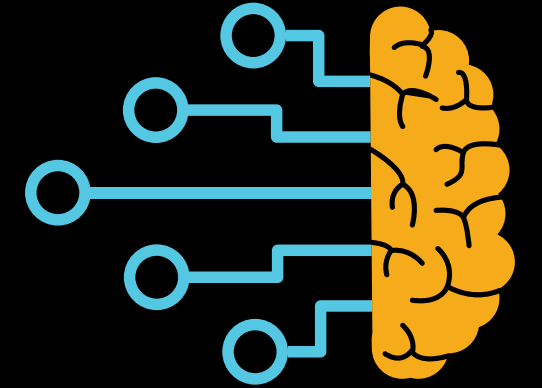


DCAO, UBA  
08-12 AGOSTO 2022



# APRENDIZAJE AUTOMÁTICO.

Fundamentos y Aplicaciones en  
Meteorología del Espacio

---

**Dra María Graciela Molina**  
FACET-UNT / CONICET  
Tucumán Space Weather Center - TSWC

<https://spaceweather.facet.unt.edu.ar/>  
IG -> @spaceweatherargentina



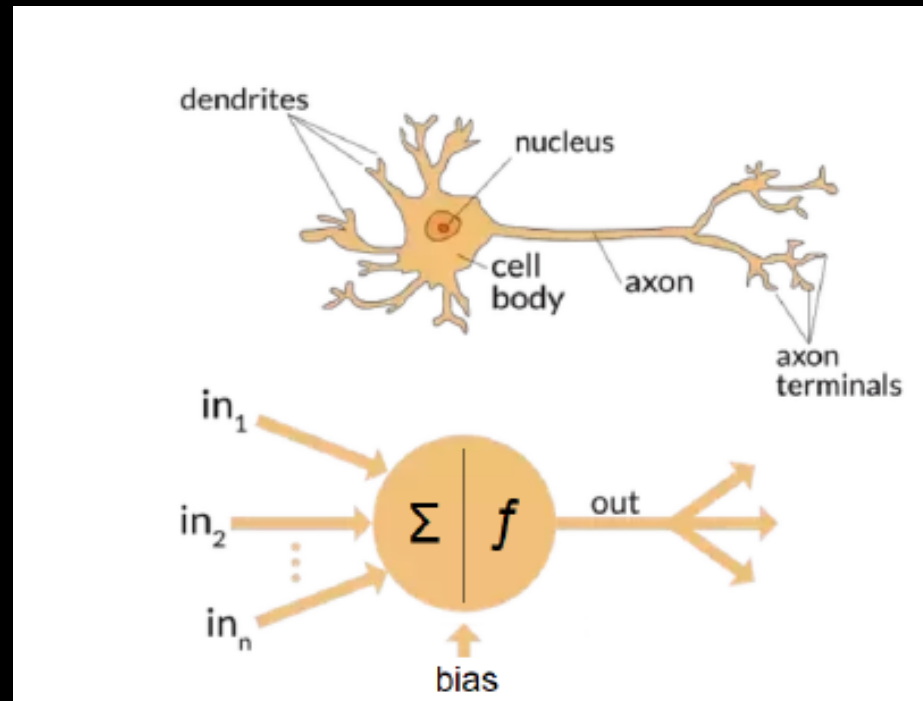
gmolina@herrera.unt.edu.ar

# Artificial Neural Networks

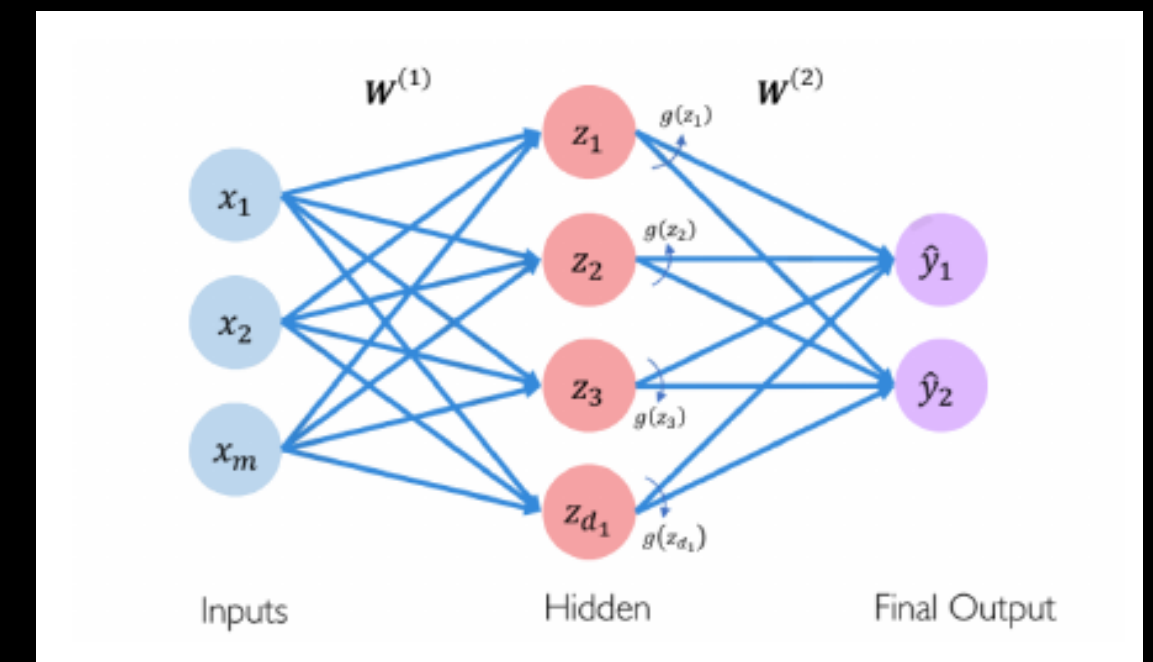
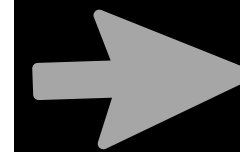
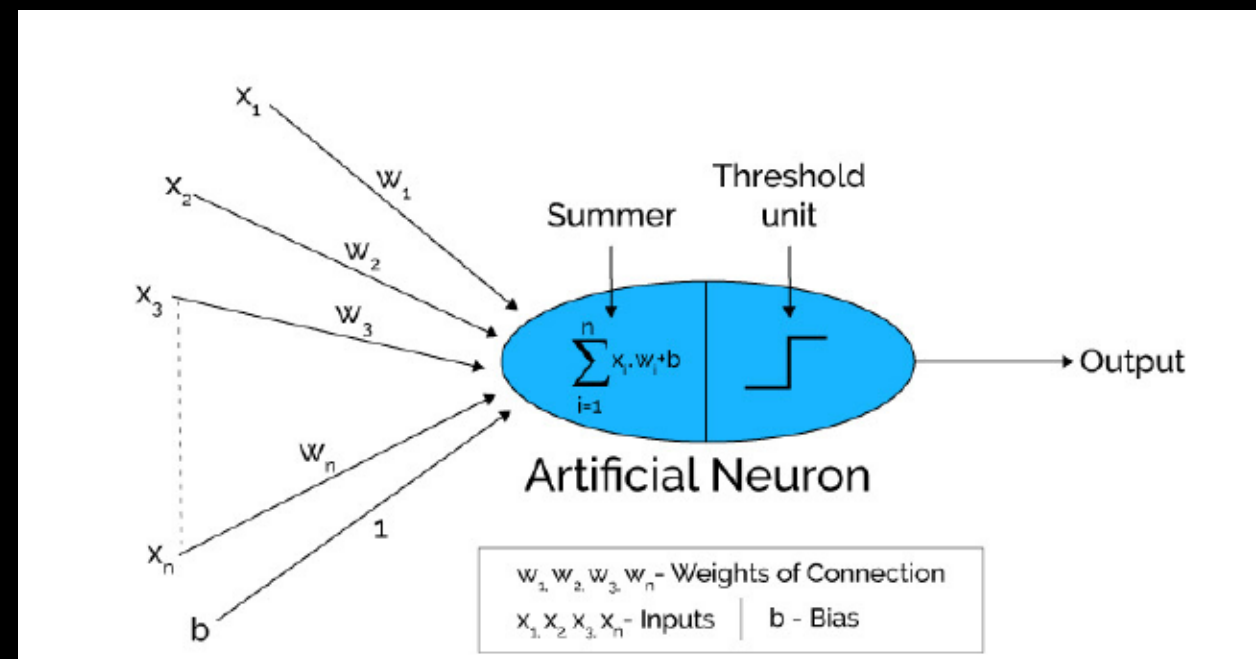


TSWC, 2022

## ANN



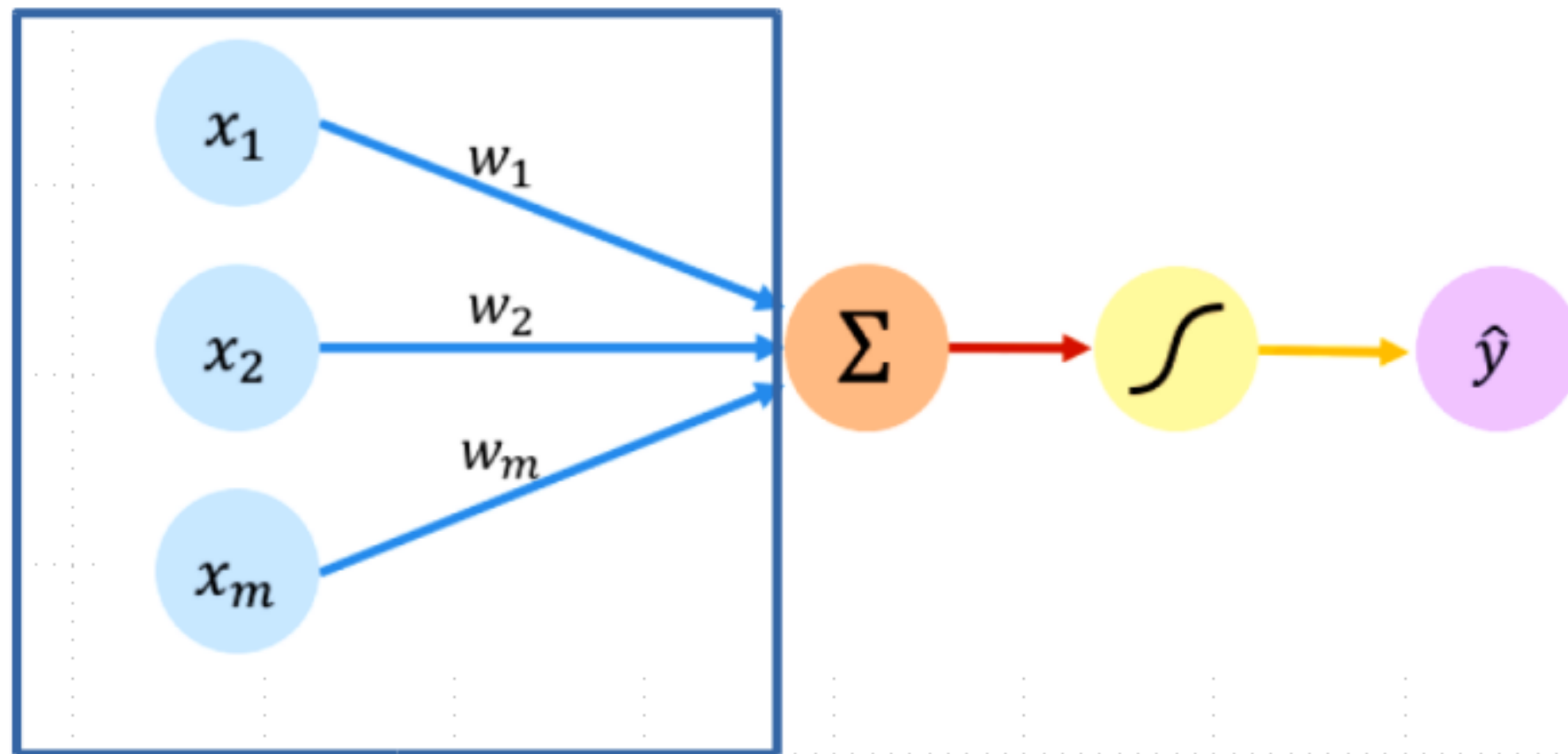
- Data-driven modeling
- Inspired in brain neural networks
- Solving wide number of complex problems: facial recognition, handwrite recognition, weather forecasting, etc.
- Black-box modelling (?) based in the composition of interconnected neurons (neural network)



# Perceptron

- Unidad fundamental para construir redes neuronales

Input  $\rightarrow$  Weights  $\rightarrow$  Sum  $\rightarrow$  Non linearity  $\rightarrow$  Output

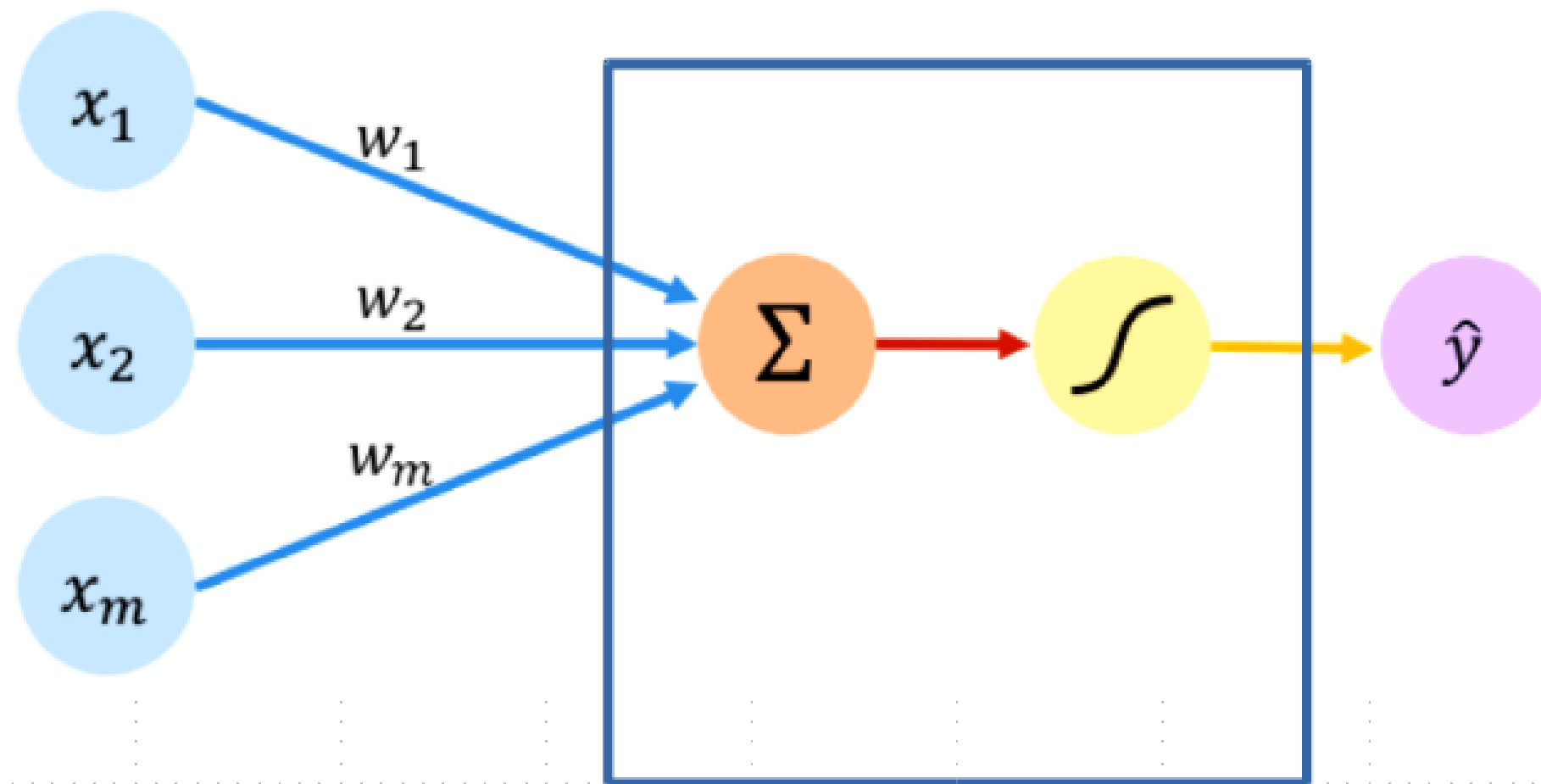


- Cada una de las entradas corresponde a una característica (feature)
- A cada una de las  $x_i$  se les aplica un peso ( $w$ ) (ponderación de los valores de la entrada)
- Los pesos son las vbles que se ajustan durante el entrenamiento



# Perceptron

Input → Weights → Sum → Non linearity → Output



Output

Linear combination of inputs

$$\hat{y} = g \left( w_0 + \sum_{i=1}^m x_i w_i \right)$$

Non-linear activation function

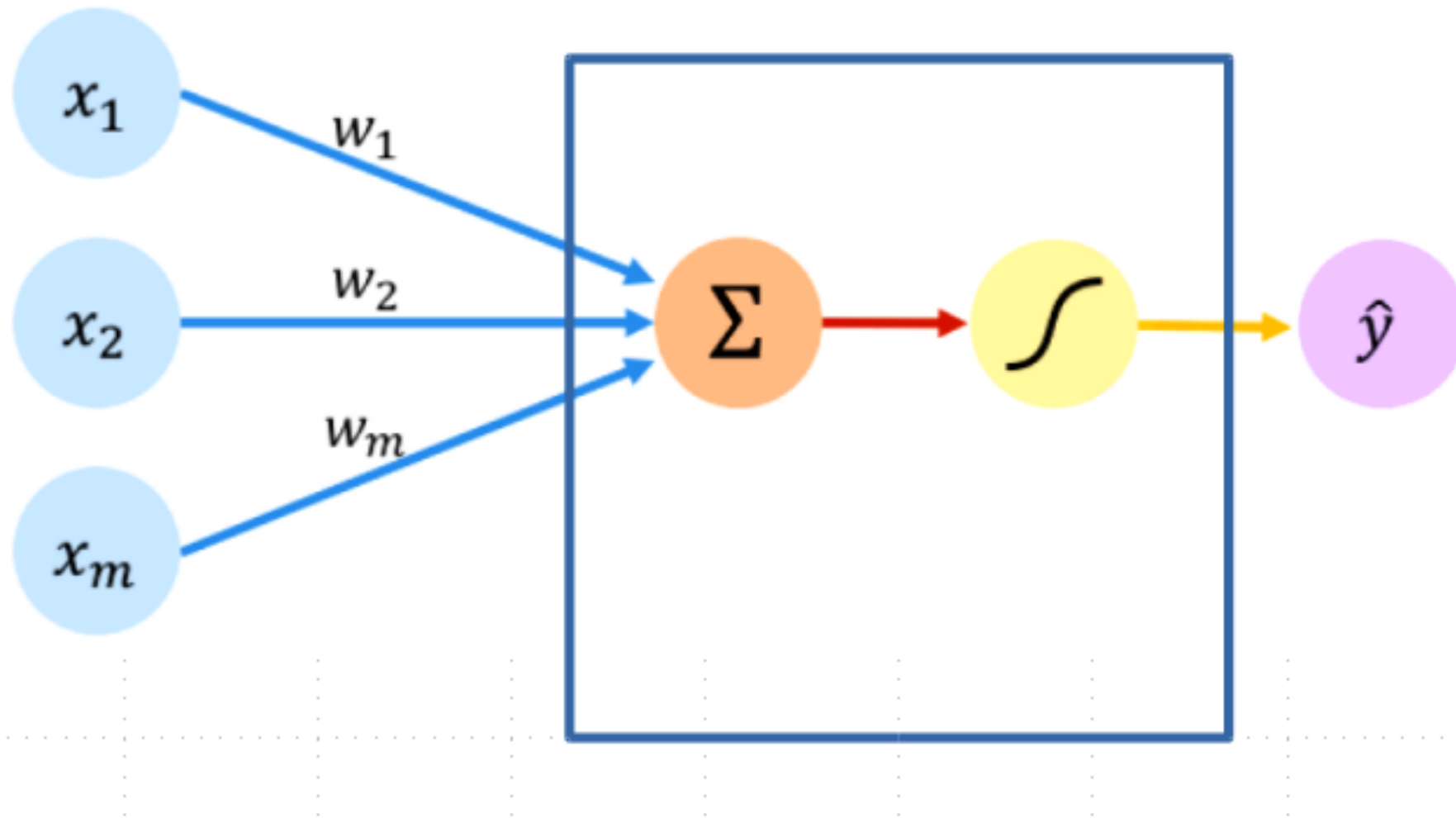
Bias

- Capa oculta (hidden layer).
- El bias es un término de ajuste
- La función de activación se aplica a la combinación lineal de las entradas ponderadas teniendo en cuenta el bias



# Perceptron

Input  $\rightarrow$  Weights  $\rightarrow$  Sum  $\rightarrow$  Non linearity  $\rightarrow$  Output



$$\hat{y} = g \left( w_0 + \sum_{i=1}^m x_i w_i \right)$$

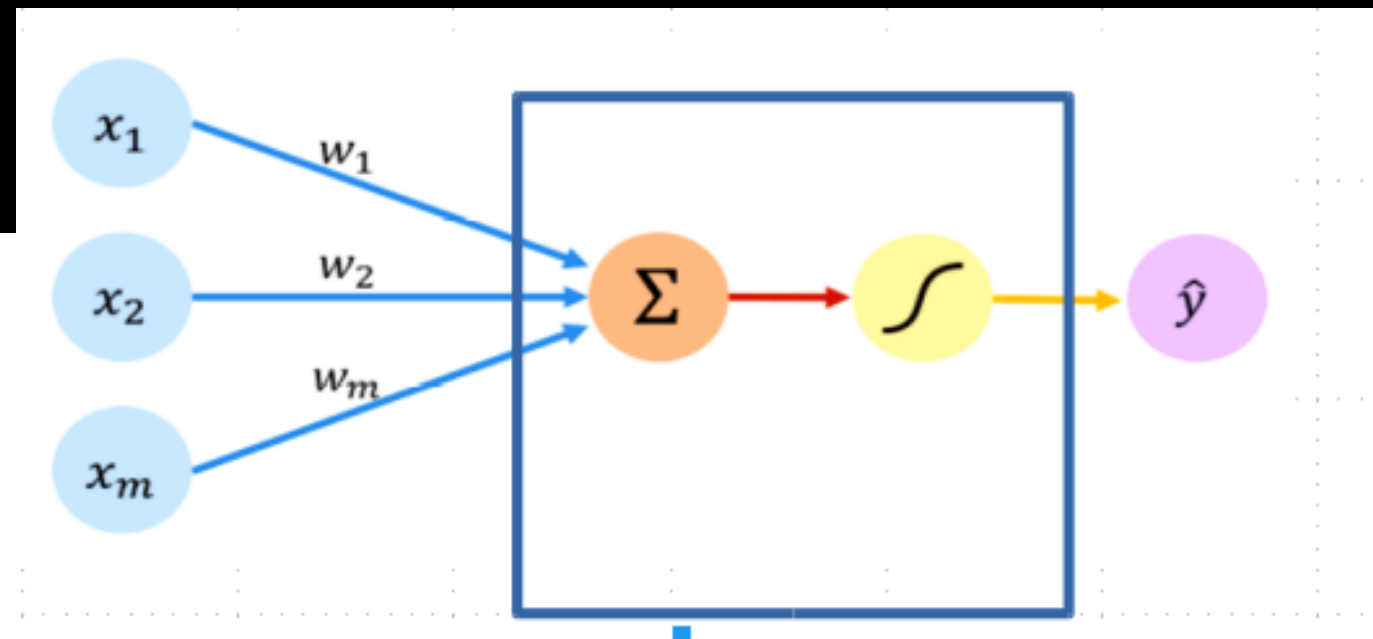
$$\hat{y} = g ( w_0 + \mathbf{X}^T \mathbf{W} )$$

$$\text{where: } \mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } \mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$



# Activation function

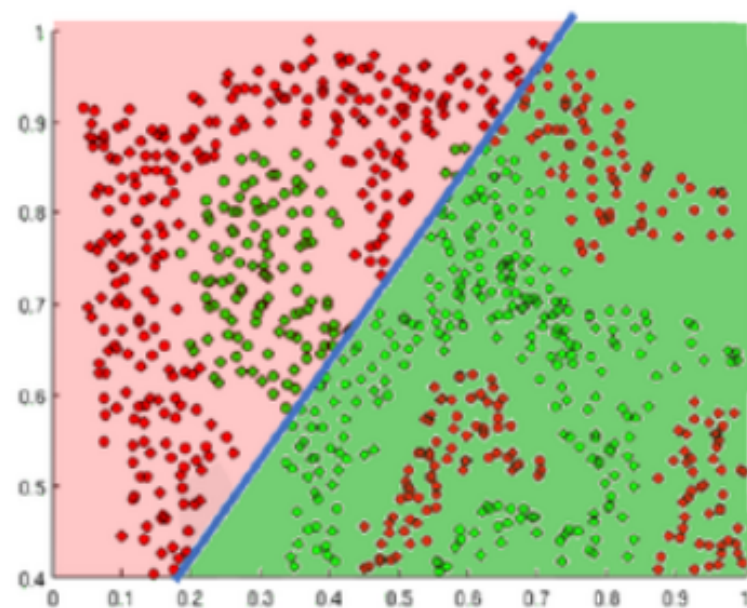
- Agrega no-linealidad al modelo



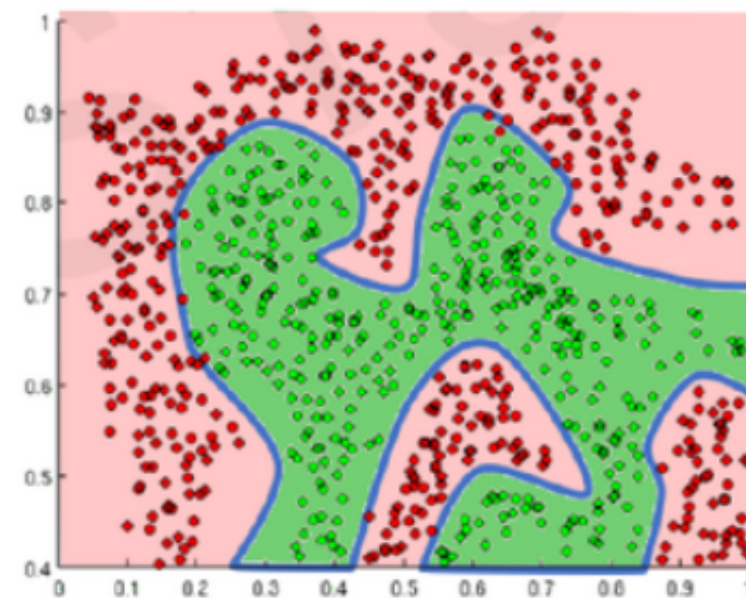
$$\hat{y} = g(w_0 + X^T W)$$

- Una de las más usadas en la función sigmoide.
- Produce resultados entre 0 y 1

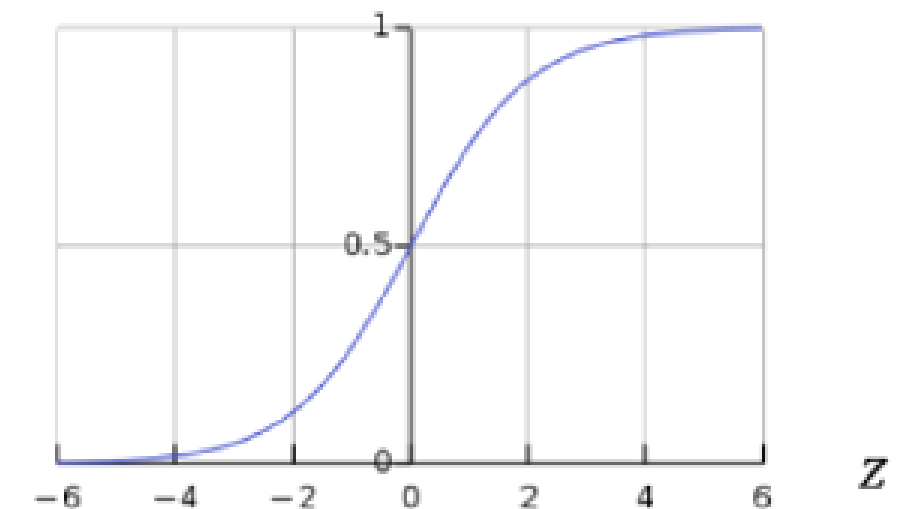
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions



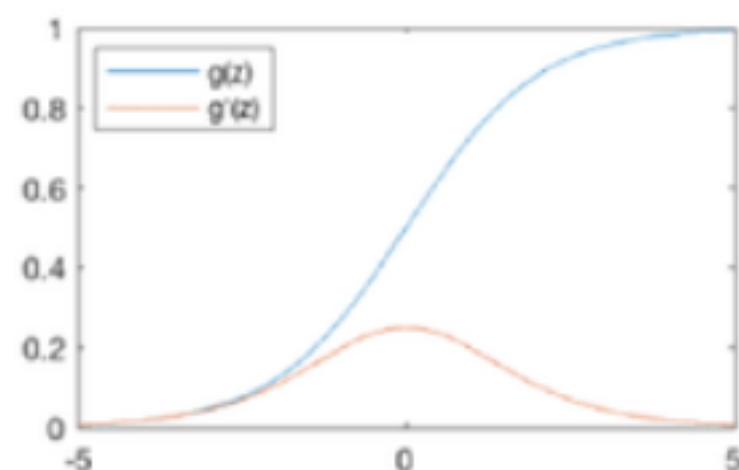




# Activation function

$$\hat{y} = g(w_0 + X^T W)$$

Sigmoid Function



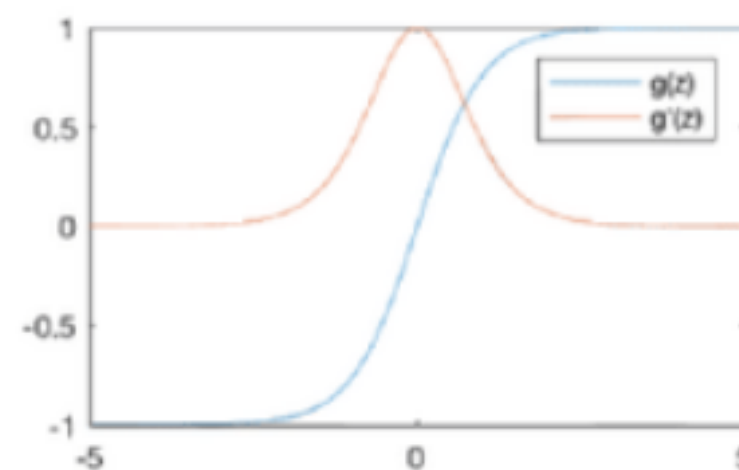
$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$



```
tf.math.sigmoid(z)
```

Hyperbolic Tangent



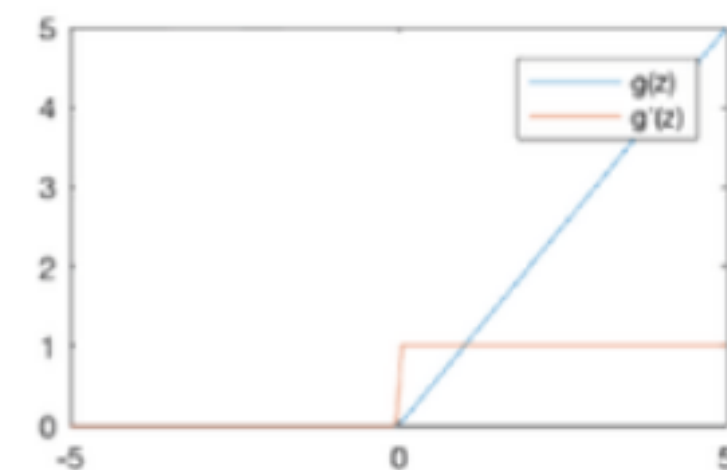
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$



```
tf.math.tanh(z)
```

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

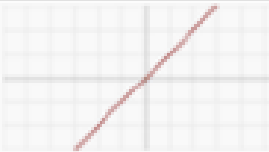


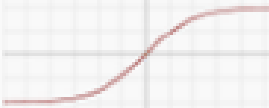
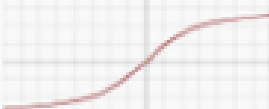
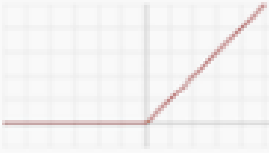
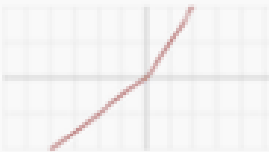
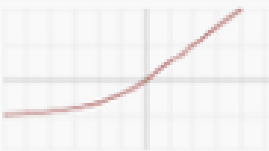
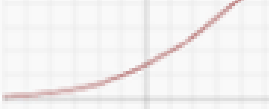


```
tf.nn.relu(z)
```

# Activation function

- Adds no-linearity to the model

$$\hat{y} = g(w_0 + X^T W)$$

Name	Plot	Equation	Derivative
Identity		$f(x) = x$	$f'(x) = 1$
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$
Logistic (a.k.a Soft step)		$f(x) = \frac{1}{1 + e^{-x}}$	$f'(x) = f(x)(1 - f(x))$
Tanh		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
ArcTan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Parametric Rectified Linear Unit (PReLU) [2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
Exponential Linear Unit (ELU) [3]		$f(x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} f(x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$
SoftPlus		$f(x) = \log_e(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$

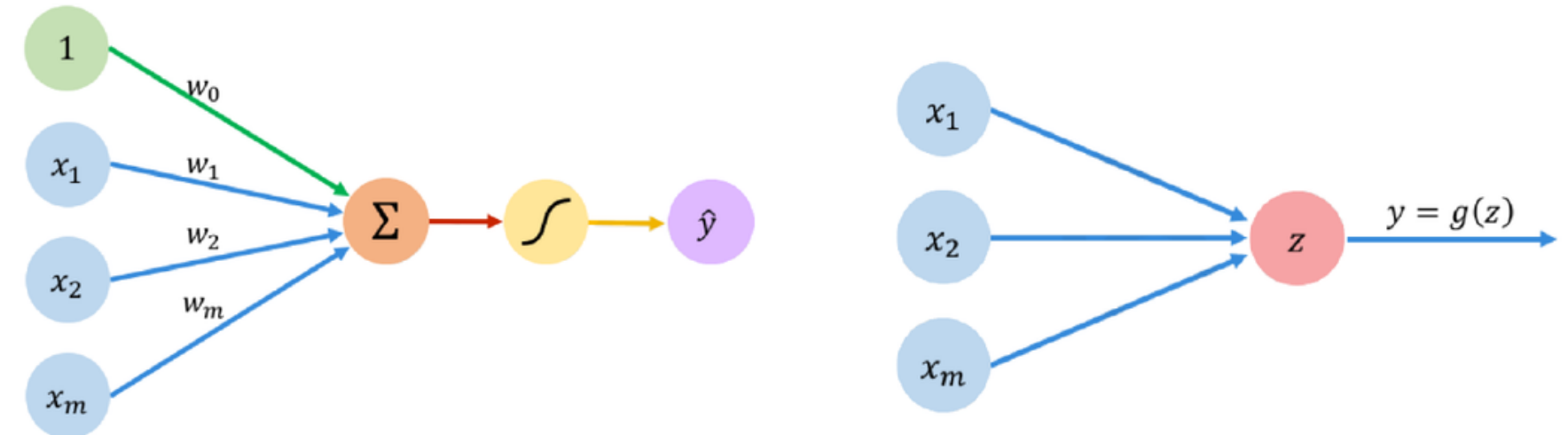




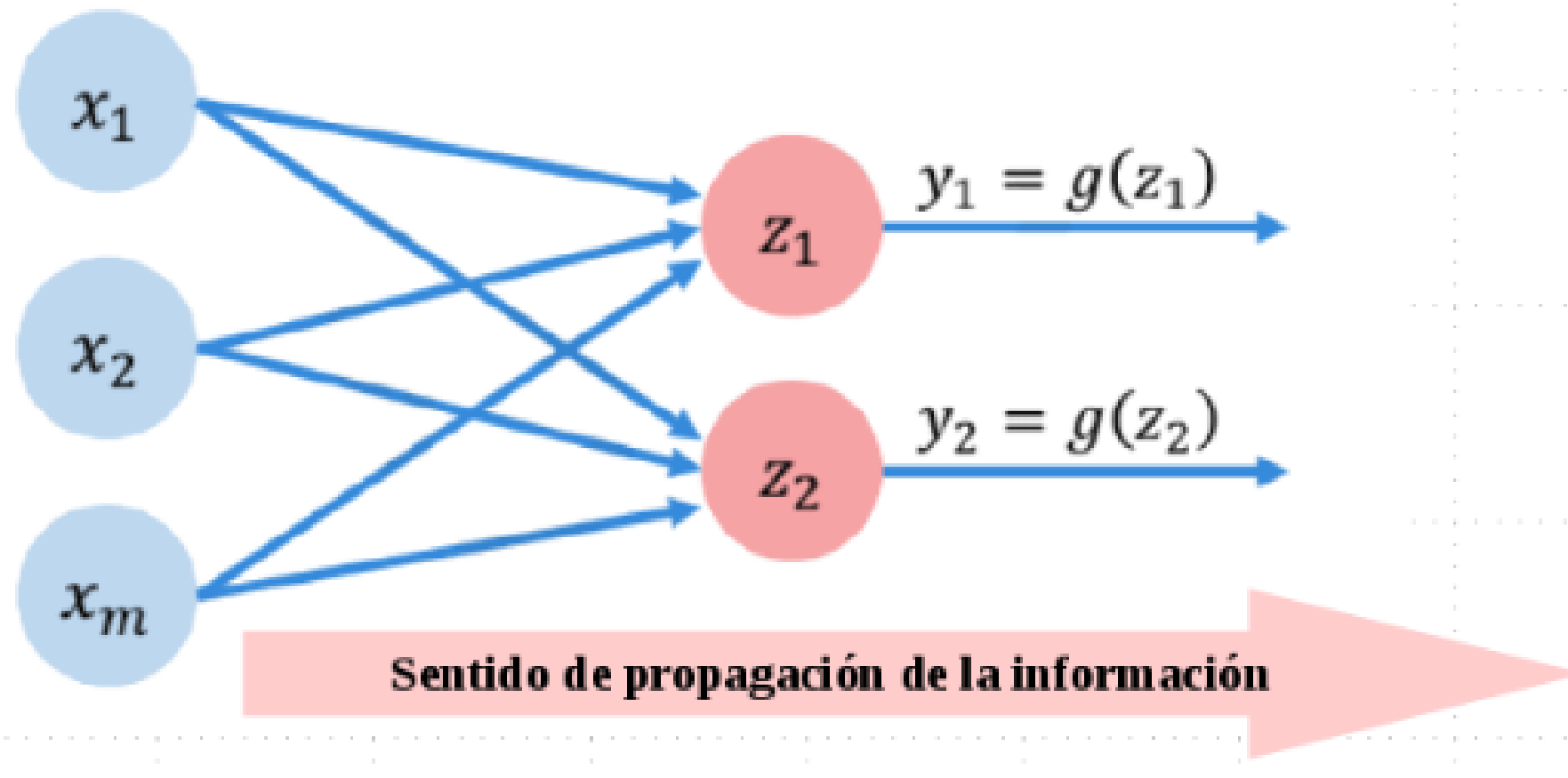
# Feed Forward (FF)

- Perceptrón simplificado (n entradas → 1 salida)

$$z = w_0 + \sum_{j=1}^m x_j w_j$$



- Armemos una red neuronal compuesta por múltiples perceptrones (n entradas → m salida)

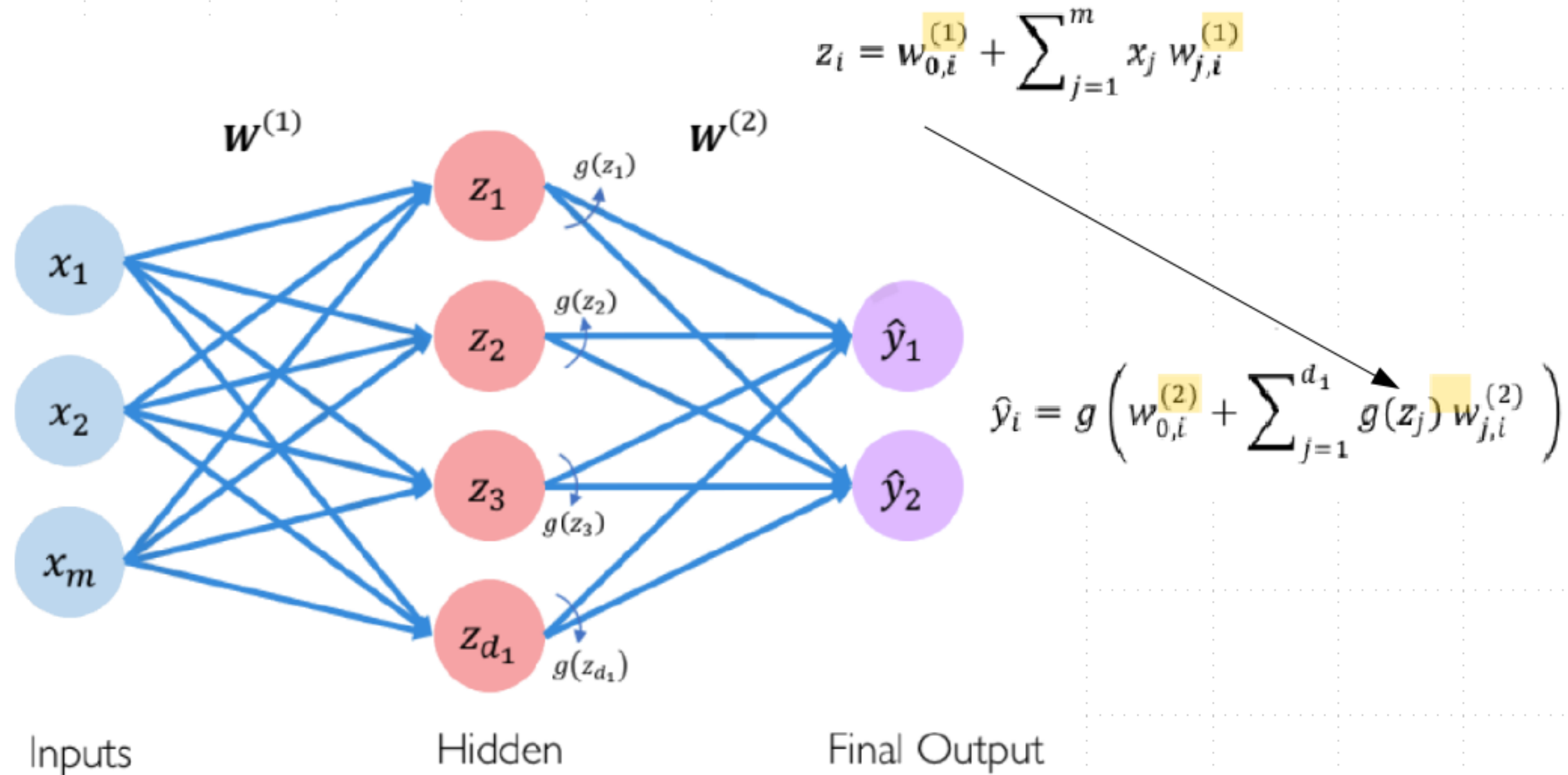


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$

Observación: todas las entradas están conectadas a todas las salidas → dense layers

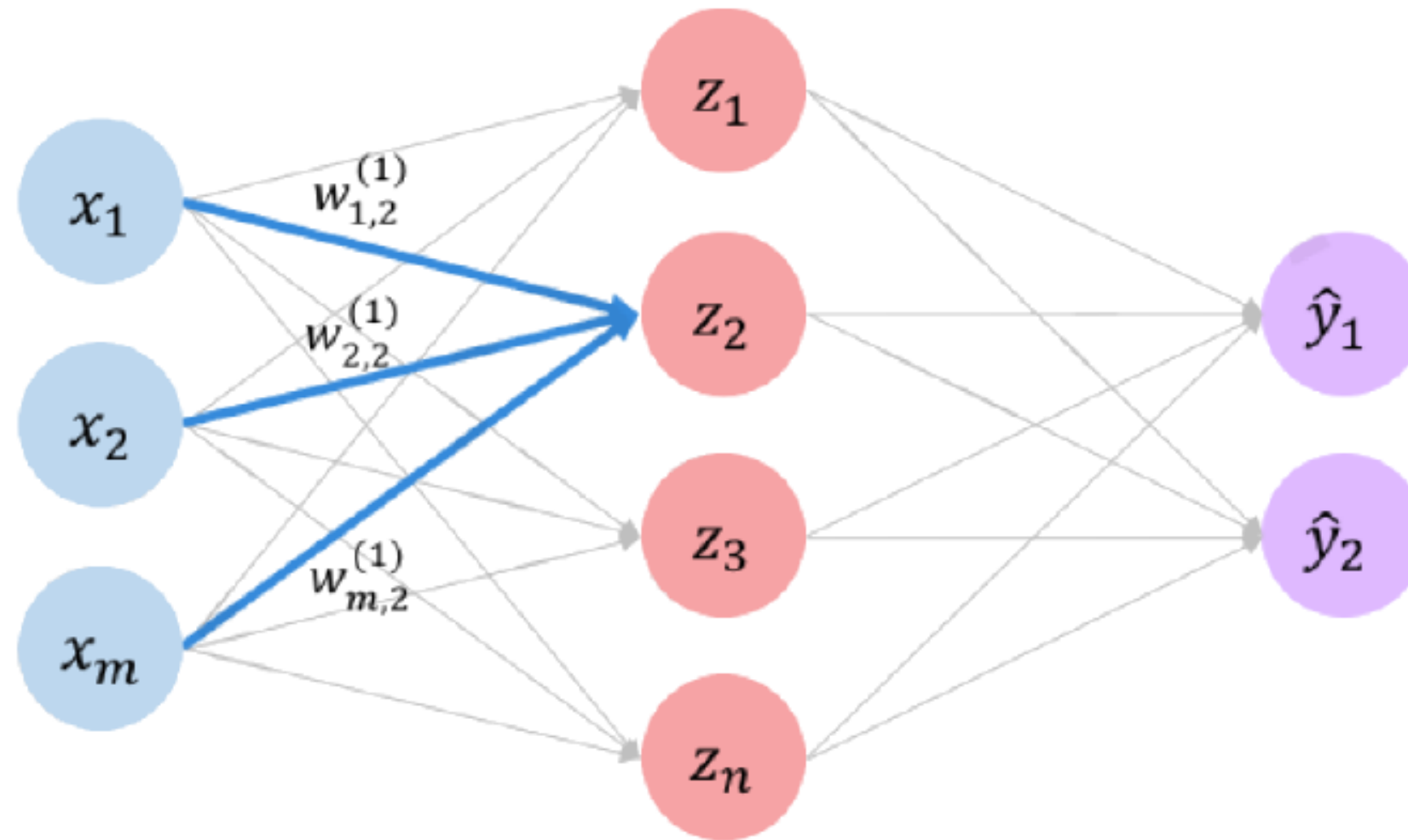


# 1-Layer Neural Network





# 1-Layer Neural Network

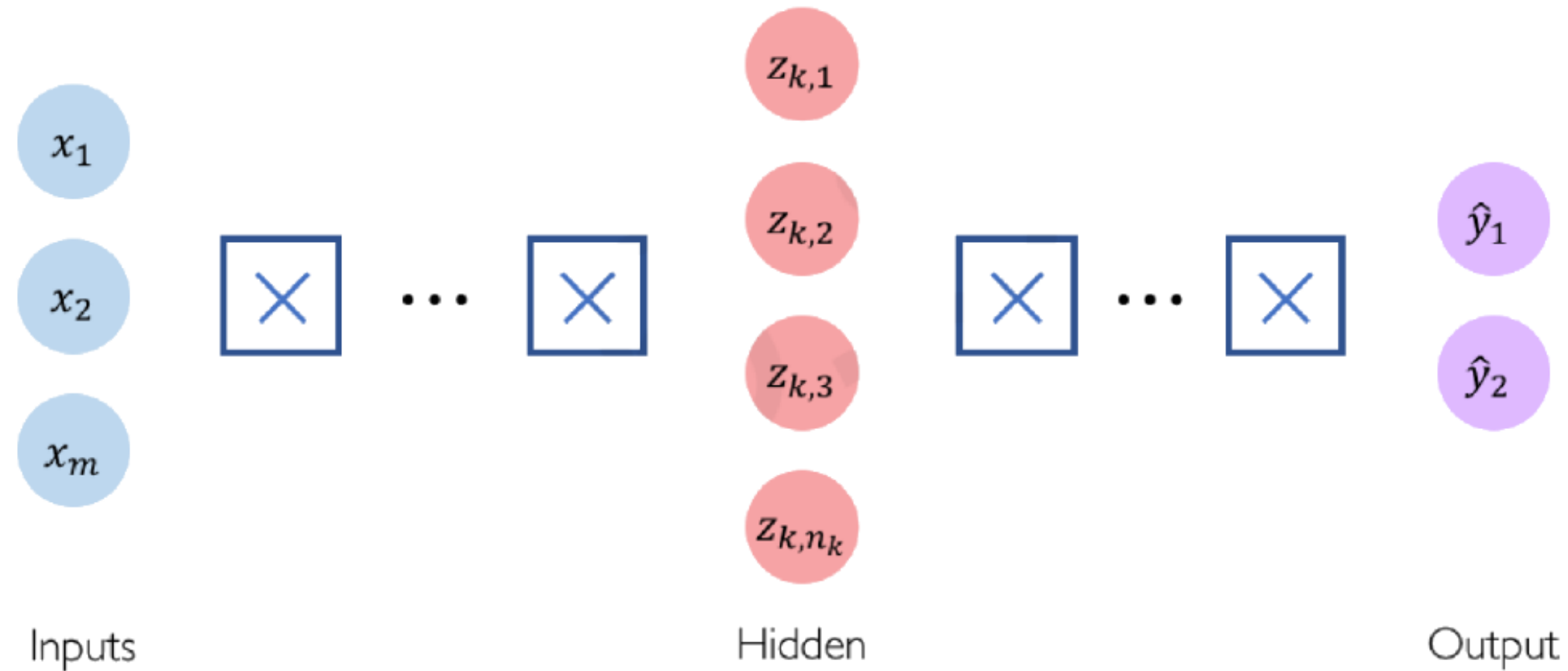


$$z_2 = w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)}$$

$$= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)}$$

# M-layers Neural Network (multi-layer)

- Deep Neural Network

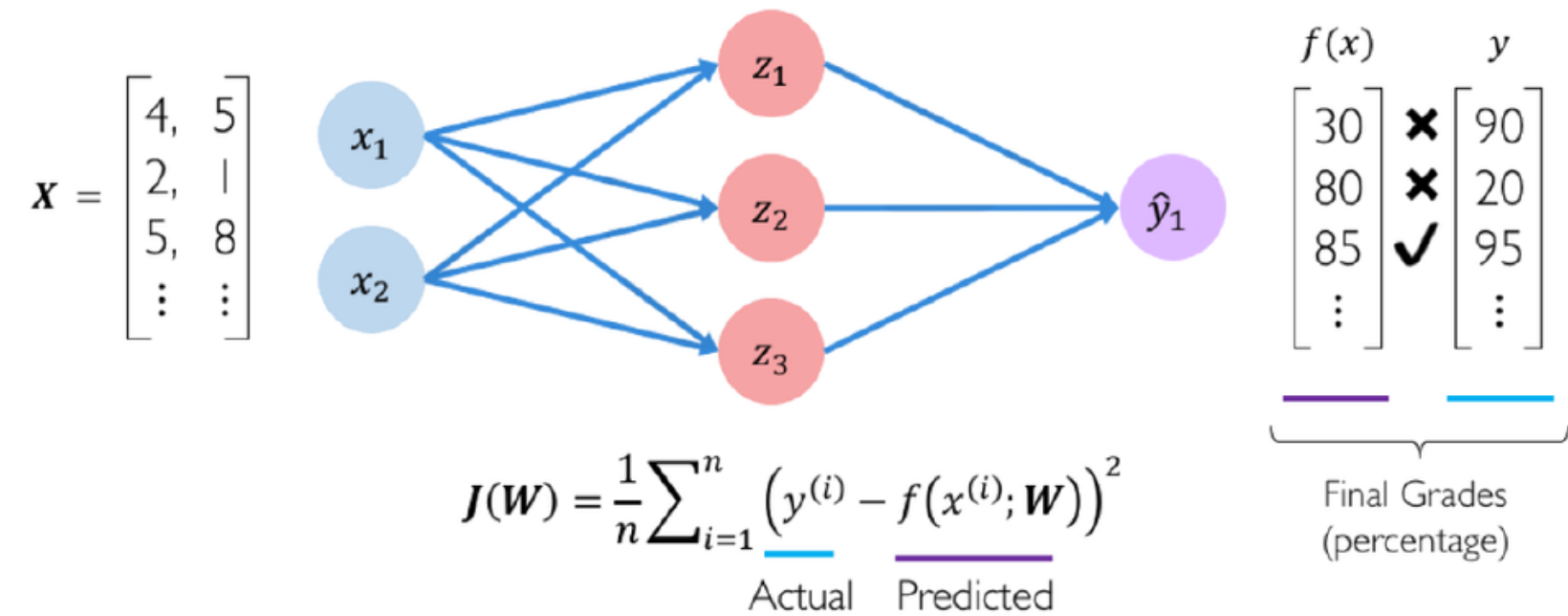


$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

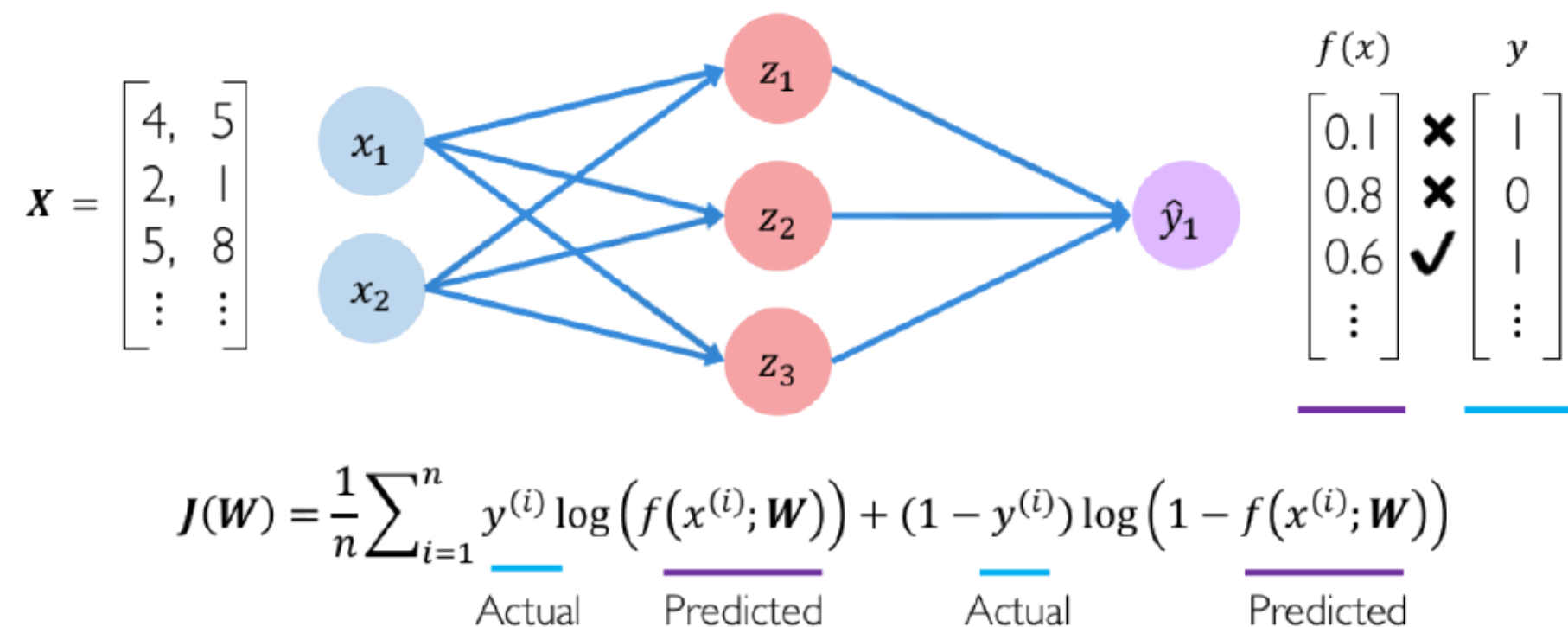


# Loss Function

Mean squared error loss can be used with regression models that output continuous real numbers

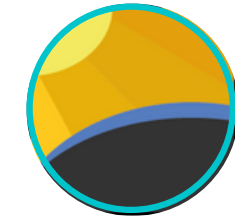


Cross entropy loss can be used with models that output a probability between 0 and 1



symbol	name	equation
$\mathcal{L}_1$	L <sub>1</sub> loss	$\ y - o\ _1$
$\mathcal{L}_2$	L <sub>2</sub> loss	$\ y - o\ _2^2$
$\mathcal{L}_1 \circ \sigma$	expectation loss	$\ y - \sigma(o)\ _1$
$\mathcal{L}_2 \circ \sigma$	regularised expectation loss <sup>1</sup>	$\ y - \sigma(o)\ _2^2$
$\mathcal{L}_\infty \circ \sigma$	Chebyshev loss	$\max_j  \sigma(o)^{(j)} - y^{(j)} $
hinge	hinge [13] (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{y}^{(j)} o^{(j)})$
hinge <sup>2</sup>	squared hinge (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{y}^{(j)} o^{(j)})^2$
hinge <sup>3</sup>	cubed hinge (margin) loss	$\sum_j \max(0, \frac{1}{2} - \hat{y}^{(j)} o^{(j)})^3$
log	log (cross entropy) loss	$-\sum_j y^{(j)} \log \sigma(o)^{(j)}$
log <sup>2</sup>	squared log loss	$-\sum_j [y^{(j)} \log \sigma(o)^{(j)}]^2$
tan	Tanimoto loss	$\frac{-\sum_j \sigma(o)^{(j)} y^{(j)}}{\ \sigma(o)\ _2^2 + \ y\ _2^2 - \sum_j \sigma(o)^{(j)} y^{(j)}}$
D <sub>CS</sub>	Cauchy-Schwarz Divergence [3]	$-\log \frac{\sum_j \sigma(o)^{(j)} y^{(j)}}{\ \sigma(o)\ _2 \ y\ _2}$

• <https://arxiv.org/pdf/1702.05659.pdf>s



# Training (loss optimization)

- Encontrar los valores de los pesos de la red neuronal tal que la pérdida o costo sea mínimo

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

- Es un problema de optimización
- Podemos usar algún método numérico de optimización

MÉTODO DEL GRADIENTE DESCENDIENTE





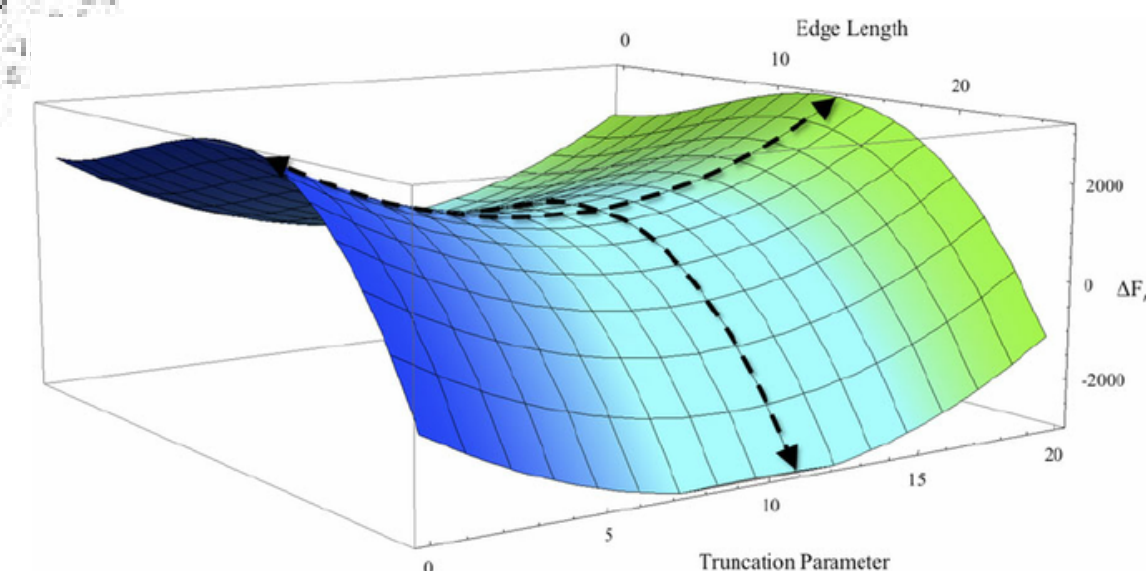
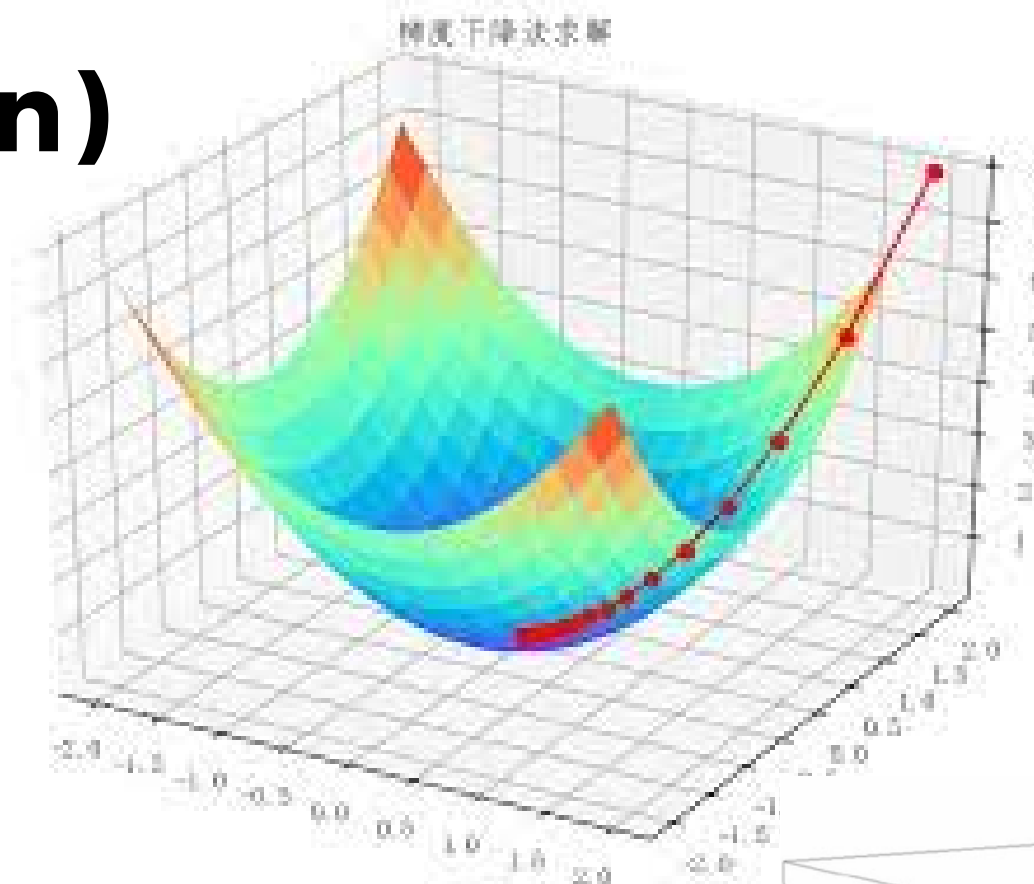
# Training (loss optimization)

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$

## Algorithm

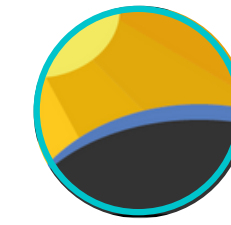
1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(W)}{\partial W}$
4. Update weights,  $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights



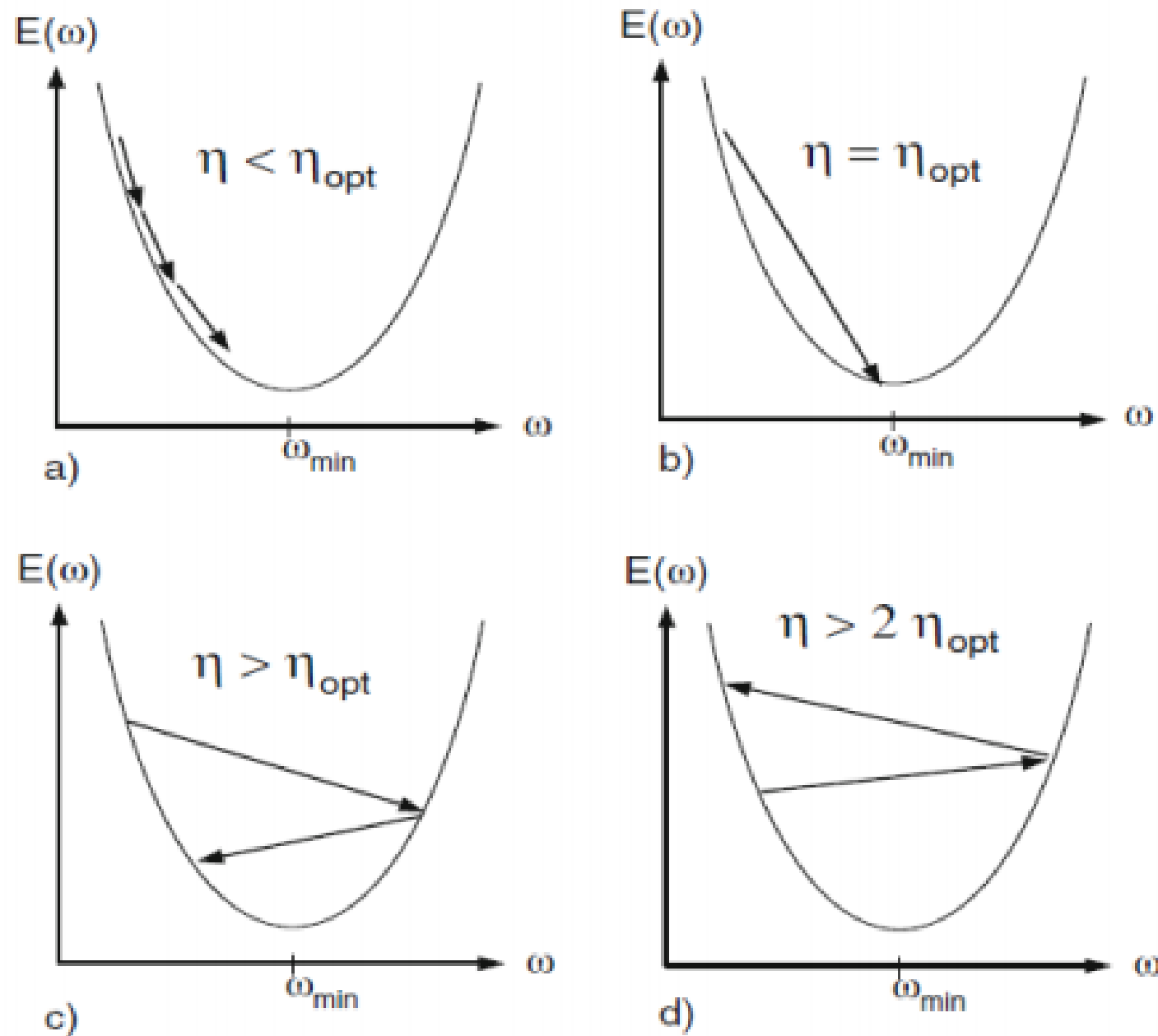
## Principales problemas:

- Mínimos locales
- Silla de montar: el g. d. un vez que llega a una región con gradiente cero, no puede escapar de allí independientemente de la calidad del mínimo.





# Training (loss optimization)








$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

How can we set the learning rate?

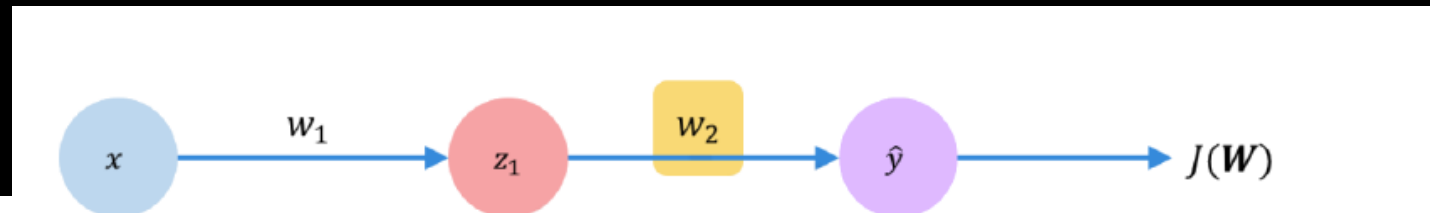
- <https://runder.io/optimizing-gradient-descent/index.html>

## Gradient Descent Algorithms

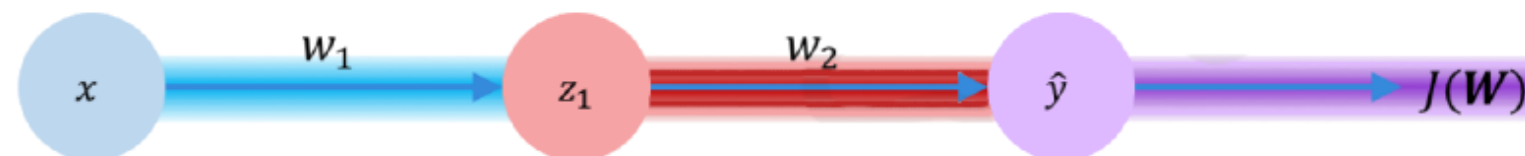
Algorithm	TF Implementation	Reference
• SGD	 <code>tf.keras.optimizers.SGD</code>	Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952.
• Adam	 <code>tf.keras.optimizers.Adam</code>	Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.
• Adadelta	 <code>tf.keras.optimizers.Adadelta</code>	Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.
• Adagrad	 <code>tf.keras.optimizers.Adagrad</code>	Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.
• RMSProp	 <code>tf.keras.optimizers.RMSProp</code>	

# Backpropagation

- Qué tanto puede afectar a la función de pérdida un pequeño cambio en uno de los pesos?



- Aplicando la regla de la cadena



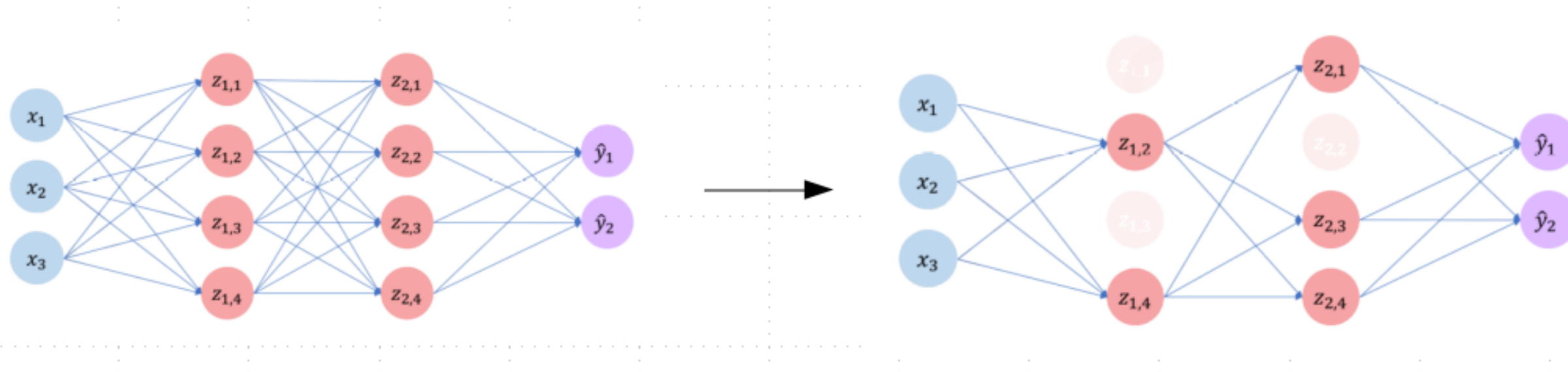
$$\frac{\partial J(W)}{\partial w_1} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial z_1}}_{\text{red}} * \underbrace{\frac{\partial z_1}{\partial w_1}}_{\text{blue}}$$

- Repetir para cada uno de los pesos en la red usando los gradientes de las capas siguiente!

- Info se propaga desde las entradas hacia adelante mediante sus parámetros hasta que logra hacer una predicción
- Luego realiza una propagación hacia atrás a lo largo de la red para ir modificando los parámetros de manera que el error final sea el mínimo.
- El error asociado a una mala predicción es usado para ajustar los parámetros (el aprendizaje puede ser visto como una optimización). Para esto la salida y es comparada con el valor esperado (target) del conjunto de entrenamiento (z). La diferencia se entre el valor esperado y el de salida se llama error o residuo.

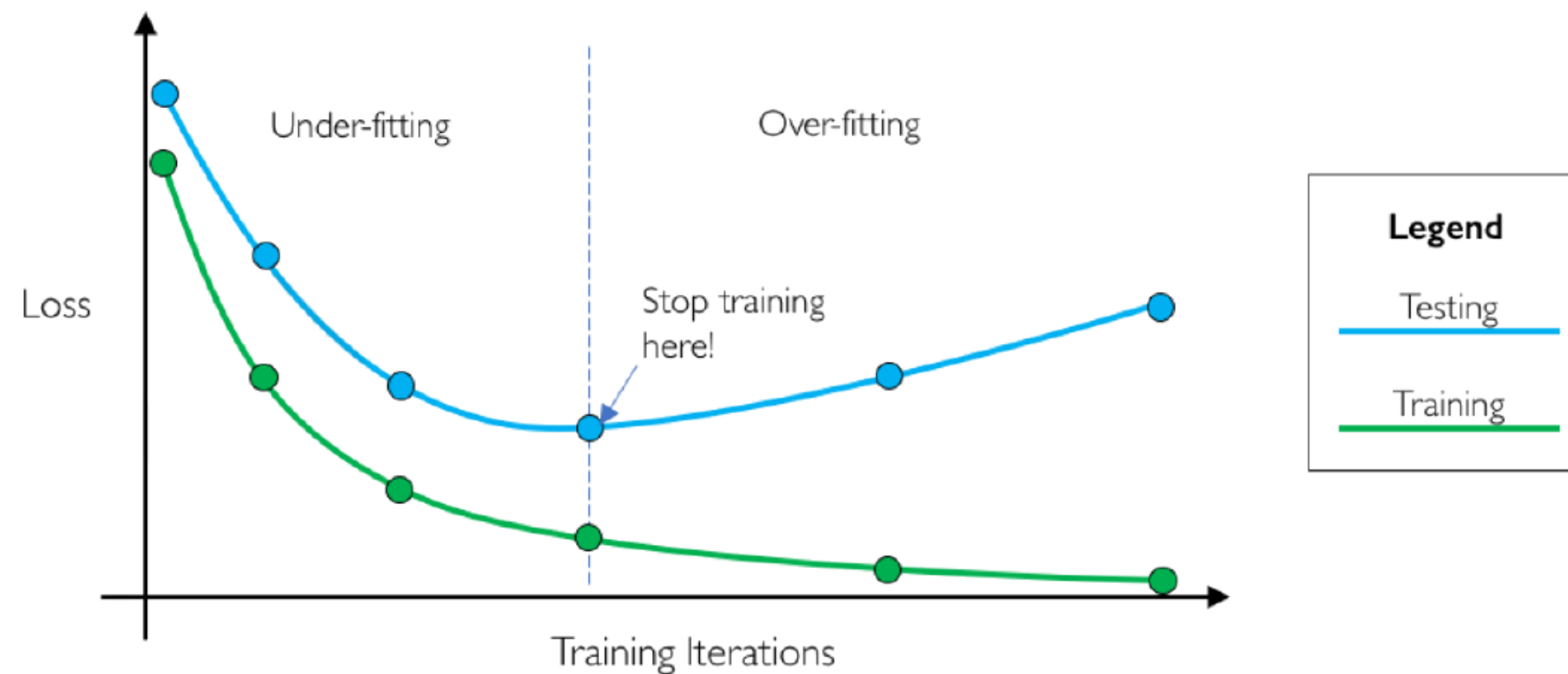
# Regularización

- Técnica que consiste en establecer restricciones al problema de optimización para evitar llegar a modelos complejos
- La idea es poder generalizar nuestro modelo para nuevos datos



## DROP OUT

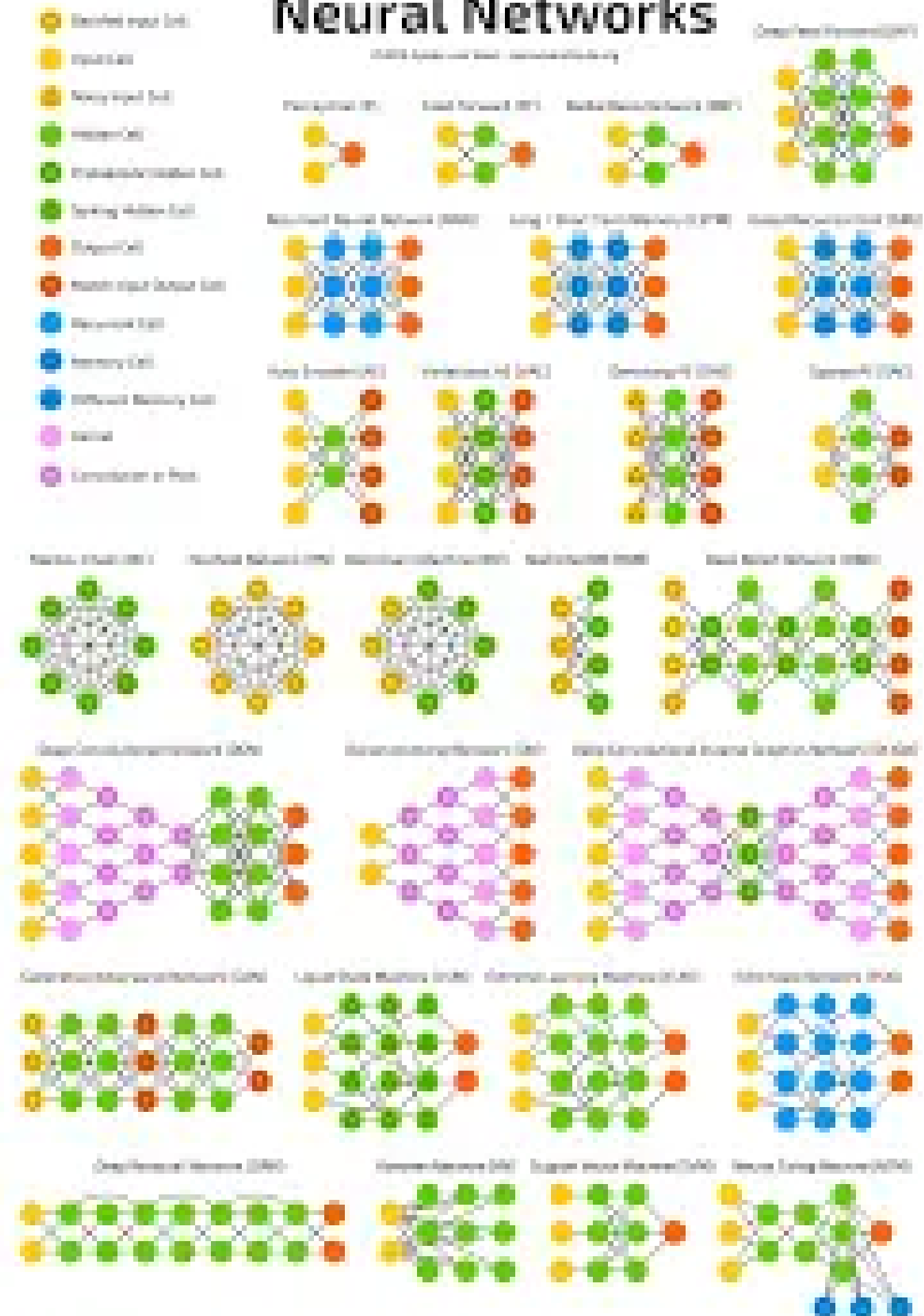
- Durante el entrenamiento de manera aleatoria se setean algunas activaciones en 0.
- << overfitting



## EARLY STOPPING

- Consiste en parar antes de llegar al overfitting observando los datos de training y testing.

\_\_\_\_\_





# HANDS-ON LAB

## TP - 1

2 problems

- Radar signal classification (mandatory)
- Flare classification (optional)

