

Ph. D. Juan David Martínez Vargas  
Ph. D. Raul Andrés Castañeda

Escuela de Ciencias Aplicadas e Ingeniería

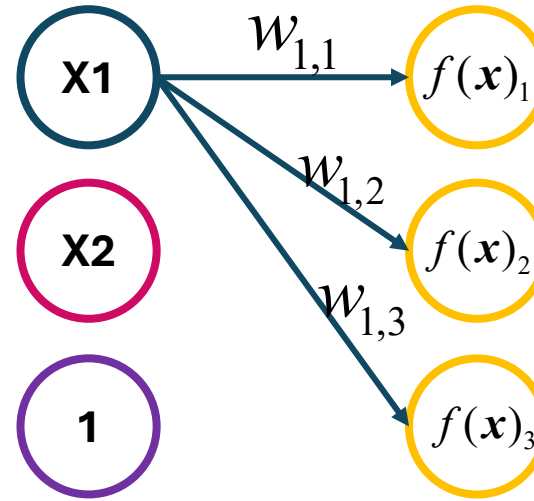
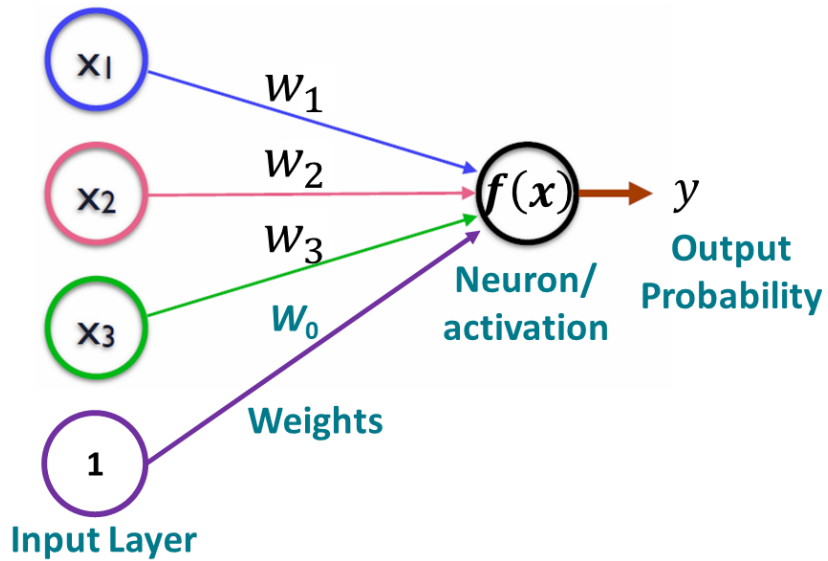
# Lecture 06

# Deep Learning

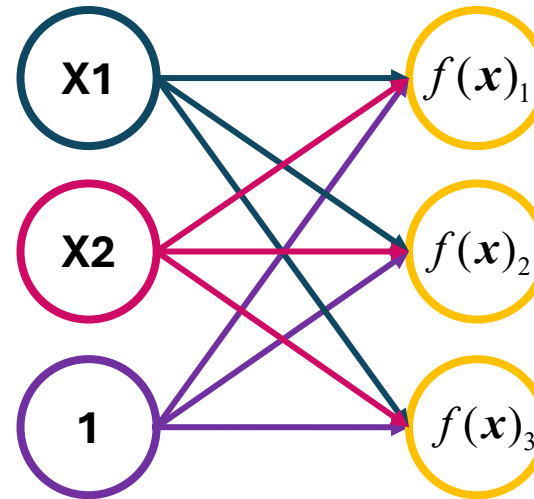
# Agenda

- **Backpropagation**
- **Métricas de evaluación**
- **Regularización**

# Backpropagation



$$f(\mathbf{x})_1 = w_{1,1}x_1 + w_{2,1}x_2 + b_1$$



$$f(\mathbf{x})_2 = w_{1,2}x_1 + w_{2,2}x_2 + b_2$$

$$f(\mathbf{x})_3 = w_{1,3}x_1 + w_{2,3}x_2 + b_3$$

$$f(x_1, x_2, x_3) = x_1 w_1 + x_2 w_2 + x_3 w_3 + w_0$$

$$f(x) = X^T \cdot W$$

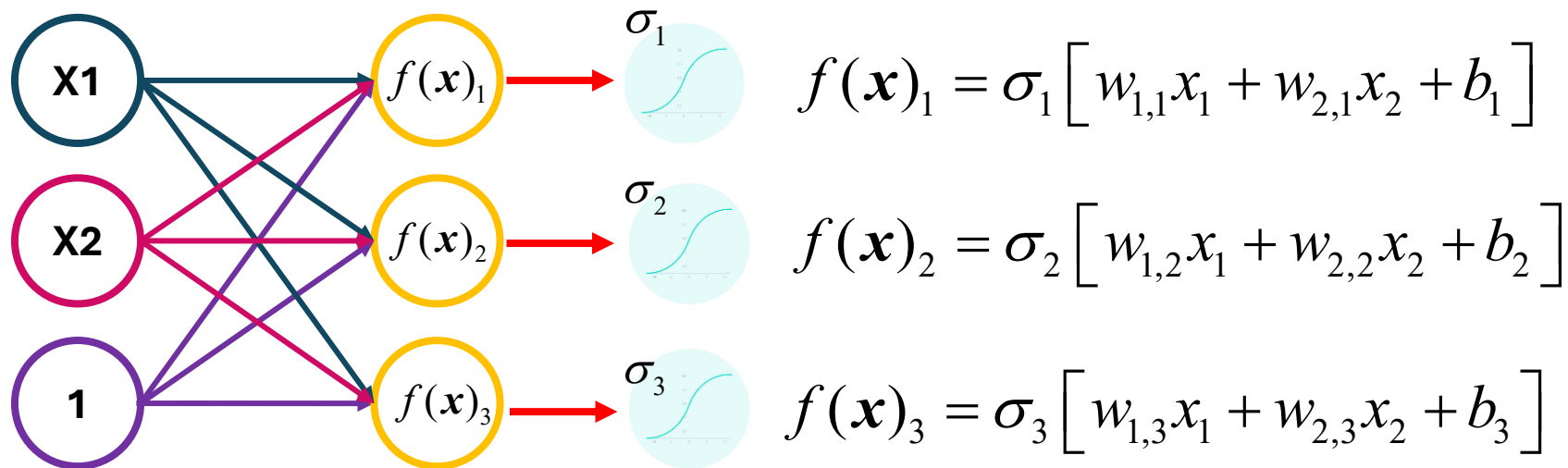
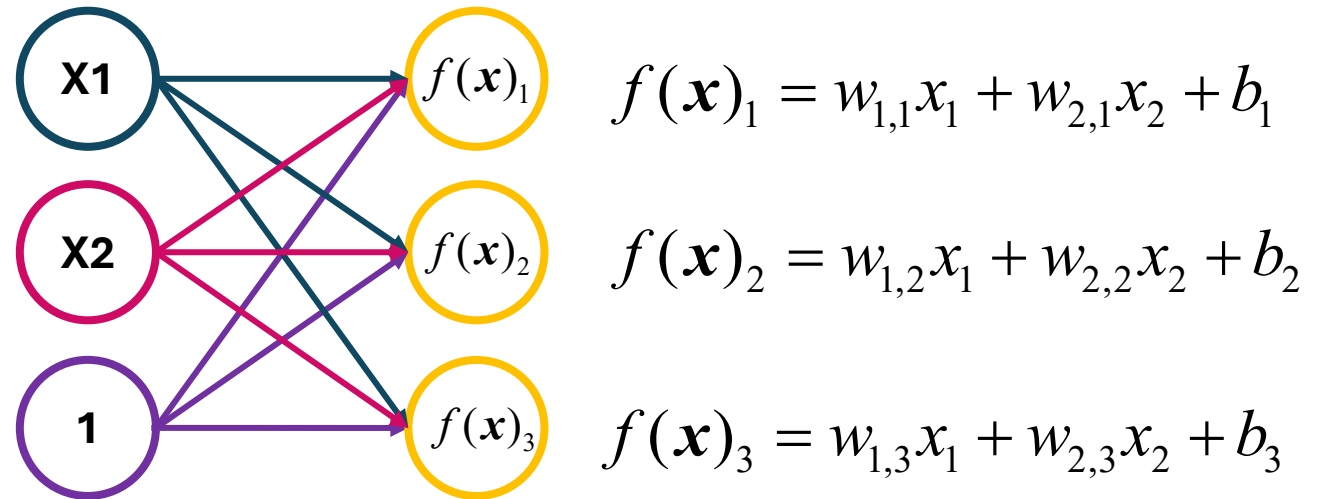
$$\hat{y} = \text{acti}[X^T \cdot W] \rightarrow \hat{y} = \sigma[X^T \cdot W]$$

# Backpropagation

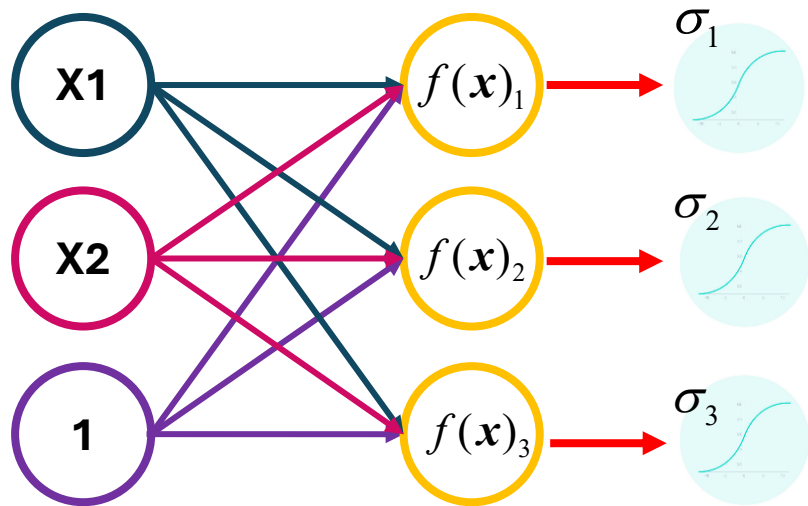
$$f(x_1, x_2, x_3) = x_1 w_1 + x_2 w_2 + x_3 w_3 + w_0$$

$$f(x) = X^T \cdot W$$

$$\hat{y} = \text{acti}[X^T \cdot W] \longrightarrow \hat{y} = \sigma[X^T \cdot W]$$



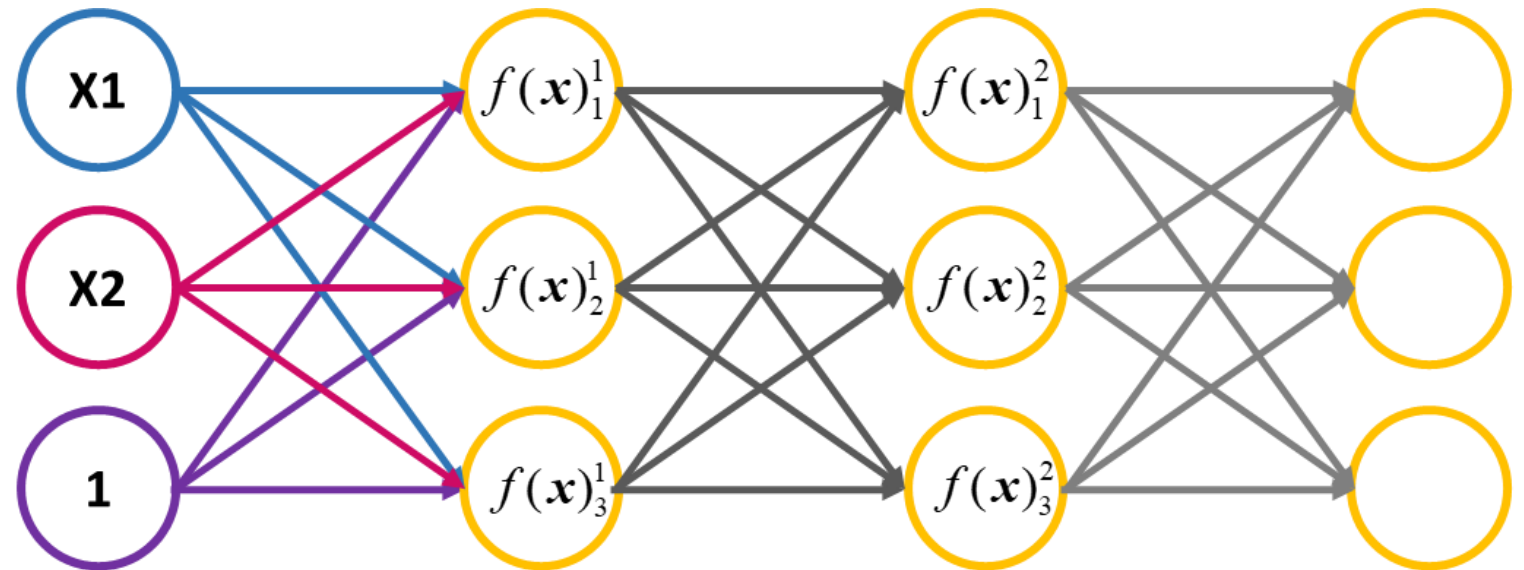
# Backpropagation



$$f(\mathbf{x})_1 = \sigma_1 [w_{1,1}x_1 + w_{2,1}x_2 + b_1]$$

$$f(\mathbf{x})_2 = \sigma_2 [w_{1,2}x_1 + w_{2,2}x_2 + b_2]$$

$$f(\mathbf{x})_3 = \sigma_3 [w_{1,3}x_1 + w_{2,3}x_2 + b_3]$$



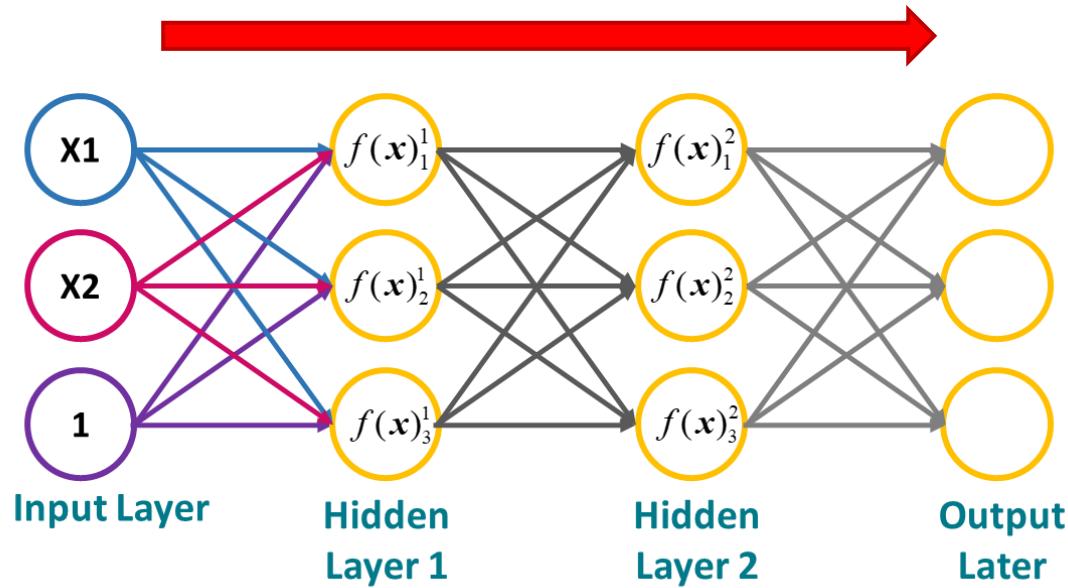
Input Layer

Hidden  
Layer 1

Hidden  
Layer 2

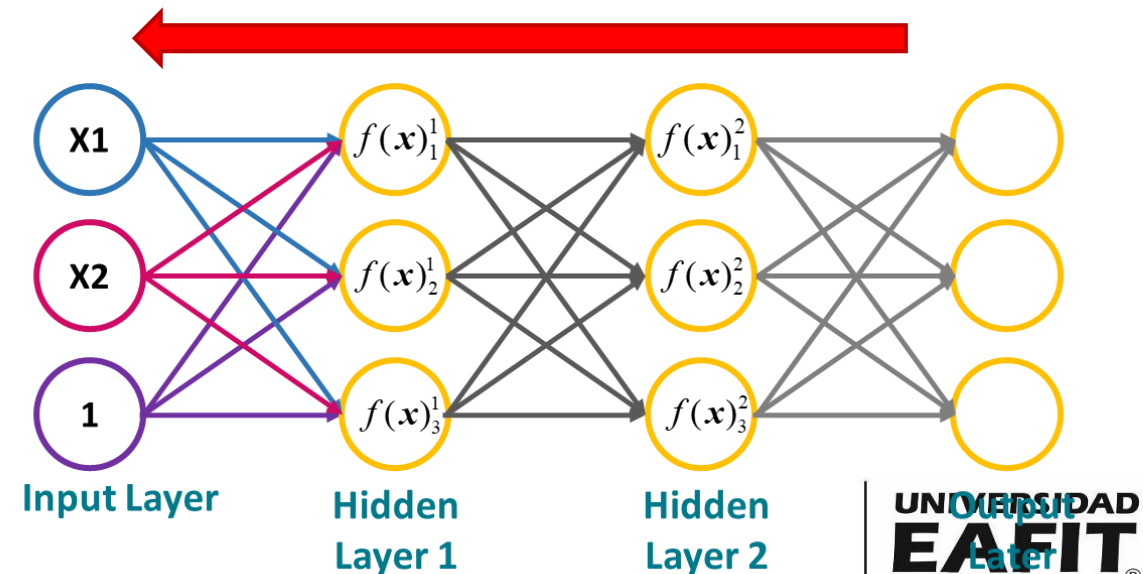
Output  
Layer

# Backpropagation



**Forward propagation** es el proceso de pasar los datos de entrada a través de la red neuronal para calcular la salida prevista. Cada capa de la red realiza una suma ponderada de sus entradas, aplica una función de activación y pasa el resultado a la siguiente capa.

**Backward propagation** es el proceso de actualización de los pesos y sesgos, en función de la pérdida calculada durante la propagación hacia adelante. Implica calcular los gradientes de la pérdida con respecto a los parámetros del modelo. Estos gradientes se utilizan para actualizar los parámetros mediante un algoritmo de optimización, como el descenso de gradientes.

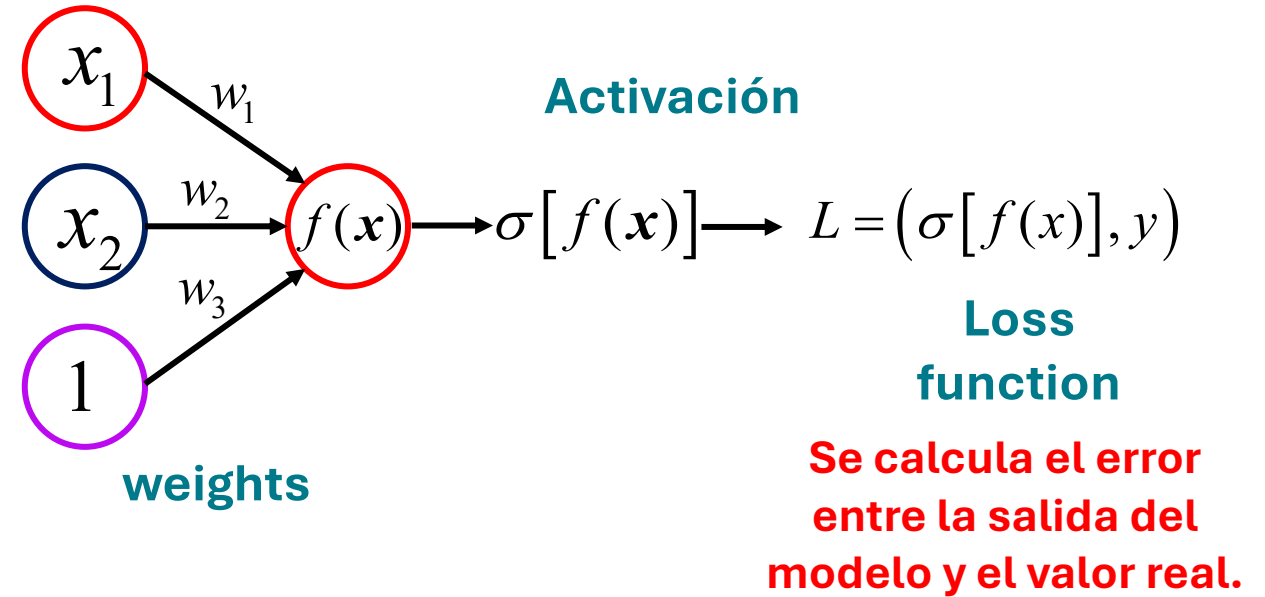


# Backpropagation

El entrenamiento de una NN se realiza en dos fases:

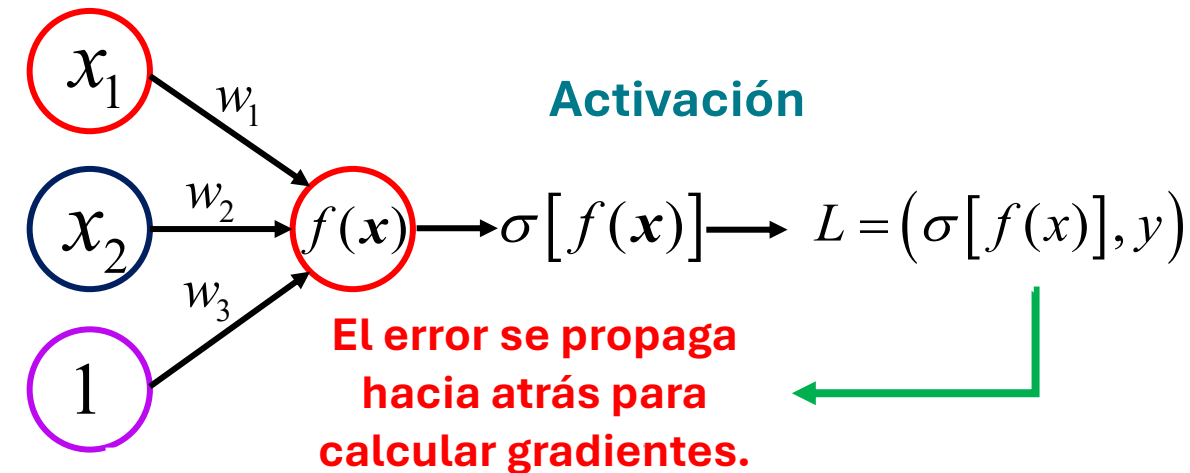
## Forward propagation:

- La señal de entrada se propaga por toda la red, capa por capa.
- Se realizan transformaciones lineales y no lineales.
- Se calculan las activaciones neuronales.
- Se obtiene la predicción del modelo.
- Genera los valores de salida de la red.

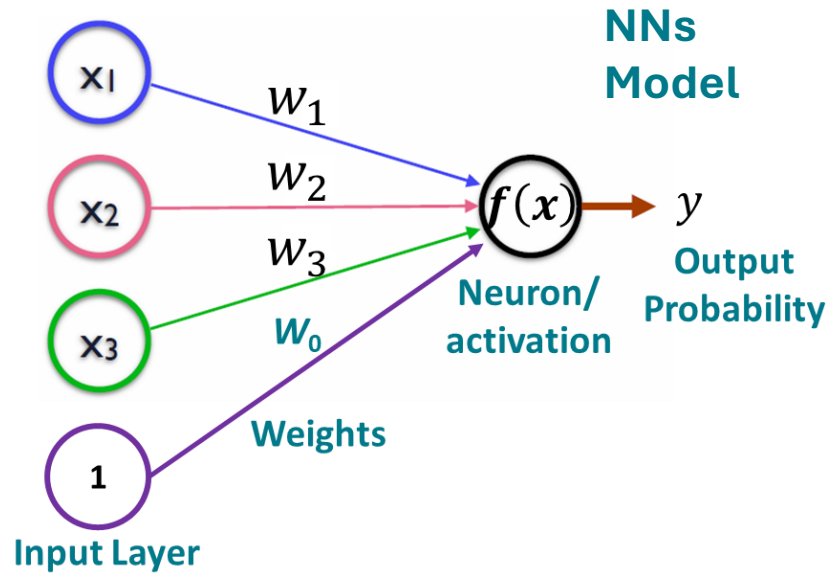


## Backward propagation:

- El error se propaga desde la capa de salida hacia la capa de entrada mediante la regla de la cadena.
- Se calculan gradientes.
- Se obtiene información para ajustar los pesos.



# Backpropagation



## Predictions

$y \rightarrow$  Real value

$\hat{y} \rightarrow$  Model prediction

## Loss function

**Mean-square error (MSE)** – Used for numerical predictions / continues data.

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

**Logistic (Cross-entropy)** – Used for probability / categorical data.

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

Find the weight that minimizes the losses (Loss function)

$$J = \frac{1}{n} \sum_{i=1}^n \wp(\hat{y}_i, y_i) \rightarrow$$

$$W = \arg \min J = \frac{1}{n} \sum_{i=1}^n L(\hat{y}_i, y_i)$$

$$W = \arg \min J = \frac{1}{n} \sum_{i=1}^n L(f(x, W)_i, y_i)$$

## Cost function

Loss quantifies error; Cost defines the optimization problem



# Backpropagation

## Backpropagation

### input:

example  $(x, y)$ , weight vector  $w$ , layered graph  $(V, E)$ ,  
activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$

### initialize:

denote layers of the graph  $V_0, \dots, V_T$  where  $V_t = \{v_{t,1}, \dots, v_{t,k_t}\}$   
define  $W_{t,i,j}$  as the weight of  $(v_{t,j}, v_{t+1,i})$   
(where we set  $W_{t,i,j} = 0$  if  $(v_{t,j}, v_{t+1,i}) \notin E$ )

### forward:

set  $o_0 = x$   
for  $t = 1, \dots, T$   
  for  $i = 1, \dots, k_t$   
    set  $a_{t,i} = \sum_{j=1}^{k_{t-1}} W_{t-1,i,j} o_{t-1,j}$   
    set  $o_{t,i} = \sigma(a_{t,i})$

### backward:

set  $\delta_T = o_T - y$   
for  $t = T - 1, T - 2, \dots, 1$   
  for  $i = 1, \dots, k_t$   
     $\delta_{t,i} = \sum_{j=1}^{k_{t+1}} W_{t+1,j,i} \delta_{t+1,j} \sigma'(a_{t+1,j})$

### output:

foreach edge  $(v_{t-1,j}, v_{t,i}) \in E$   
  set the partial derivative to  $\delta_{t,i} \sigma'(a_{t,i}) o_{t-1,j}$

[Link](#)

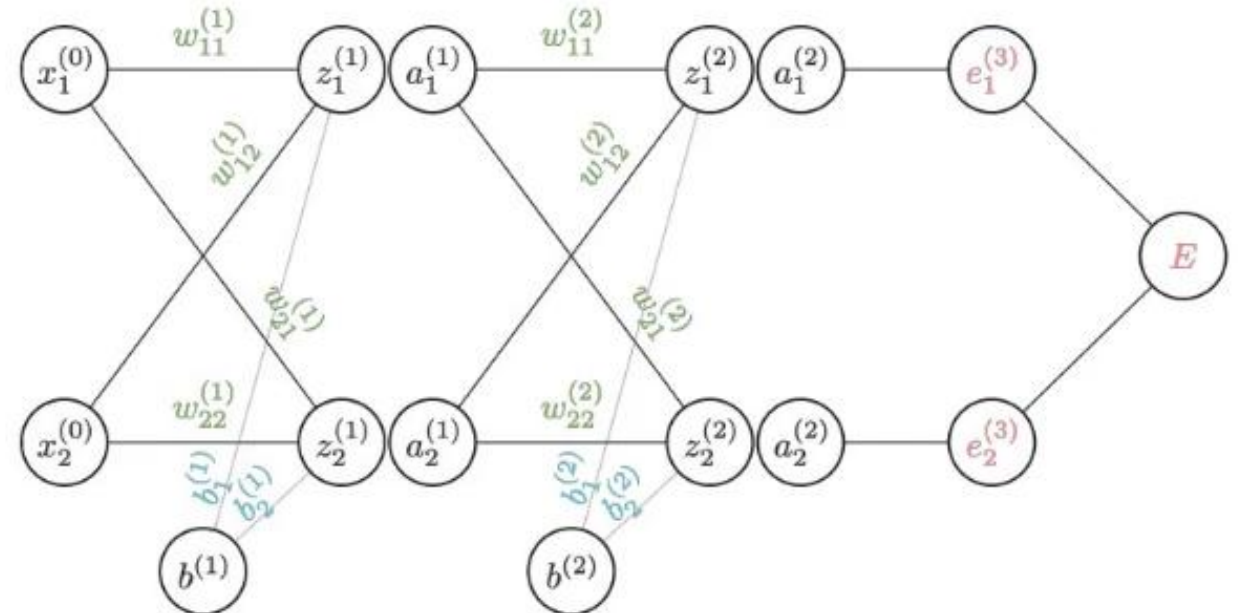
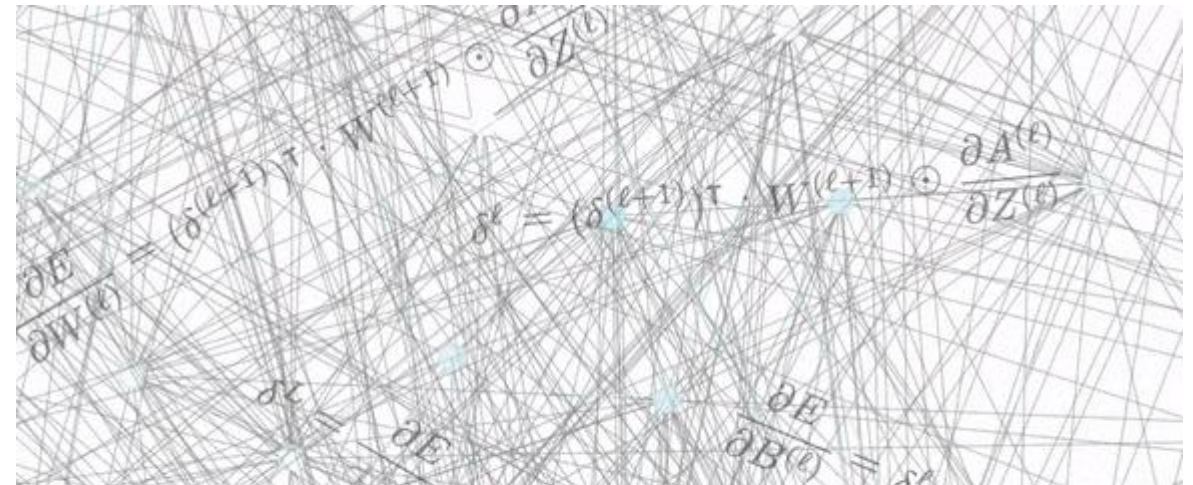
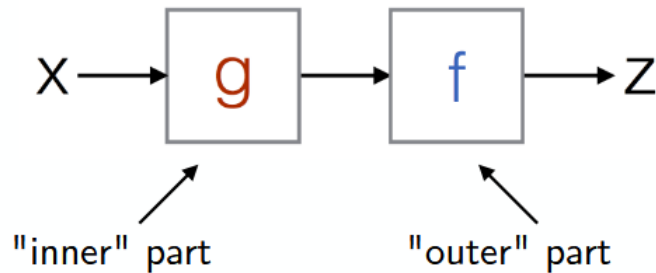


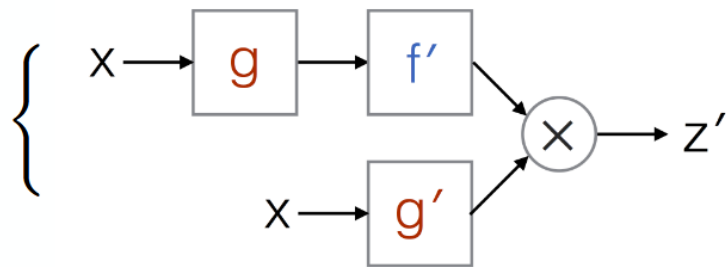
Figure 1: The example neural network setup (Image by Author)

## Regla de la cadena

$$F(x) = f(g(x)) = z$$



$$F'(x) = f'(g(x)) g'(x) = z'$$



$$\frac{d}{dx} [f(g(x))] = \frac{df}{dg} \cdot \frac{dg}{dx}$$

- Imagina que quieres convertir *pesos* a *dólares*, pero primero conviertes pesos a *euros*, y luego euros a dólares.
- Cada conversión es una función:

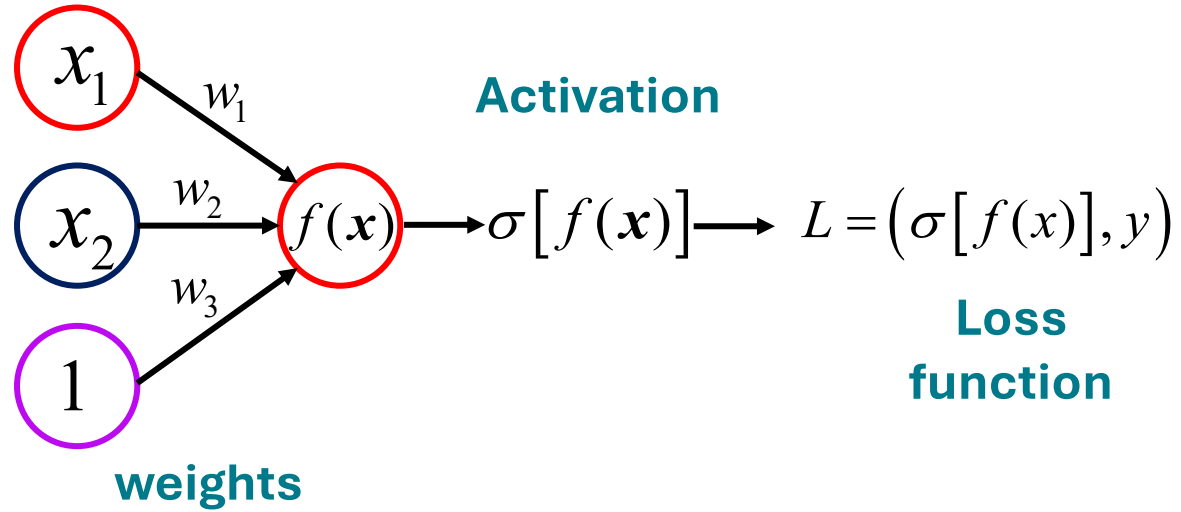
$$\text{Euros} = f(\text{Pesos}), \quad \text{Dólares} = g(\text{Euros})$$

- El cambio total de pesos a dólares se obtiene multiplicando las tasas de cambio

$$\frac{d(\text{Dólares})}{d(\text{Pesos})} = \frac{d(\text{Dólares})}{d(\text{Euros})} \cdot \frac{d(\text{Euros})}{d(\text{Pesos})}$$

- **Interpretación:** la regla de la cadena multiplica los cambios intermedios para obtener el cambio total.

# Backpropagation



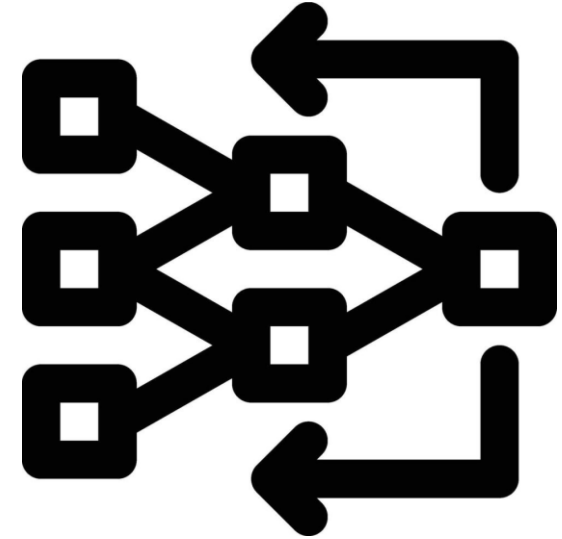
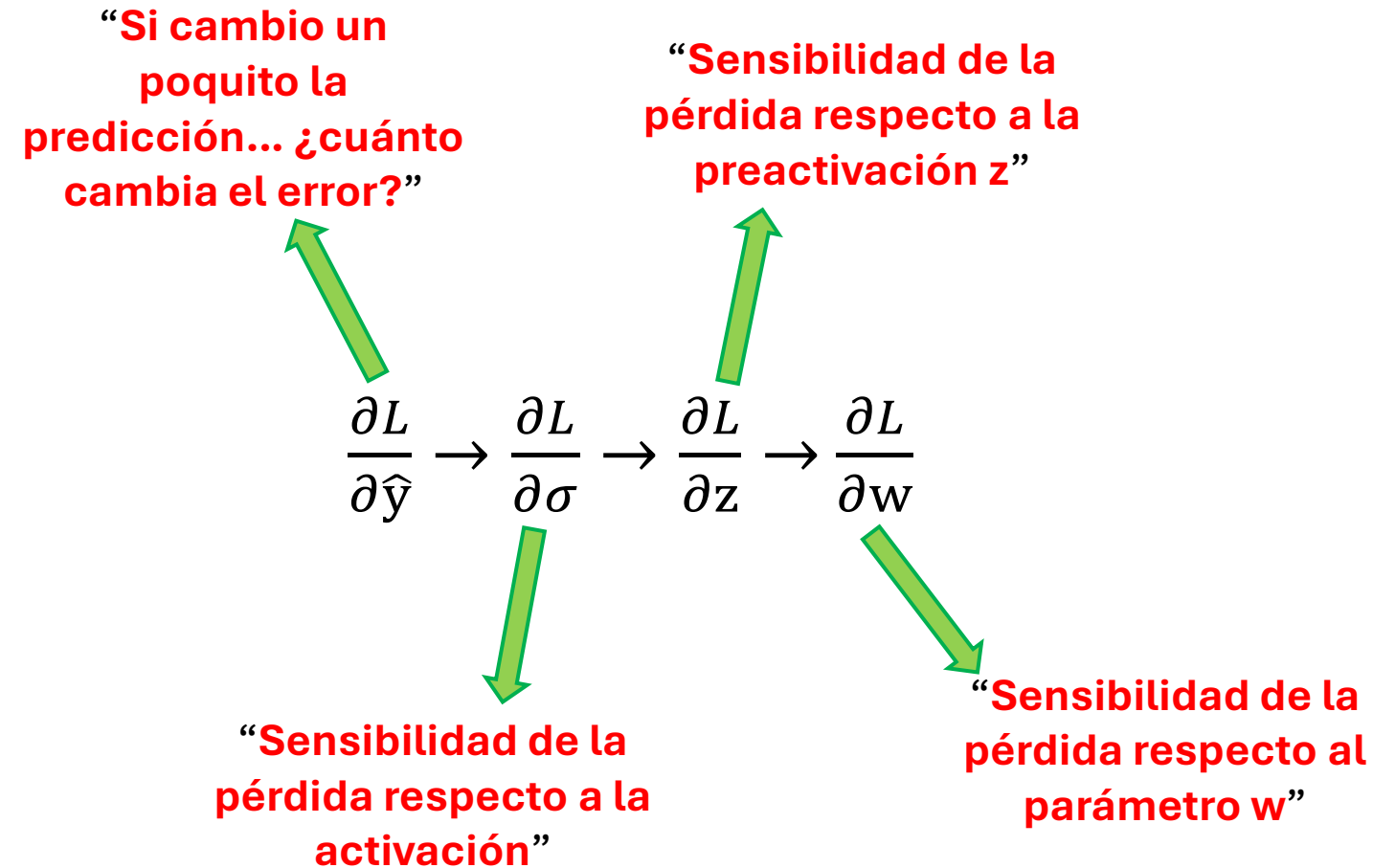
$$w_{n+1} = w - \alpha \frac{\partial L}{\partial w}$$

$$W \rightarrow Z \rightarrow \sigma \rightarrow L$$

$$\begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \rightarrow z = w_1 x_1 + w_2 x_2 + \dots + b \rightarrow \hat{y} = \sigma(z) \rightarrow L = (\hat{y}, y)$$

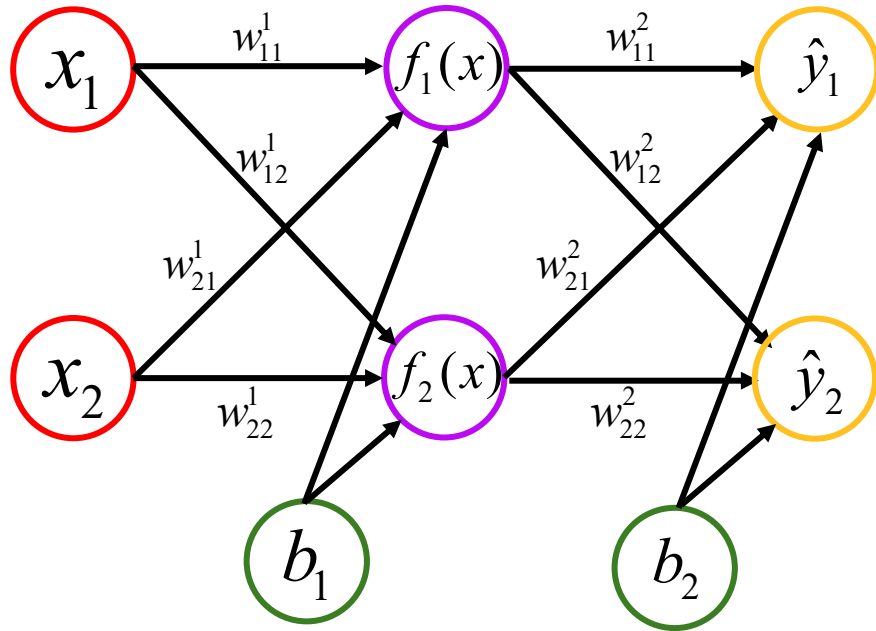
$$L = (\sigma[f(x)], y)$$

# Backpropagation



$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}$$

# Backpropagation



Supongamos sigmoidea como función de activación y una función de pérdida MSE.

$$f(x) = \frac{1}{1 + e^{-x}} \quad L = \frac{1}{2}(y - \hat{y})^2$$

$$\frac{df(x)}{dx} = \frac{d}{dx} \left[ \frac{1}{1 + e^{-x}} \right] \rightarrow \frac{df(x)}{dx} = \frac{\frac{d(1)}{dx}(1 + e^{-x}) - \frac{d(1 + e^{-x})}{dx}(1)}{(1 + e^{-x})^2} \rightarrow \frac{df(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2}$$

Expresémoslo alternatively por razones de esfuerzo computacional

$$\frac{df(x)}{dx} = \frac{1 - 1 + e^{-x}}{(1 + e^{-x})^2} \rightarrow \frac{df(x)}{dx} = \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \frac{1}{(1 + e^{-x})^2} \rightarrow \frac{df(x)}{dx} = \frac{1}{(1 + e^{-x})} - \frac{1}{(1 + e^{-x})^2}$$

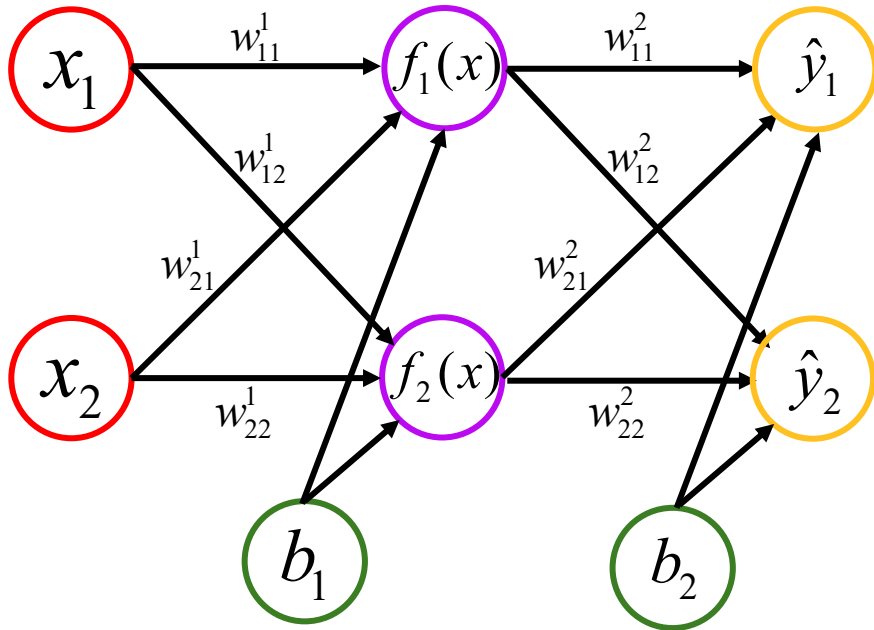
$$\frac{df(x)}{dx} = \frac{1}{(1 + e^{-x})} \left[ 1 - \frac{1}{(1 + e^{-x})} \right] \rightarrow \boxed{f(x)[1 - f(x)]}$$

**IMPORTANTE**

$$\frac{d}{dx} \left[ \frac{g(x)}{f(x)} \right] = \frac{g'(x)f(x) - g(x)f'(x)}{f(x)^2}$$

$$w_{n+1} = w - \alpha \frac{\partial L}{\partial w}$$

# Backpropagation



$$L = \frac{1}{2}(y - \hat{y})^2 \rightarrow \frac{\partial L}{\partial \hat{y}} = \frac{1}{2} \frac{d(y - \hat{y})^2}{d\hat{y}}$$

$$\frac{1}{2} \frac{d(y - \hat{y})^2}{d\hat{y}} = \frac{1}{2} (-2(y - \hat{y}))$$

$$\frac{\partial L}{\partial \hat{y}} = \boxed{(\hat{y} - y)}$$

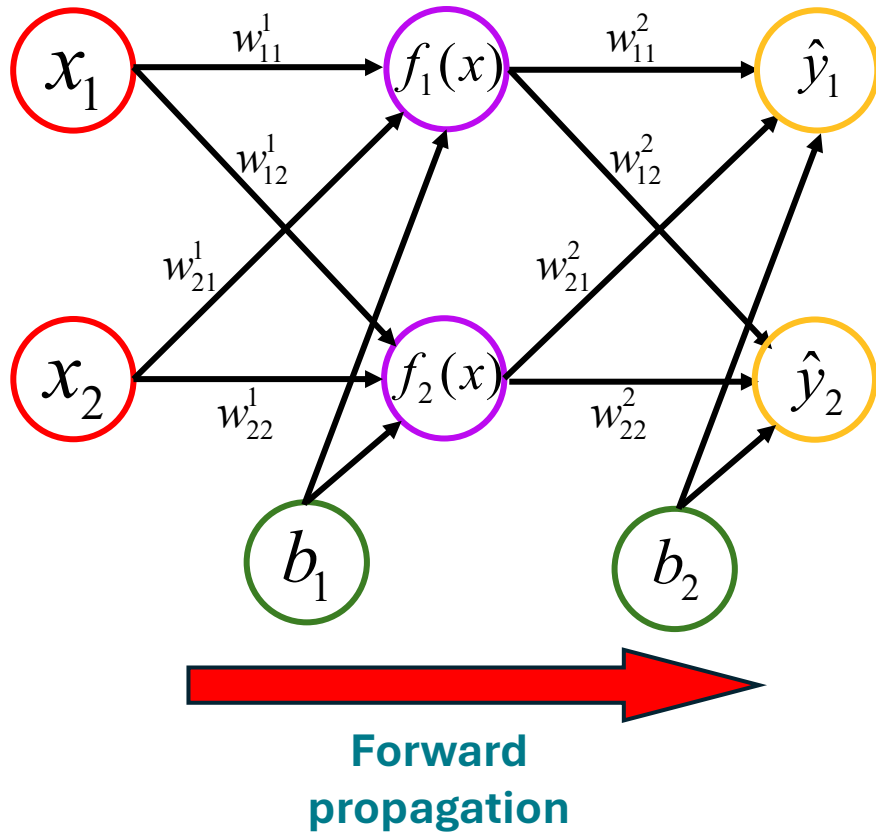
Inputs / Outputs / Learning Rate

$x_1 = 0.05$	$y_1 = 0.01$	$\alpha = 0.5$
$x_2 = 0.10$	$y_2 = 0.99$	

Initialize values

$w_{11}^1 = 0.15$	$w_{11}^2 = 0.40$	$b_1 = 0.35$
$w_{12}^1 = 0.25$	$w_{12}^2 = 0.50$	
$w_{21}^1 = 0.20$	$w_{21}^2 = 0.45$	$b_2 = 0.60$
$w_{22}^1 = 0.30$	$w_{22}^2 = 0.55$	

# Backpropagation



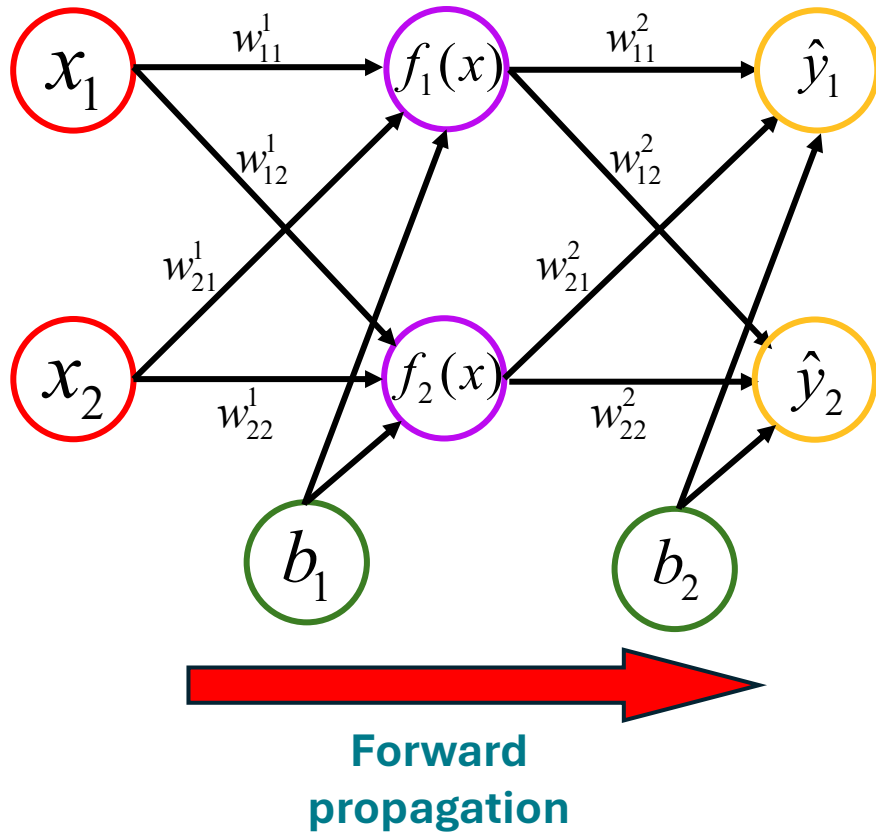
$$f_1(x) = x_1 w_{11}^1 + x_2 w_{21}^1 + b_1$$

$$f_1(x) = (0.05)(0.15) + (0.10)(0.20) + 0.35 \rightarrow f_1(x) = 0.3775$$

$$\sigma[f_1(x)] = \frac{1}{1 + e^{-f_1(x)}} \rightarrow \sigma[f_1(x)] = \frac{1}{1 + e^{-0.3775}}$$

$$\sigma[f_1(x)] = 0.59326999211$$

# Backpropagation



$$f_1(x) = x_1 w_{11}^1 + x_2 w_{21}^1 + b_1$$

$$f_1(x) = (0.05)(0.15) + (0.10)(0.20) + 0.35 \rightarrow f_1(x) = 0.3775$$

$$\sigma[f_1(x)] = \frac{1}{1 + e^{-f_1(x)}} \rightarrow \sigma[f_1(x)] = \frac{1}{1 + e^{-0.3775}}$$

$$\sigma[f_1(x)] = 0.59326999211$$

$$f_2(x) = x_1 w_{12}^1 + x_2 w_{22}^1 + b_1$$

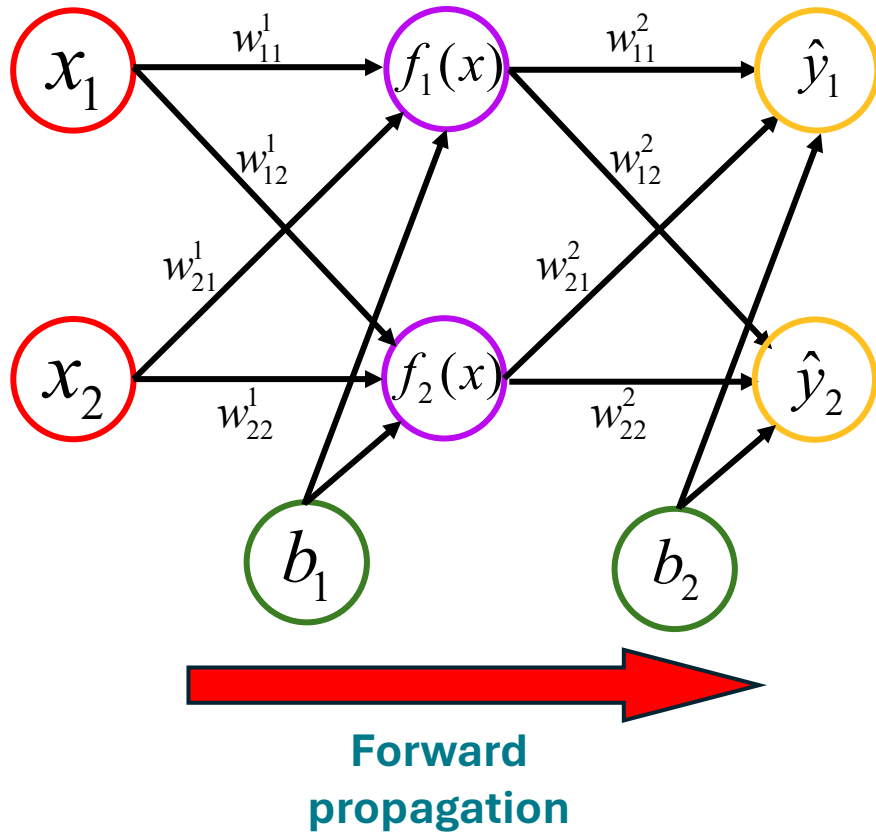
$$f_2(x) = (0.05)(0.25) + (0.10)(0.30) + 0.35 \rightarrow f_2(x) = 0.3925$$

$$\sigma[f_2(x)] = \frac{1}{1 + e^{-f_2(x)}} \rightarrow \sigma[f_2(x)] = \frac{1}{1 + e^{-0.3925}}$$

$$\sigma[f_2(x)] = 0.59688437826$$



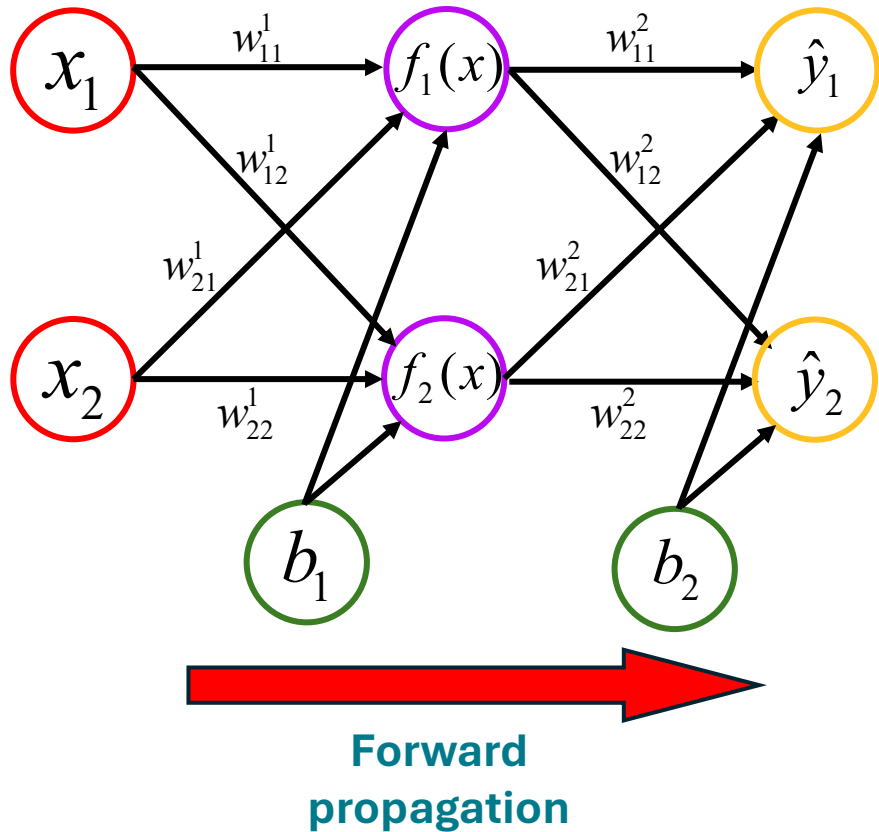
# Backpropagation



$$\sigma[f_1(x)] = 0.59326999211$$

$$\sigma[f_2(x)] = 0.59688437826$$

# Backpropagation



$$\sigma[f_1(x)] = 0.59326999211$$

$$\sigma[f_2(x)] = 0.59688437826$$

$$\hat{y}_1 = \sigma[f_1(x)]w_{11}^2 + \sigma[f_2(x)]w_{21}^2 + b_2$$

$$\hat{y}_1 = (0.5932)(0.40) + (0.5968)(0.45) + 0.60$$

$$\hat{y}_1 = 1.10590596706$$

$$\sigma[\hat{y}_1] = \frac{1}{1 + e^{-\hat{y}_1}} \rightarrow \sigma[\hat{y}_1] = \frac{1}{1 + e^{-1.1059}} \rightarrow \sigma[\hat{y}_1] = 0.75136506955$$

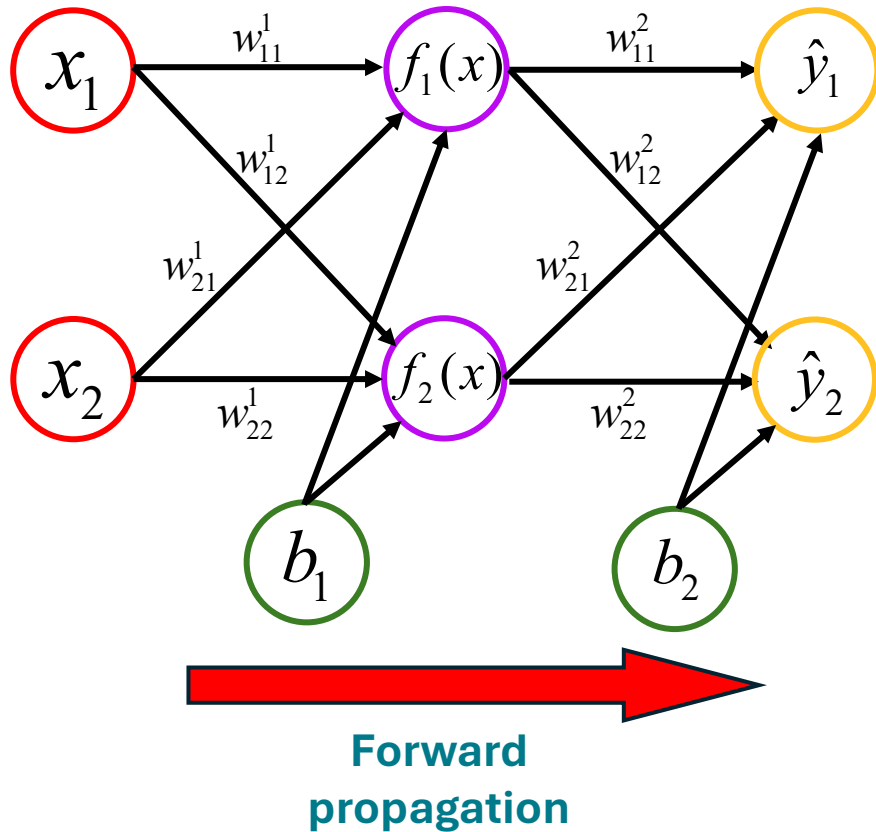
$$\hat{y}_2 = \sigma[f_1(x)]w_{12}^2 + \sigma[f_2(x)]w_{22}^2 + b_2$$

$$\hat{y}_2 = (0.5932)(0.50) + (0.5968)(0.55) + 0.60$$

$$\hat{y}_2 = 1.22492140407$$

$$\sigma[\hat{y}_2] = \frac{1}{1 + e^{-\hat{y}_2}} \rightarrow \sigma[\hat{y}_2] = \frac{1}{1 + e^{-1.2249}} \rightarrow \sigma[\hat{y}_2] = 0.77292846532$$

# Backpropagation



$$E_T = \sum_{i=1}^n E_i = \frac{1}{2}(y_1 - \hat{y}_1)^2 + \frac{1}{2}(y_2 - \hat{y}_2)^2$$

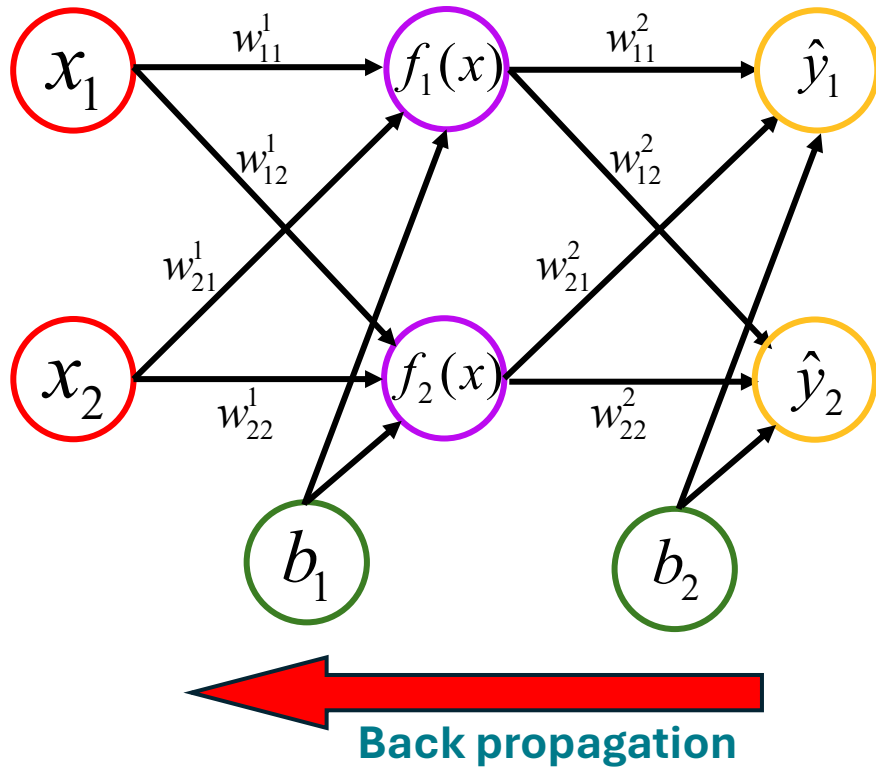
$$E_T = \frac{1}{2}(0.01 - 0.7513)^2 + \frac{1}{2}(0.99 - 0.7729)^2 = 0.29837110876$$

**Aunque la salida para el segundo dato predicha es “similar”, para el primer dato no lo es.. El modelo presenta un error significativo respecto a los valores objetivo.**

$$\sigma[\hat{y}_1] = 0.75136506955 \quad \hat{y}_1$$

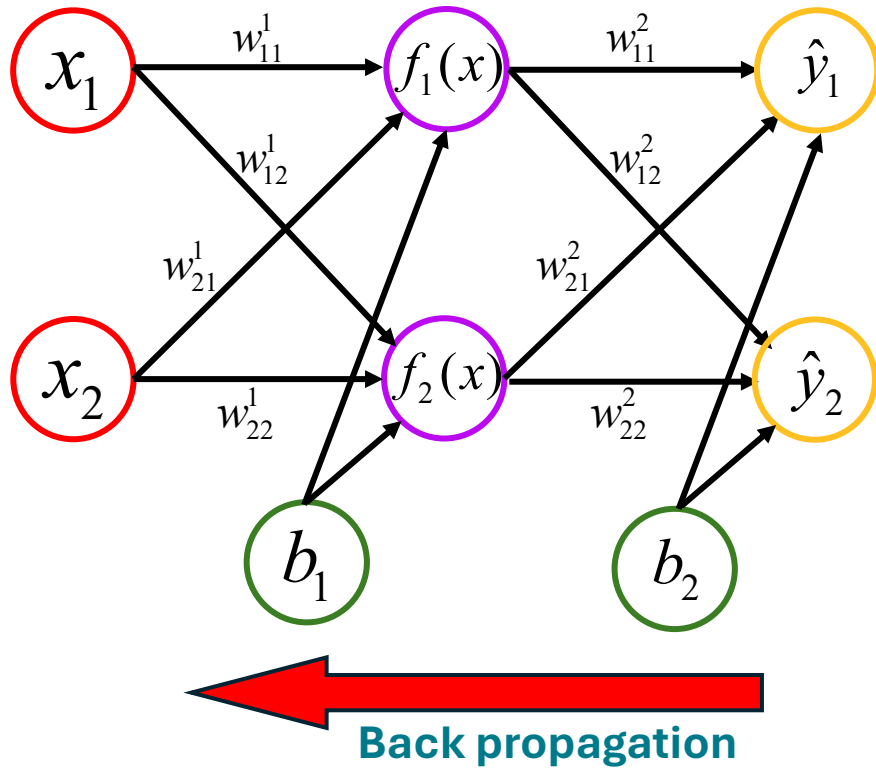
$$\sigma[\hat{y}_2] = 0.77292846532 \quad \hat{y}_2$$

# Backpropagation



$$E_T = \frac{1}{2}(0.01 - 0.7573)^2 + \frac{1}{2}(0.99 - 0.7729)^2 = 0.29837110876$$

# Backpropagation



$$E_T = \frac{1}{2}(0.01 - 0.7573)^2 + \frac{1}{2}(0.99 - 0.7729)^2 = 0.29837110876$$

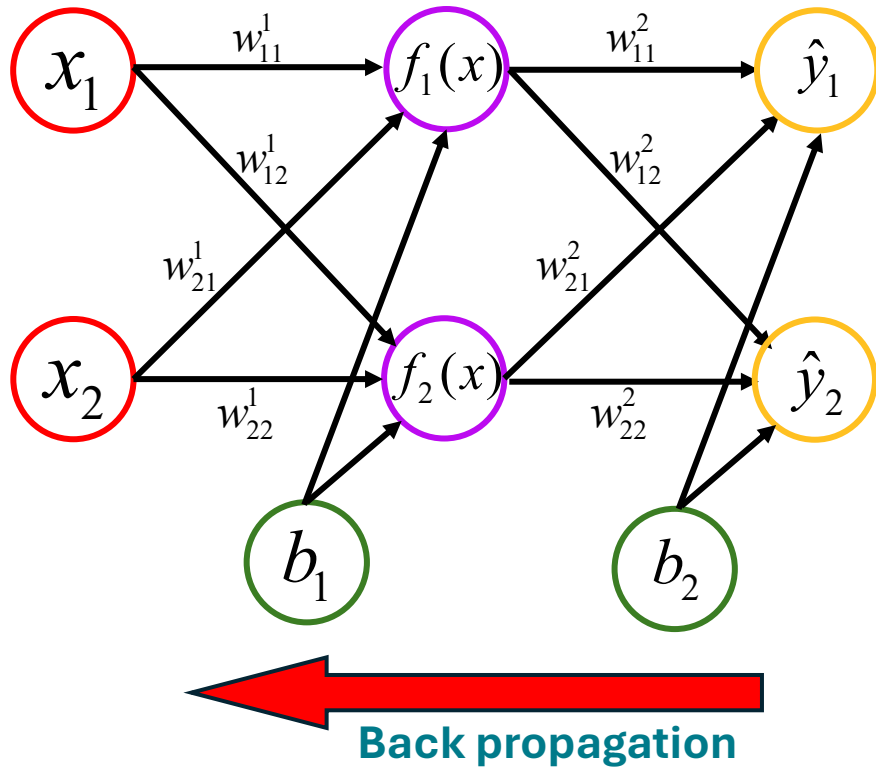
**Backpropagation calcula gradientes para cada parámetro**

$$[w_{11}^2, w_{12}^2, w_{21}^2, w_{22}^2]$$

$$[w_{11}^1, w_{12}^1, w_{21}^1, w_{22}^1]$$

$$w_{n+1} \leftarrow w - \alpha \frac{\partial E}{\partial w} \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}$$

# Backpropagation



$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}$$

$$[w_{11}^2, w_{12}^2, w_{21}^2, w_{22}^2]$$

$$w_{11}^{2'} = w_{11}^2 - (0.5) \frac{\partial L}{\partial w_{11}^2}$$

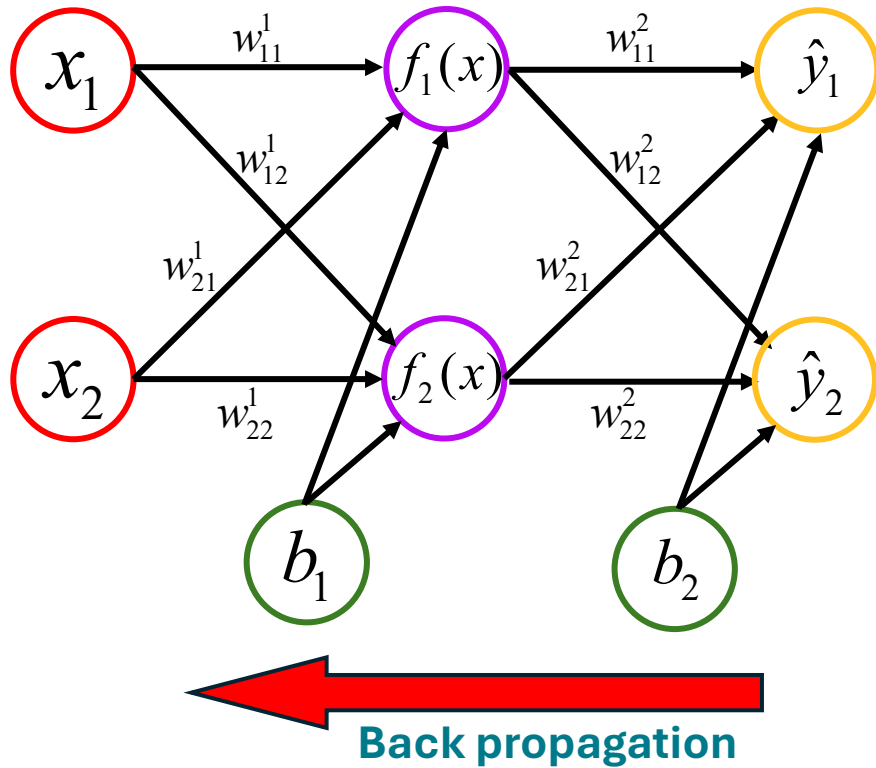
$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}^2} \rightarrow \frac{\partial L}{\partial w_{11}^2} = (\hat{y}_1 - y) (\sigma(z_1) [1 - \sigma(z_1)]) p$$

Diagram illustrating the backpropagation of error gradients through the network layers:

- Output Layer:** The error gradient  $\frac{\partial L}{\partial \hat{y}_1} = (\hat{y}_1 - y)$  is calculated.
- Hidden Layer:** The error gradient  $\frac{\partial \hat{y}_1}{\partial z_1} = \sigma(z_1) [1 - \sigma(z_1)]$  is calculated.
- Input Layer:** The error gradient  $\frac{\partial z_1}{\partial w_{11}^2} = \frac{\partial}{\partial w_{11}^2} [\sigma[f_1(x)]w_{11}^2 + \sigma[f_2(x)]w_{21}^2 + b_2]$  is calculated.

The final error gradient  $\frac{\partial L}{\partial w_{11}^2}$  is the product of these three gradients, as shown in the equation above.

# Backpropagation



$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \frac{\partial z}{\partial w}$$

$$[w_{11}^2, w_{12}^2, w_{21}^2, w_{22}^2]$$

$$w_{11}^{2'} = w_{11}^2 - (0.5) \frac{\partial L}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_1} \frac{\partial z_1}{\partial w_{11}^2} \rightarrow \frac{\partial L}{\partial w_{11}^2} = (\hat{y}_1 - y)(\sigma(z_1)[1 - \sigma(z_1)])p$$

$$\frac{\partial L}{\partial w_{11}^2} = (0.7513 - 0.01)(0.7513)[1 - 0.7513])(0.5932)$$

$$\frac{\partial L}{\partial w_{11}^2} = (0.7513 - 0.01)(0.7513)[1 - 0.7513])(0.5932)$$

$$\frac{\partial L}{\partial w_{11}^2} = 0.082167041$$

$$w_{11}^{2'} = 0.40 - (0.5)(0.082167041)$$

$$w_{11}^{2'} = 0.35891$$

$$w_{11}^{2'} = 0.35891$$

$$w_{21}^{2'} = 0.511301270$$

$$w_{12}^{2'} = 0.408666186$$

$$w_{22}^{2'} = 0.561370121$$

## Actividad

Implemente el notebook `backpropagation_example.ipynb` para realizar las siguientes actividades:

1. Sin modificar el código, determine los valores de los pesos tras una segunda iteración.

**Importante:** realice este cálculo **sin alterar el código**.

2. Modifique el código para calcular el error después de dos iteraciones.



# Sobre y sub entrenamiento

**Sub-entrenamiento (Underfitting):** El modelo no es lo suficientemente bueno para describir las relaciones entre los datos de entrada  $\mathbf{X}$  y la variable de salida  $\mathbf{y}$ .



- El modelo es muy simple para capturar patrones importantes en los datos.
- El modelo tendrá un rendimiento bajo en los conjuntos de entrenamiento y validación.

**Sobreentrenamiento (Overfitting):** El modelo memoriza o imita los datos de entrenamiento, y falla generalizando para datos desconocidos (datos de test).



- El modelo es demasiado complejo.
- El modelo aprende patrones de ruido y no de relaciones entre variables.
- Tendrá un rendimiento muy bueno en los datos de entrenamiento pero muy malo en los datos de validación o prueba.

# Sobre y sub entrenamiento

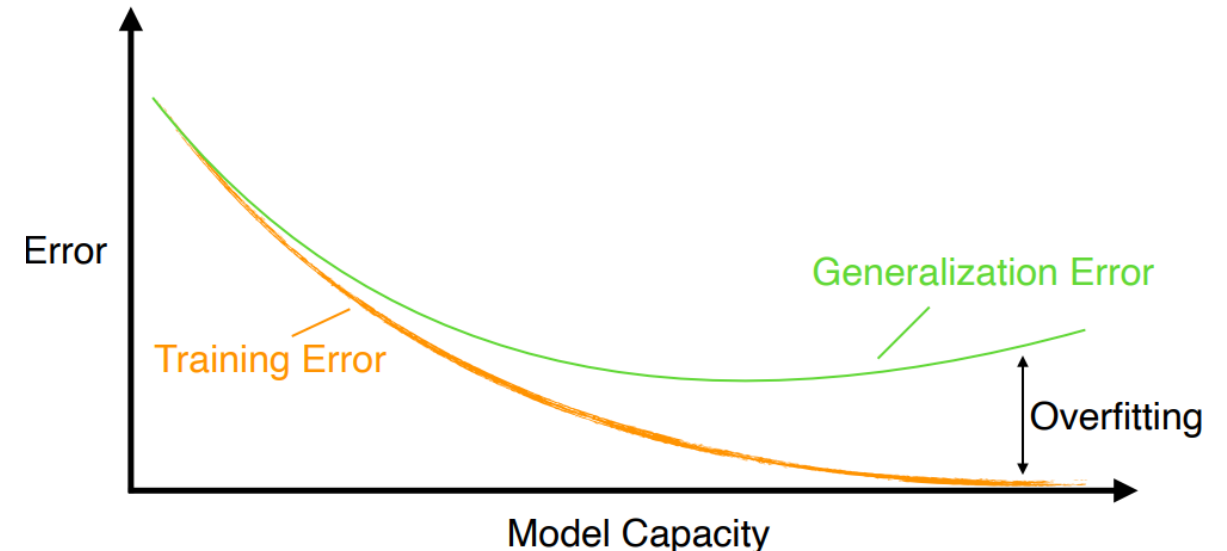
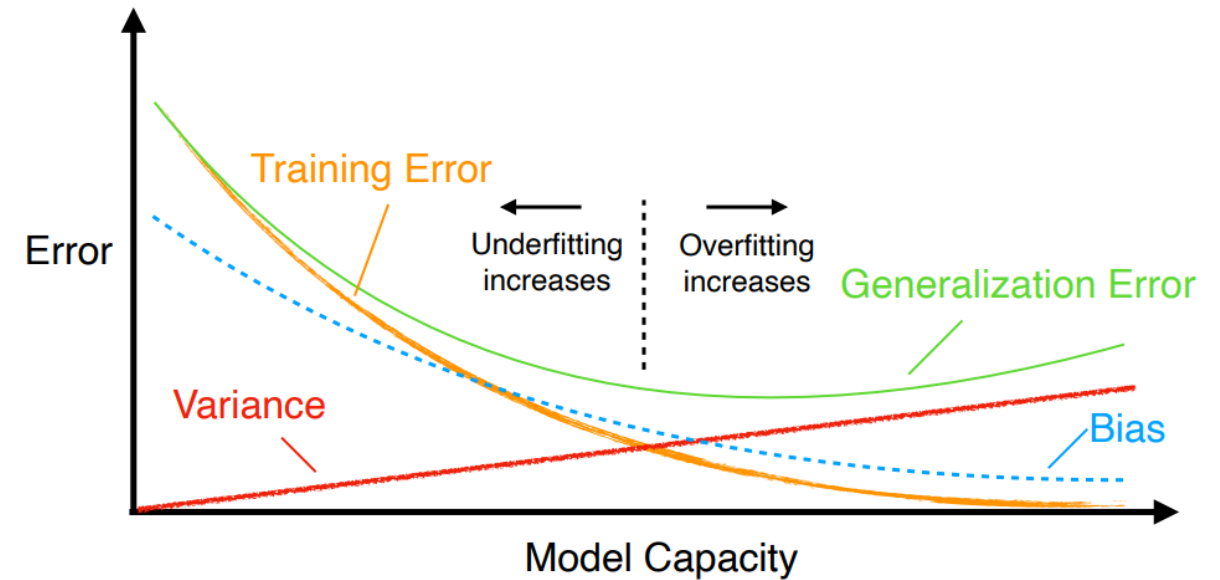
**Buen ajuste:** El modelo captura las relaciones generales entre los datos de entrada  $\mathbf{X}$  y la variable de salida  $\mathbf{y}$ .



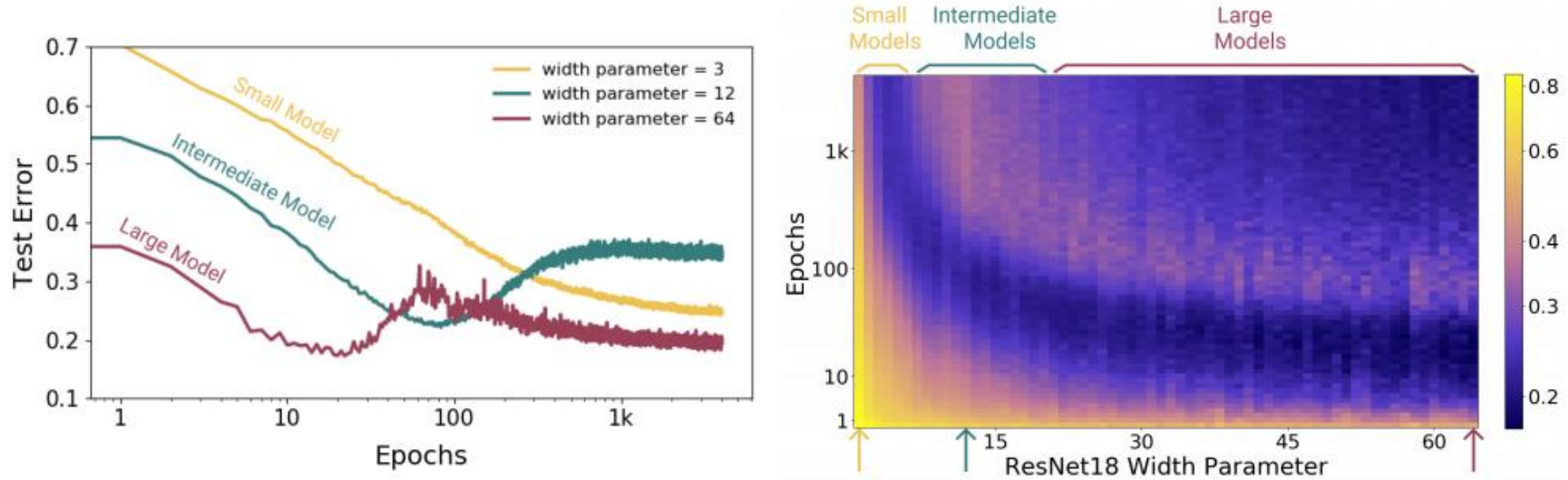
- El modelo no es ni muy simple ni muy complejo.
- El modelo descifra las relaciones subyacentes de los datos de entrenamiento y no el ruido.
- El modelo tendrá un buen rendimiento en los conjuntos de entrenamiento, validación y prueba.

# Sobre y sub entrenamiento

- **Model Capacity: que tan complejos es el modelo.**
  - A la izquierda → modelos simples.
  - A la derecha → modelos complejos.
- El error de entrenamiento **siempre disminuye** cuando aumentas la capacidad.
- El Generalization error este es **el error que realmente importa** (test / datos nuevos)..
- **Variance** mide sensibilidad al ruido.
  - Variance alta = modelo demasiado inestable.
  - Modelos simples → variance baja
  - Modelos complejos → variance alta
  - Modelo complejo = pequeñas variaciones en datos, grandes cambios en predicción.
- Usualmente se usa el error del conjunto de Test como estimador del error de generalización.

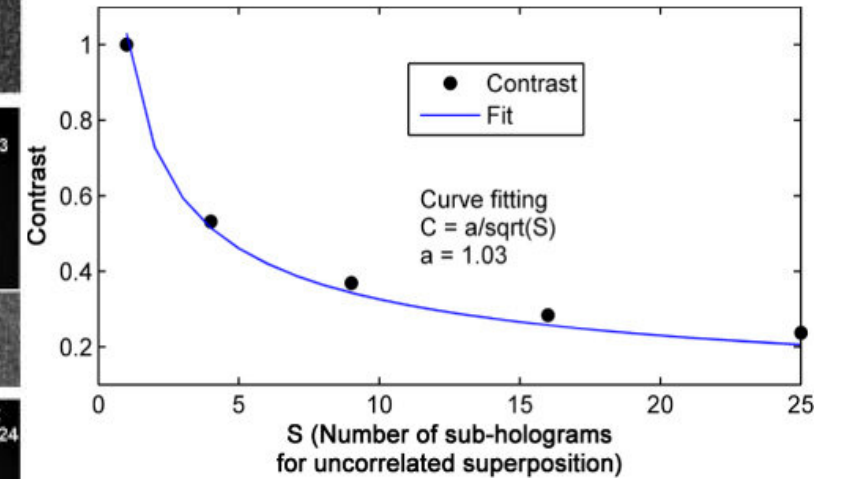
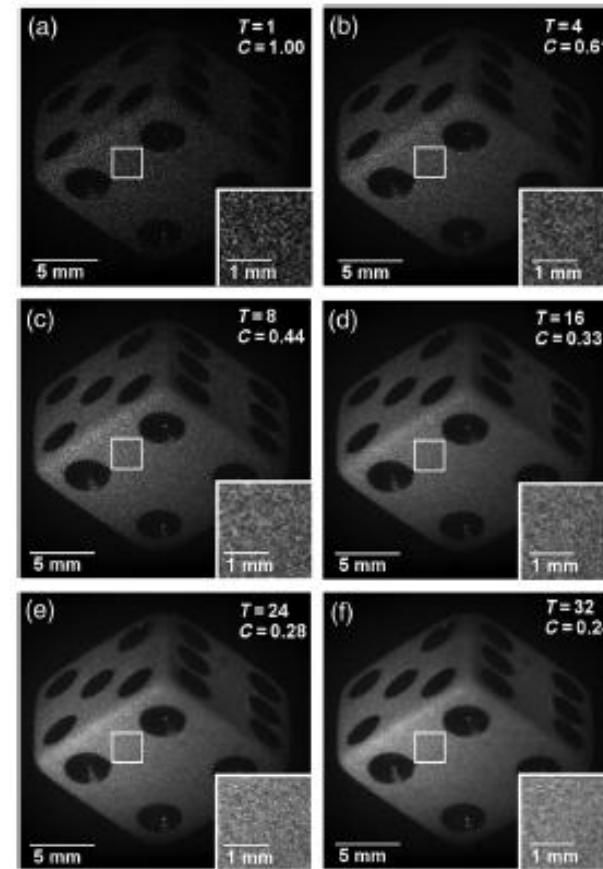
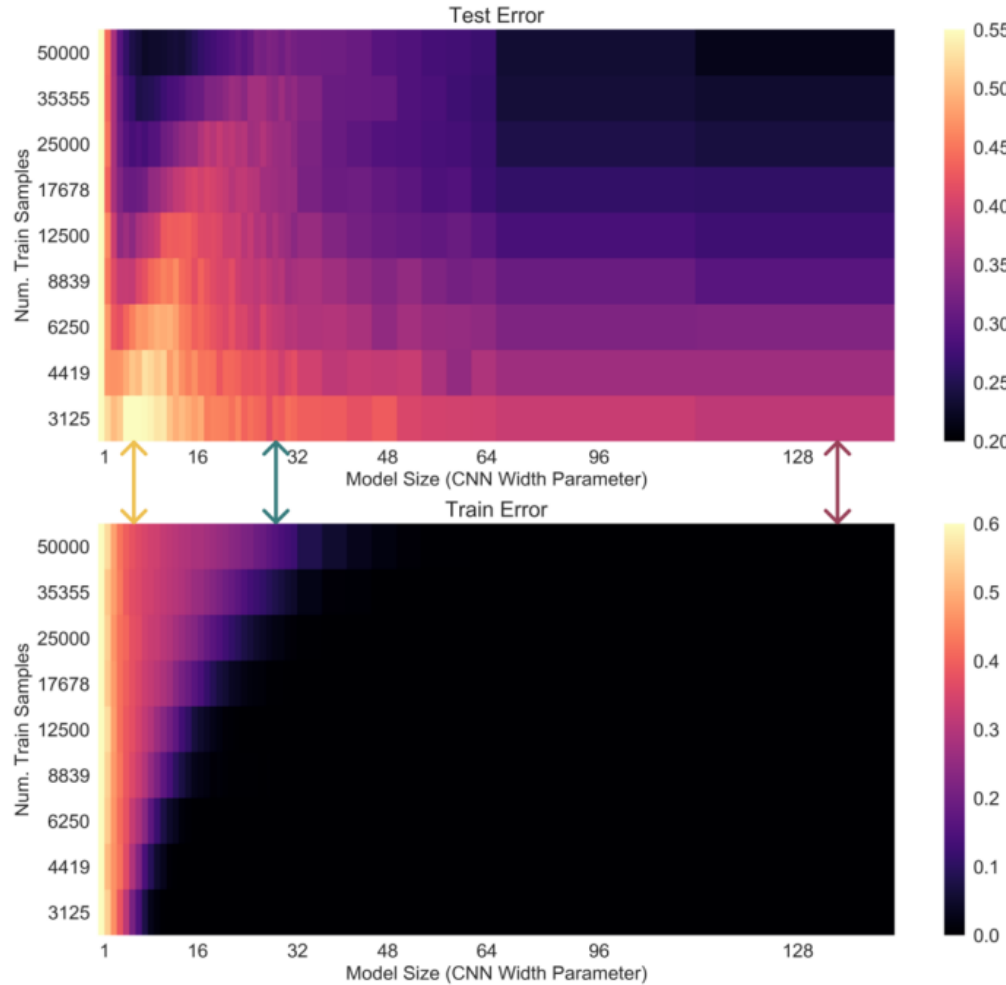


# Sobre y sub entrenamiento

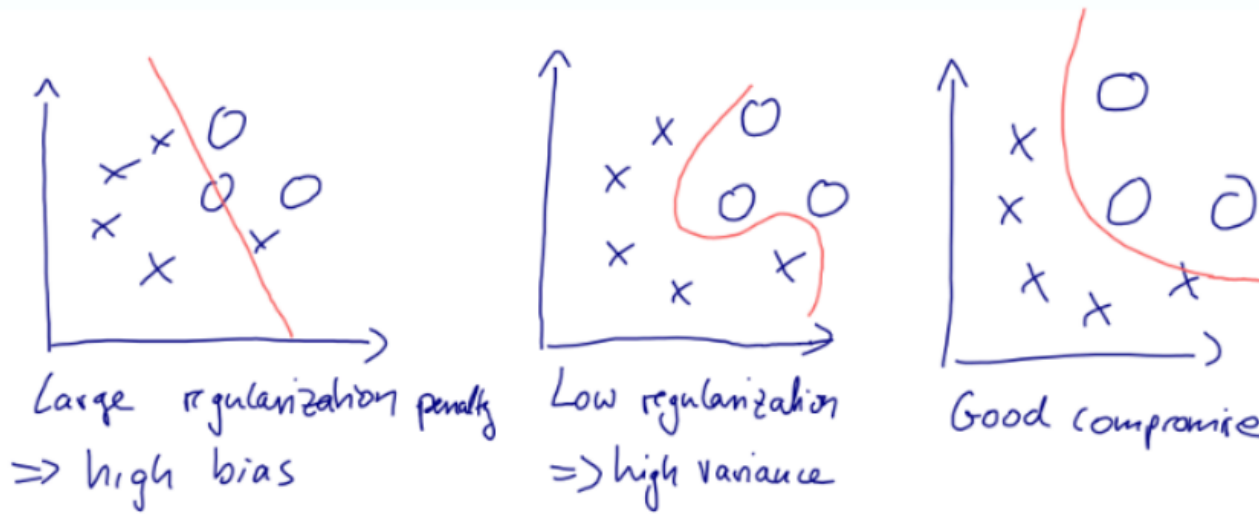


- En la región crítica, solo un modelo se ajusta bien a los datos y es muy sensible al ruido
- Modelos sobre parametrizados: muchos modelos se ajustan bien a los datos, SGD encuentra uno que memoriza los datos de entrenamiento pero que también se desempeña bien en los datos de prueba.

# Sobre y sub entrenamiento



# Conjunto de datos



## Dataset

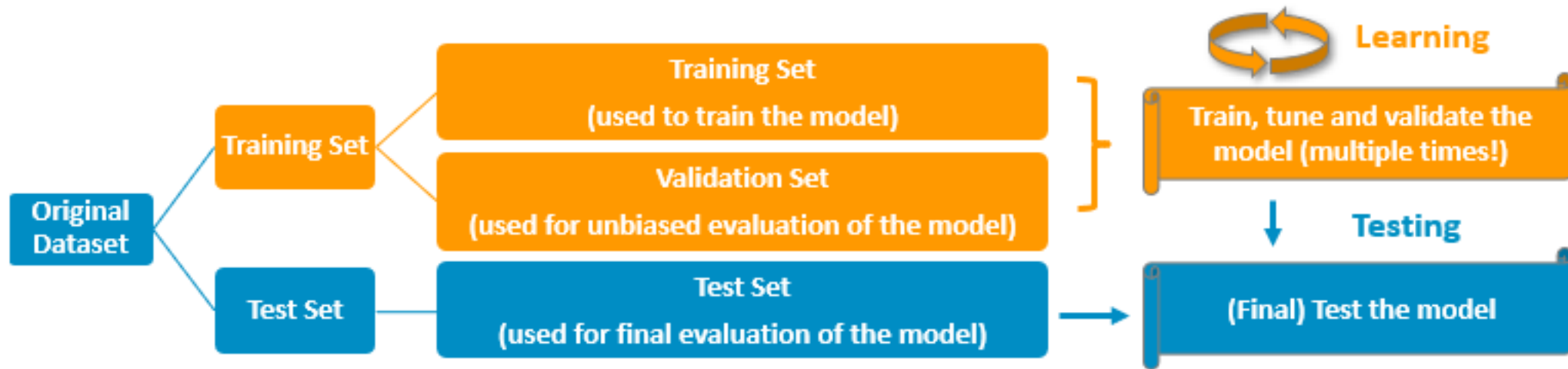


Regla común: 80 / 15 / 5

- Training set  $\rightarrow$  aprendizaje del modelo
- Validation set  $\rightarrow$  estimación del desempeño
- Test set  $\rightarrow$  evaluación final

🎯 Test = estimador del error de generalización

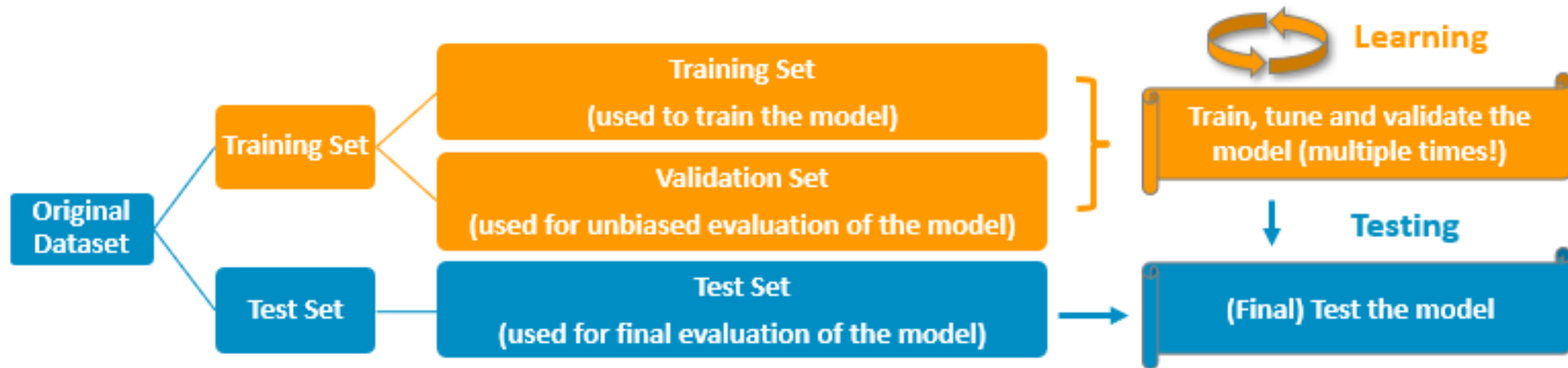
# Conjunto de datos



El **conjunto de test** no se utiliza para el aprendizaje (entrenamiento), solo se usa para asegurar que el modelo **generaliza** bien en **nuevos datos “desconocidos”**.



# Conjunto de datos



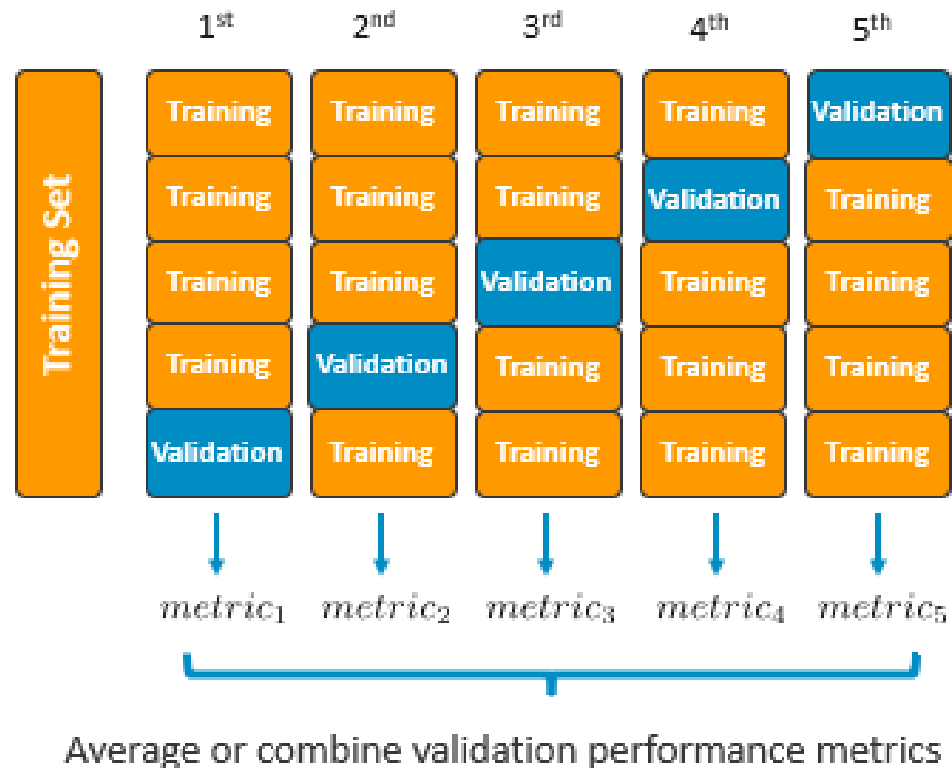
	bad_weather	is_rush_hour	mile_distance	urban_address	late
0	0.0	1.0	5.00	1.0	0.0
1	1.0	0.0	7.00	0.0	1.0
2	0.0	1.0	2.00	1.0	0.0
3	1.0	1.0	4.20	1.0	0.0
4	0.0	0.0	7.80	0.0	1.0
5	1.0	0.0	3.90	1.0	0.0
6	0.0	1.0	4.00	1.0	0.0
7	1.0	1.0	2.00	0.0	0.0
8	0.0	0.0	3.50	0.0	1.0
9	1.0	0.0	2.60	1.0	0.0
10	0.0	0.0	4.10	0.0	1.0

The table shows 11 data points (rows 0-10) with 5 columns: bad\_weather, is\_rush\_hour, mile\_distance, urban\_address, and late. The data is partitioned into three sets: Training Set (rows 0-5), Validation Set (rows 6-8), and Test Set (rows 9-10).

Generalmente se “barajan” los datos antes de hacer las particiones para evitar sesgos en los datos resultantes



# Conjunto de datos

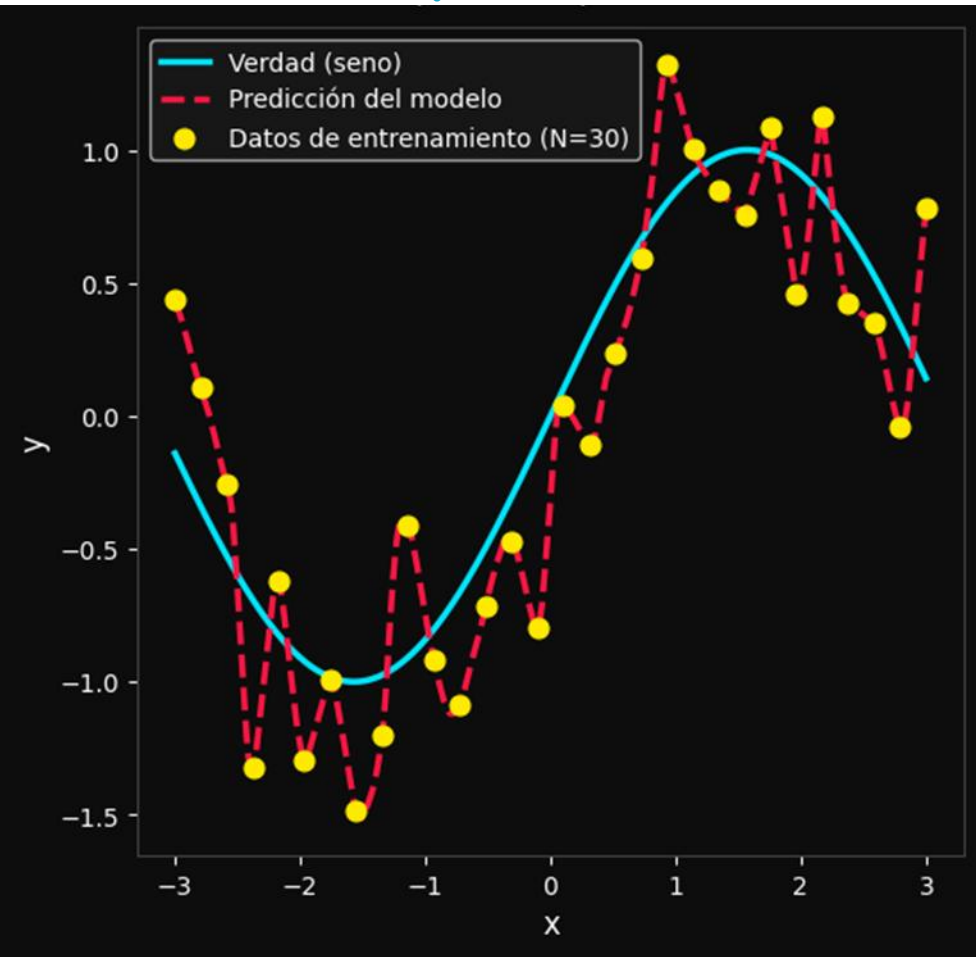
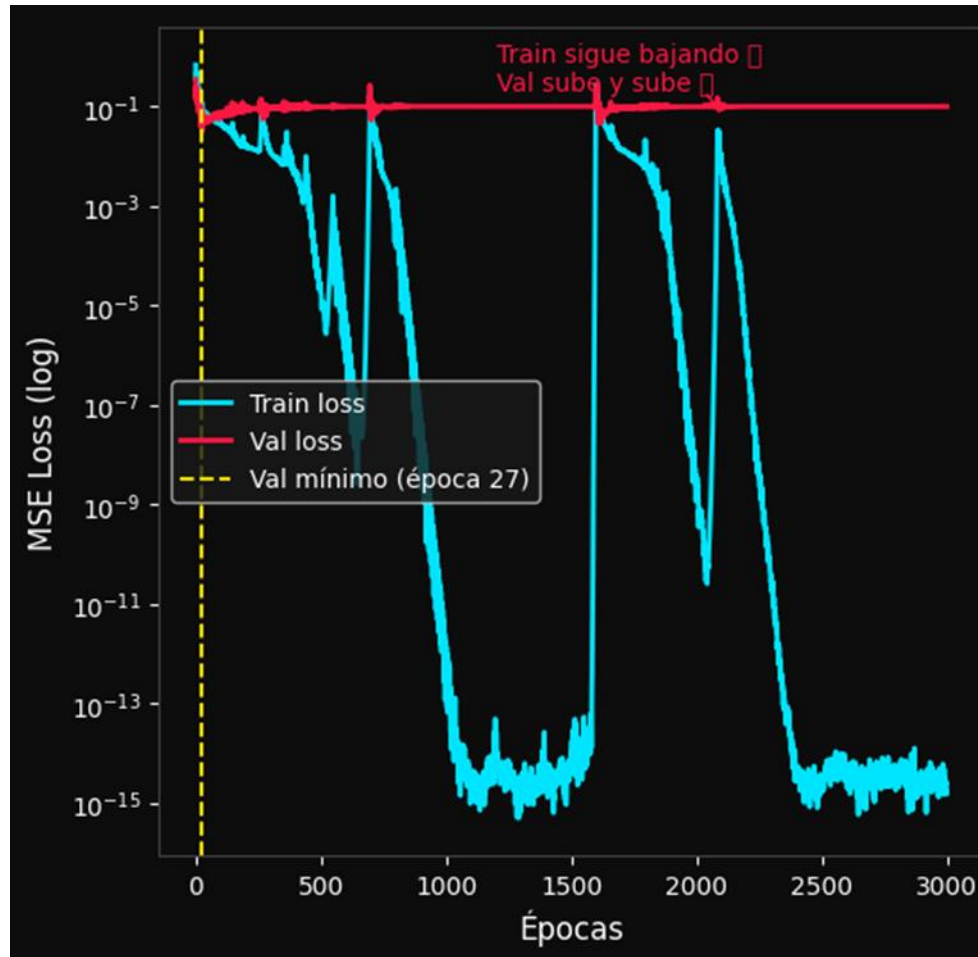


**K-fold cross-validation:** Es una técnica de validación para ver que tan bien generaliza un modelo a un conjunto de validación independiente.

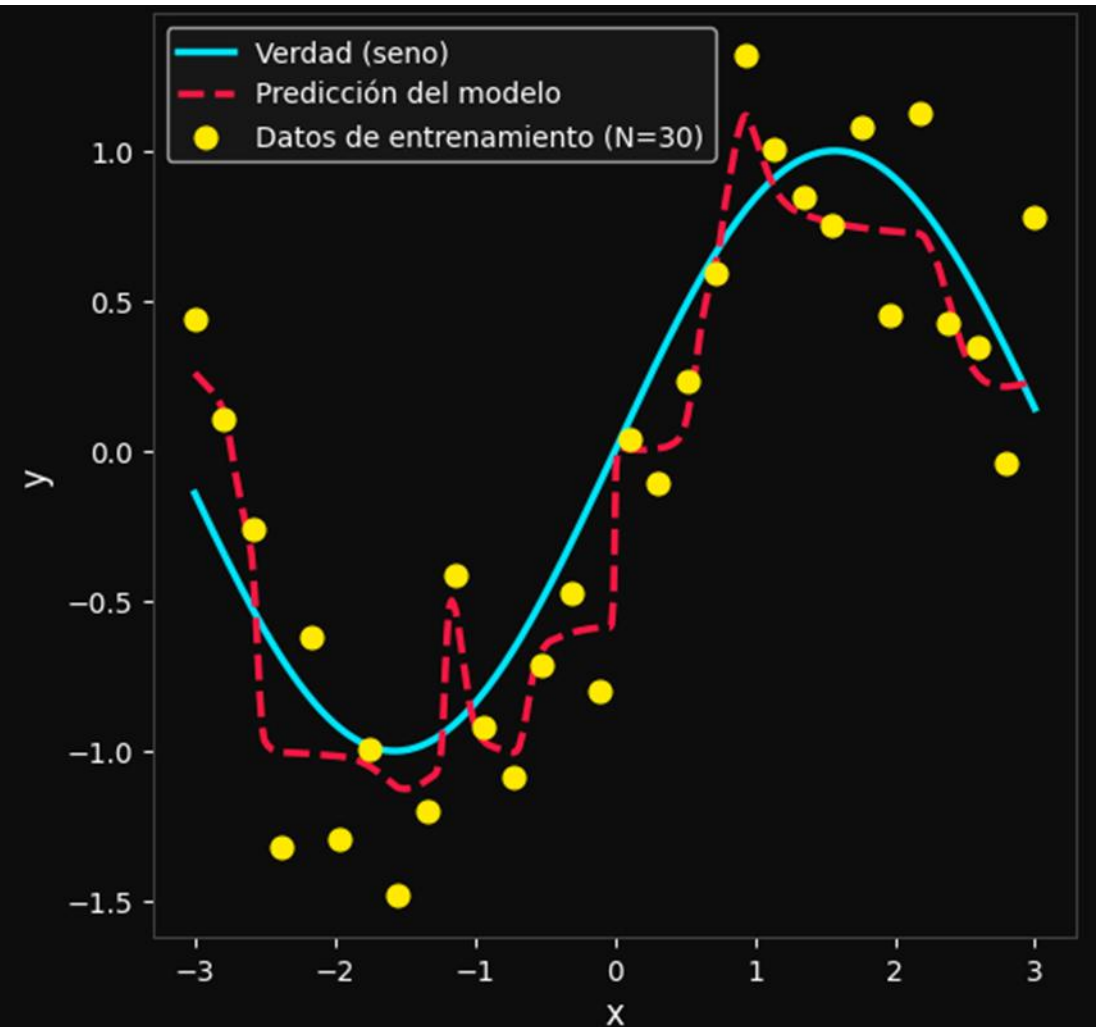
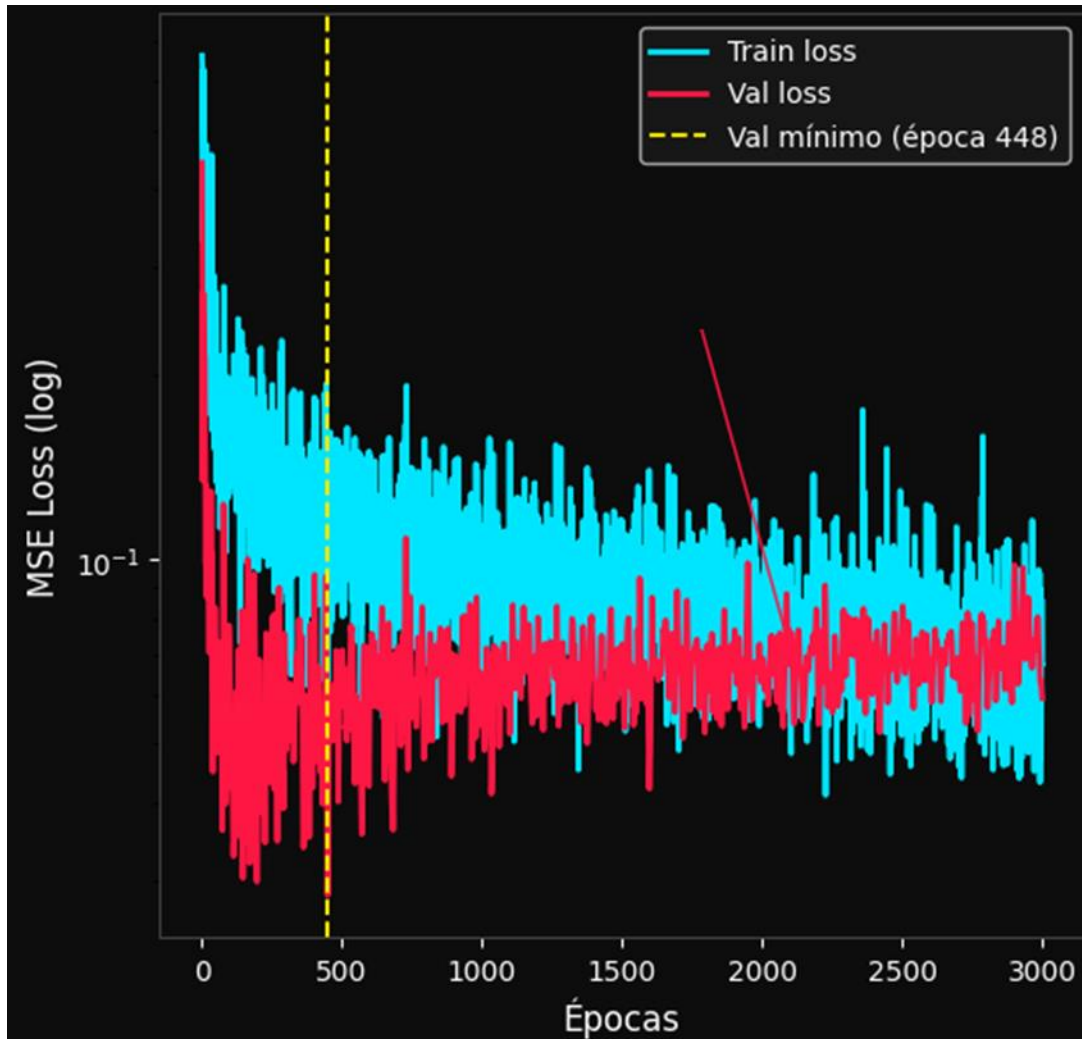
Se utilizan **K muestras reservadas** para validar el modelo, cada vez entrenando con las muestras restantes:

- Dividir el conjunto de entrenamiento en K grupos (folds).
- Repetir K veces el siguiente procedimiento:
  - Reservar el  $K^{th}$  fold para **validación**
  - Entrenar el modelo en los folds restantes
  - Calcular el desempeño en el fold de validación
- Combinar la métrica de rendimiento calculada

# Sobre entrenamiento / Dropout



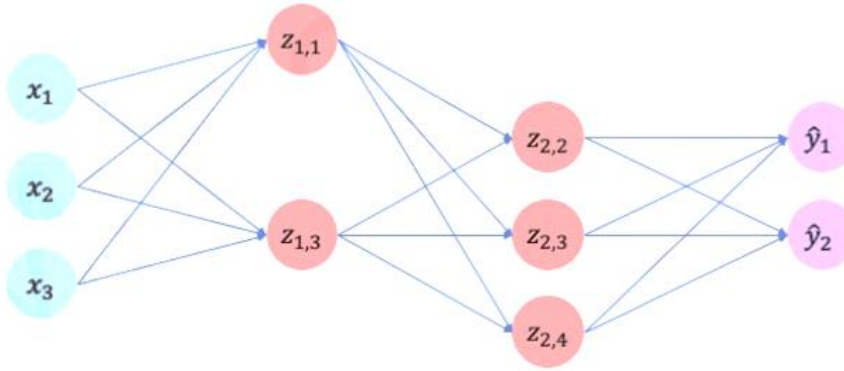
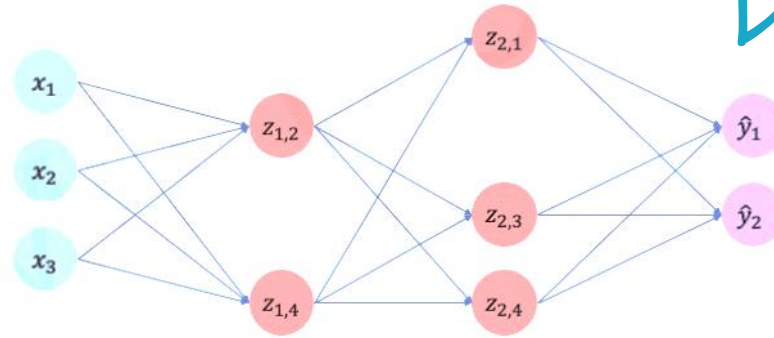
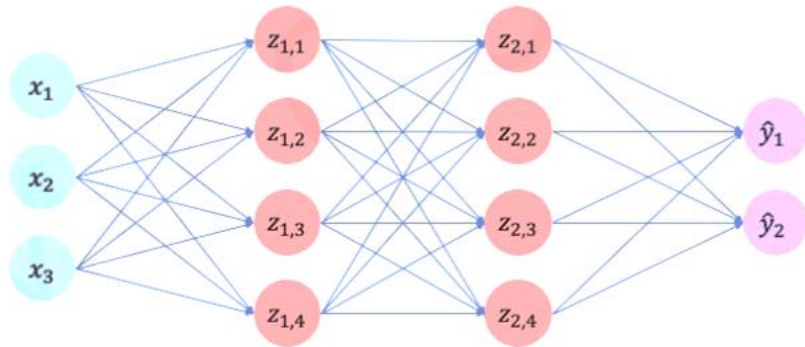
# Sobre entrenamiento / Dropout



# Sobre entrenamiento / Dropout

```
class RedConDropout(nn.Module):  
    def __init__(self):  
        super().__init__()  
        self.net = nn.Sequential(  
            nn.Linear(1, 512), nn.ReLU(), nn.Dropout(0.4),  
            nn.Linear(512, 512), nn.ReLU(), nn.Dropout(0.4),  
            nn.Linear(512, 512), nn.ReLU(), nn.Dropout(0.4),  
            nn.Linear(512, 512), nn.ReLU(), nn.Dropout(0.4),  
            nn.Linear(512, 512), nn.ReLU(), nn.Dropout(0.4),  
            nn.Linear(512, 1),  
        )  
  
    def forward(self, x):  
        return self.net(x)
```

# Sobre entrenamiento / Dropout



Durante el entrenamiento, poner aleatoriamente algunas activaciones en 0.

- Normalmente se "eliminan" el 50% de las activaciones en la capa
- Obliga a la red a no depender de un solo nodo

Artículos de investigaciones originales:

Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. arXiv preprint arXiv:1207.0580.

Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. The Journal of Machine Learning Research, 15(1), 1929-1958.

# Sobre entrenamiento / Dropout

¿Cómo se eliminan nodos eficientemente?

Muestreo de Bernoulli (durante el entrenamiento):

- $p :=$  probabilidad de eliminación
- $v :=$  muestra aleatoria de una distribución uniforme en el rango  $[0, 1]$
- $\forall i \in v : u_i := 0$  si  $u_i < p$  si no 1
- $a := a \odot v$  ( $p \times 100\%$  de las activaciones  $a$  será 0)

Después del entrenamiento, durante la “inferencia”, se deben escalar las activaciones a través de:  $a := a \odot (1 - p)$

¿Por qué es necesario?

## Parámetros vs Hiperparámetros

### Parámetros

- *weights (weight parameters)*
- *biases (bias units)*

### Hiperparámetros

- *minibatch size*
- *data normalization schemes*
- *number of epochs*
- *number of hidden layers*
- *number of hidden units*
- *learning rates*
- *(random seed, **why?**)*
- *loss function*
- *various weights (weighting terms)*
- *activation function types*
- *regularization schemes*
- *weight initialization schemes*
- *optimization algorithm type*

(En su mayoría sin explicación científica, principalmente ingeniería; se necesita probar muchas cosas → *graduate student descent*)