{[l]} + \eta \delta^{[l]} a^{T[l-1]}$$

$$b^{[l]} \leftarrow b^{[l]} + \eta \delta^{[l]}$$

Training Algorithm

$$a^{[0]} = x$$

foreach $l \in \{1, \dots, L\}$

$$z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]}, da^{[l]} = \varphi^{[l]}(z^{[l]})$$

foreach $l \in \{L, \dots, 1\}$

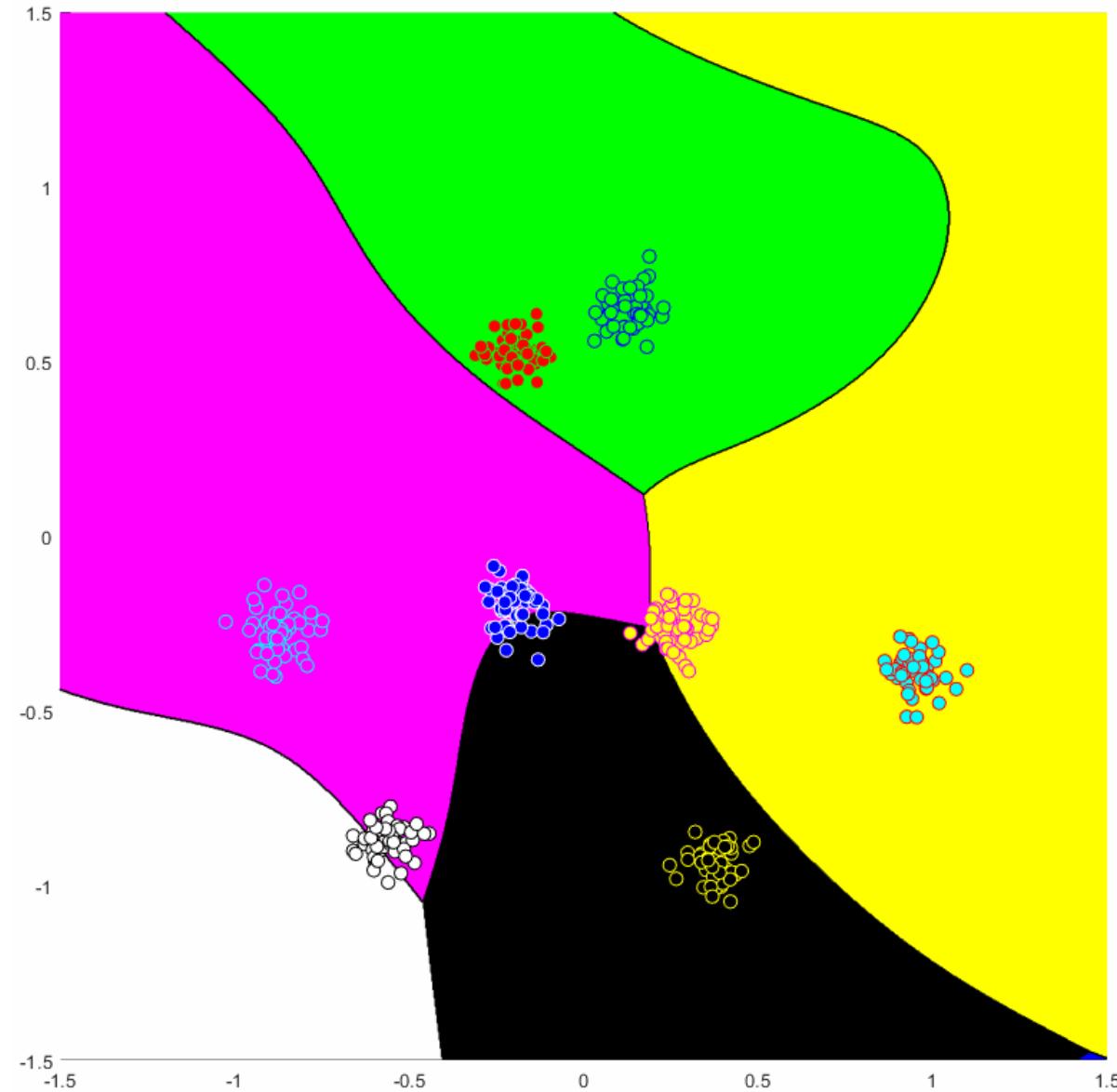
$$\delta^{[l]} = \begin{cases} (y - a^{[l]}) \odot da^{[l]} & \text{Si } l = L \\ [w^{T[l+1]} \delta^{[l+1]}] \odot da^{[l]} & \text{Si } 0 < l < L \end{cases}$$

foreach $l \in \{1, \dots, L\}$

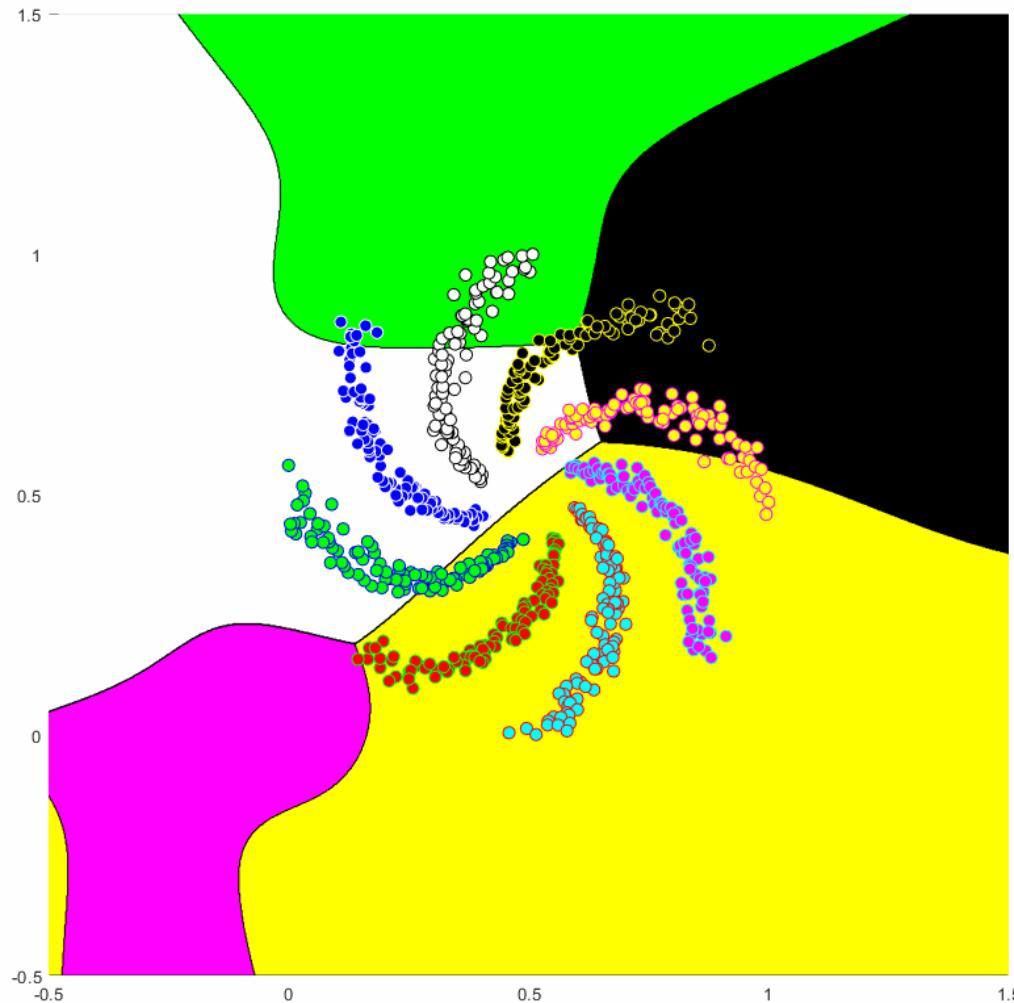
$$w^{[l]} \leftarrow w^{[l]} + \eta \delta^{[l]} a^{T[l-1]}$$

$$b^{[l]} \leftarrow b^{[l]} + \eta \delta^{[l]}$$

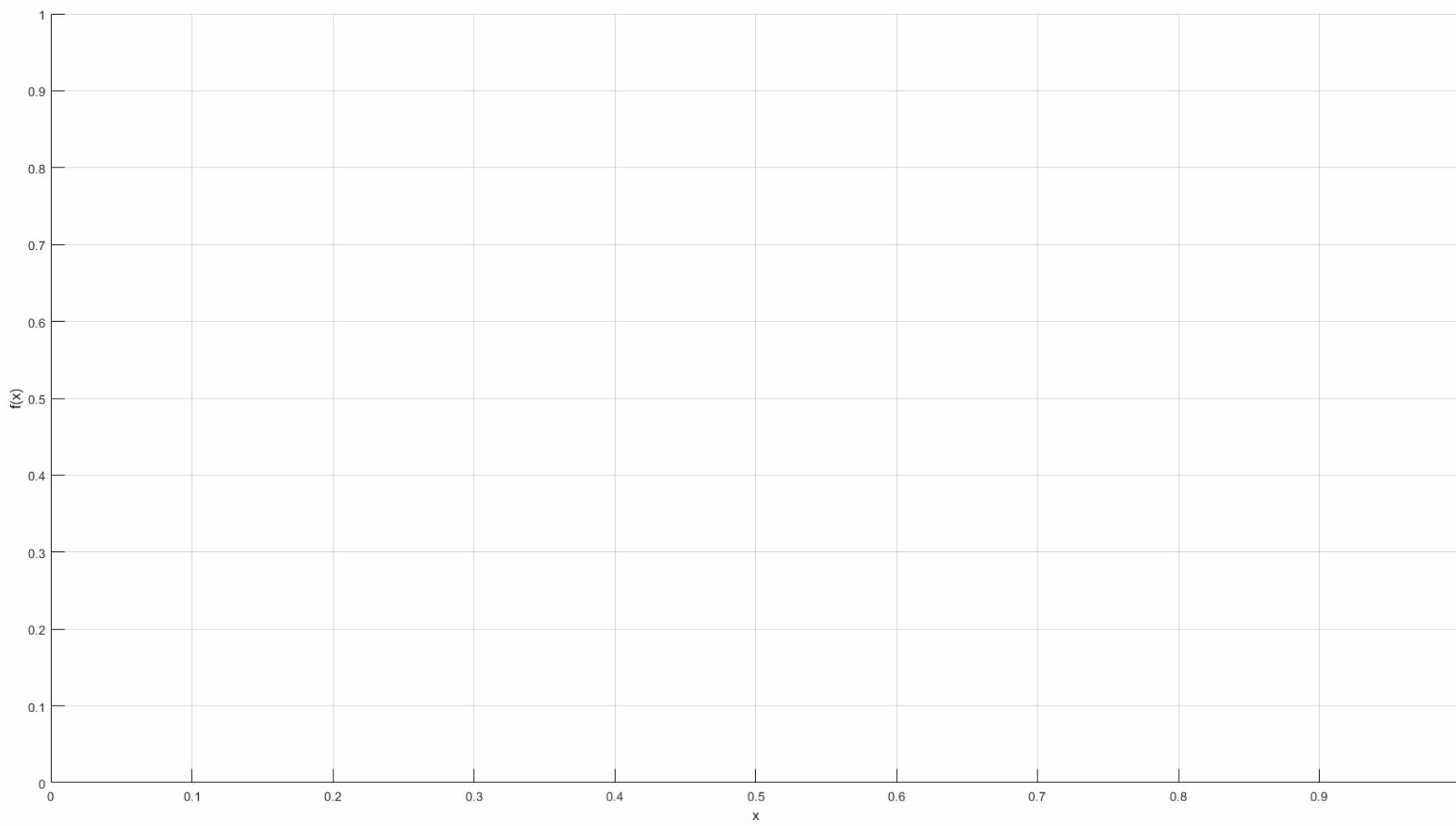
MLP with Softmax

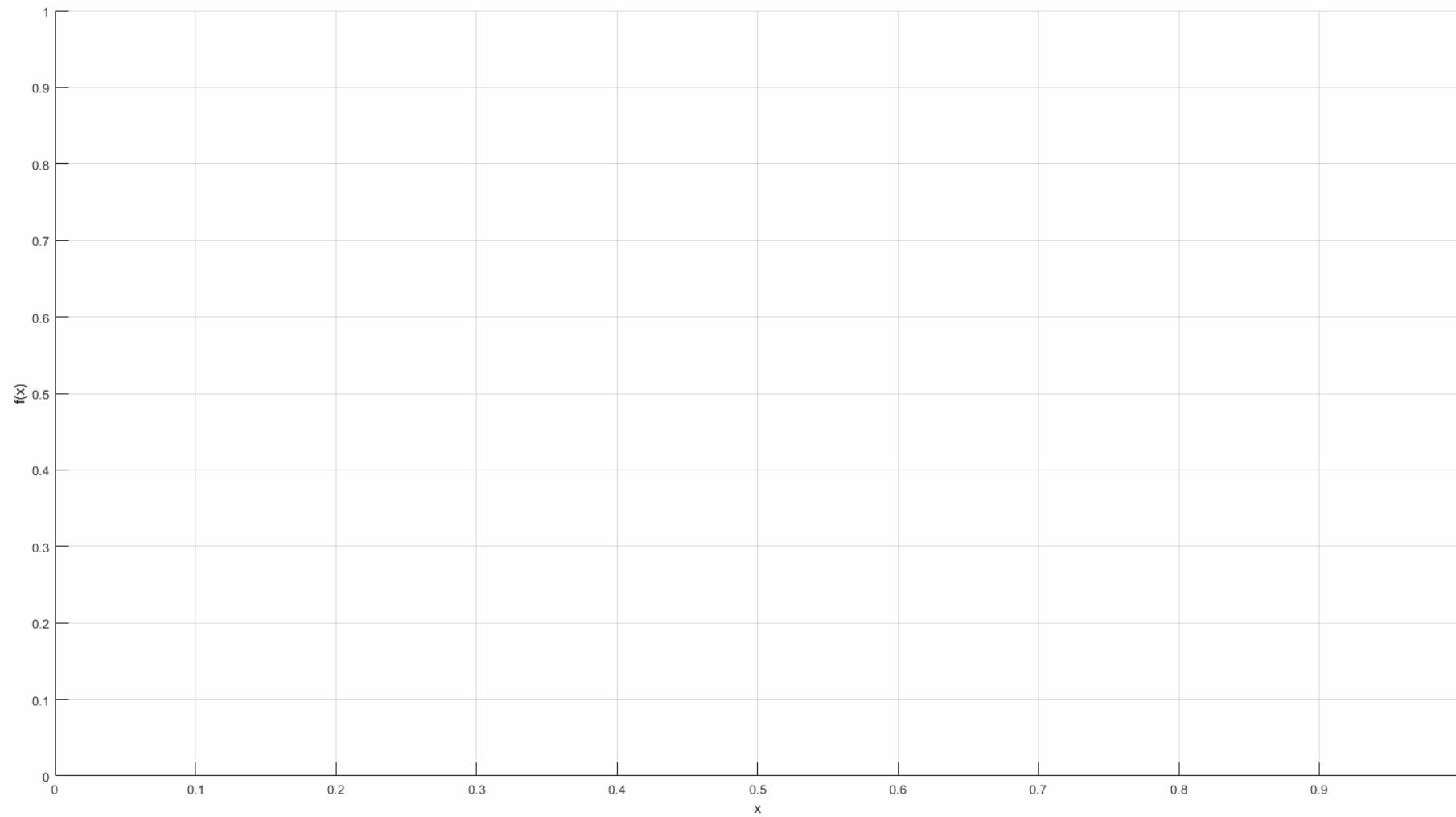


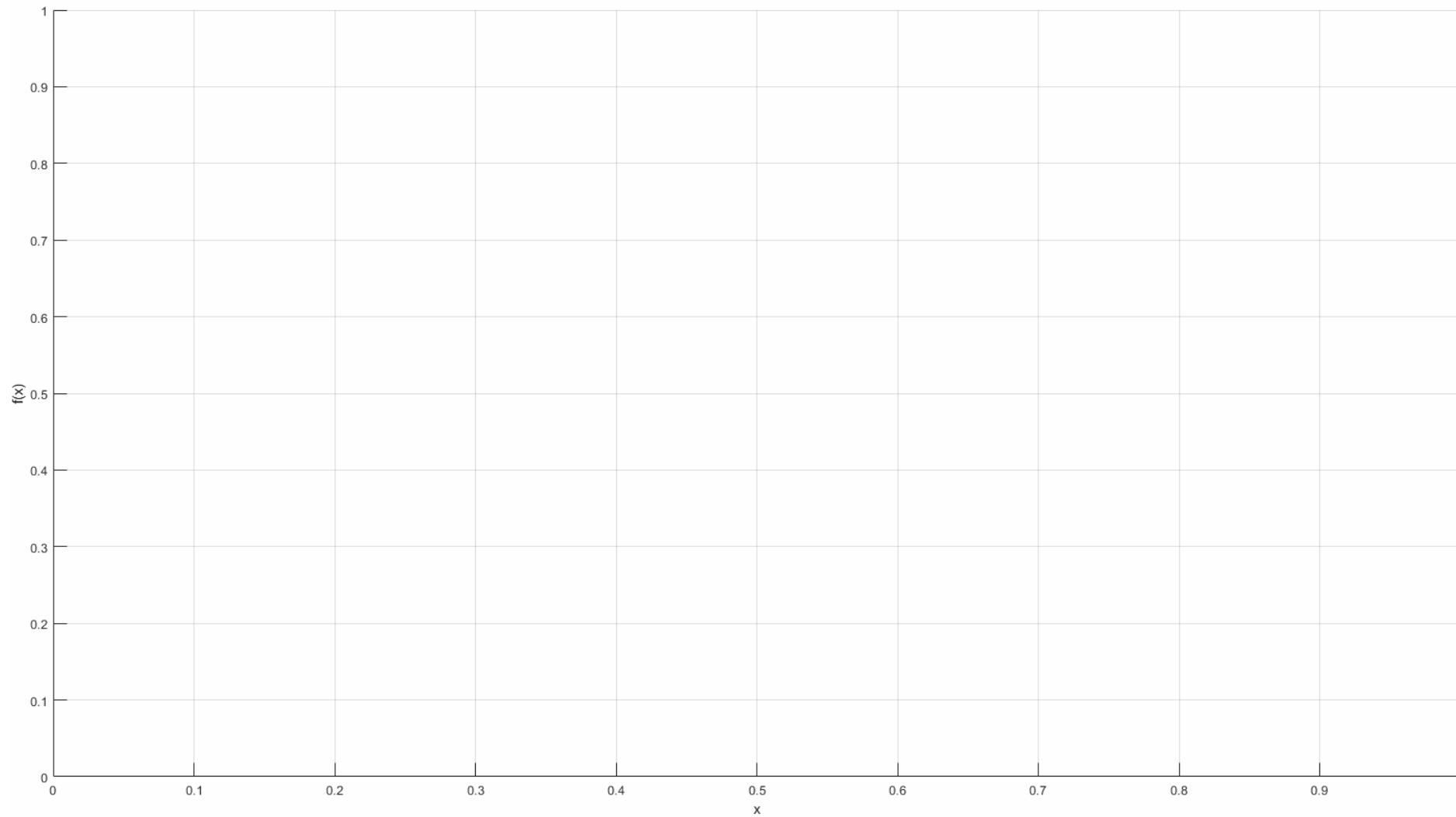
MLP with Softmax



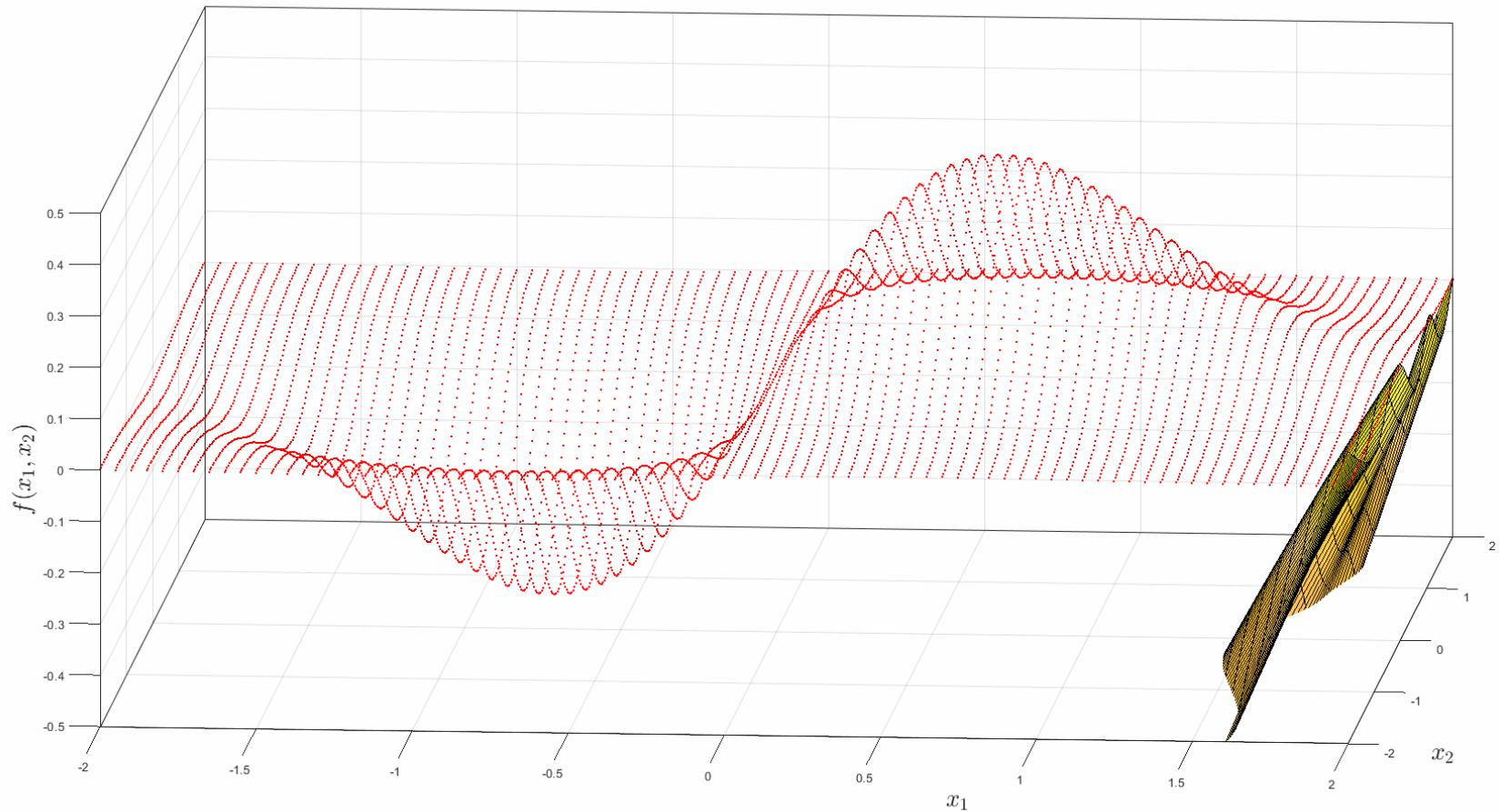
Dr. Carlos Villaseñor



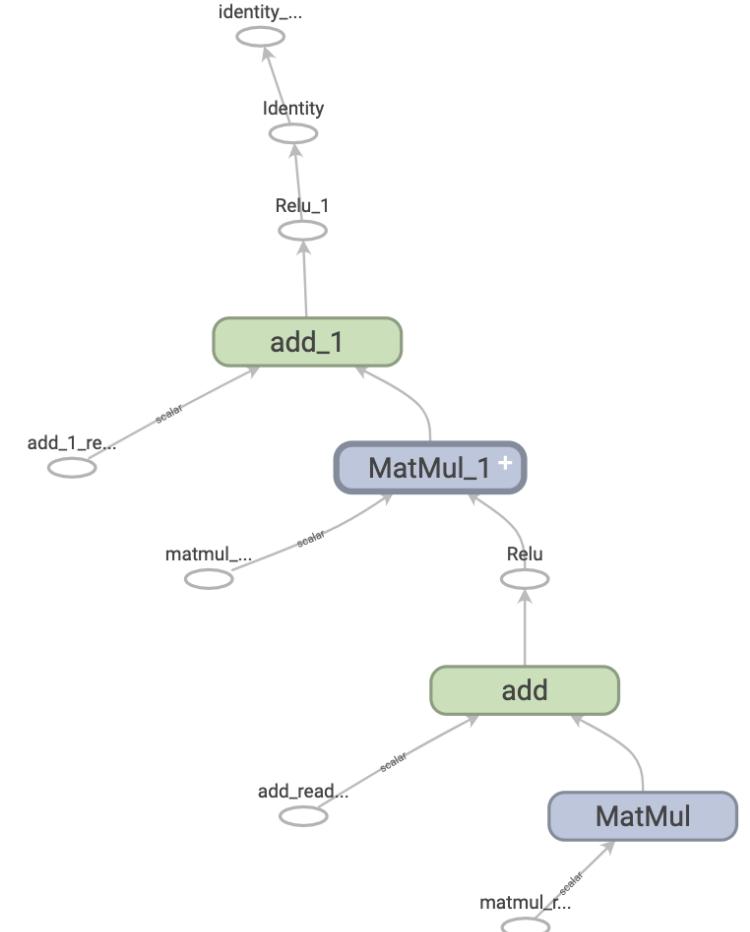
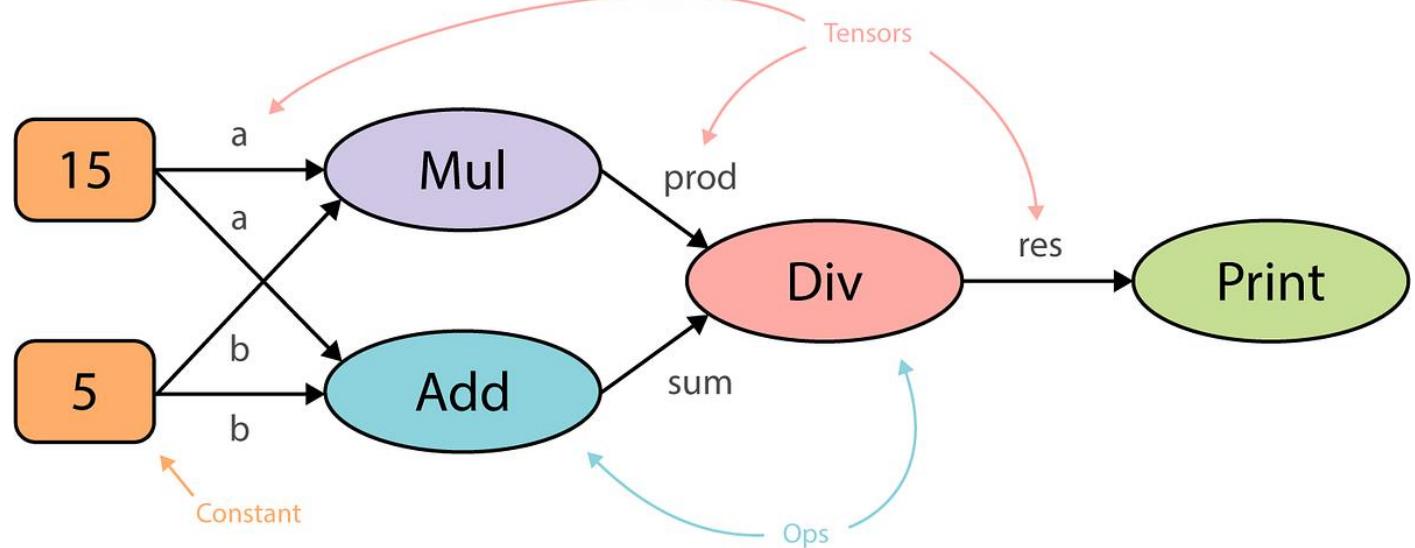




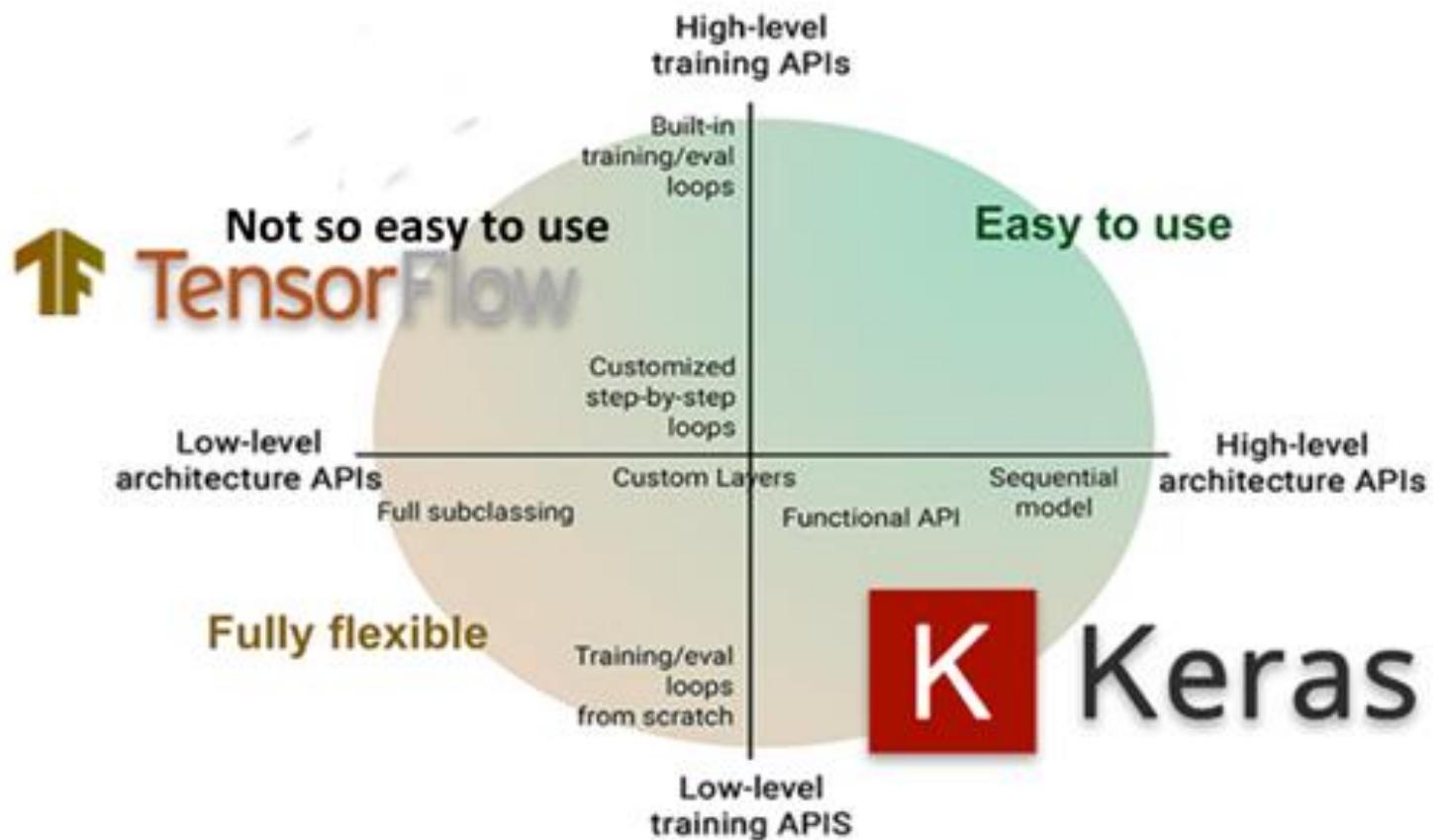
Backpropagation for Regression



Keras/Tensorflow



Keras/Tensorflow



Keras/Tensorflow

	Keras 	TensorFlow 	PyTorch 
Level of API	high-level API ¹	Both high & low level APIs	Lower-level API ²
Speed	Slow	High	High
Architecture	Simple, more readable and concise	Not very easy to use	Complex ³
Debugging	No need to debug	Difficult to debugging	Good debugging capabilities
Dataset Compatibility	Slow & Small	Fast speed & large	Fast speed & large datasets
Popularity Rank	1	2	3
Uniqueness	Multiple back-end support	Object Detection Functionality	Flexibility & Short Training Duration
Created By	Not a library on its own	Created by Google	Created by Facebook ⁴
Ease of use	User-friendly	Incomprehensive API	Integrated with Python language
Computational graphs used	Static graphs	Static graphs	Dynamic computation graphs ⁵