

# Intro to Deep Learning

## Big Data y Machine Learning para Economía Aplicada

Ignacio Sarmiento-Barbieri

Universidad de los Andes

# Deep Learning: Intro

- ▶ Linear Models may miss the nonlinearities that best approximate  $f^*(x)$
- ▶ Neural networks are simple models.
- ▶ The model has **linear combinations** of inputs that are passed through **nonlinear activation functions** called nodes (or, in reference to the human brain, neurons).

# Agenda

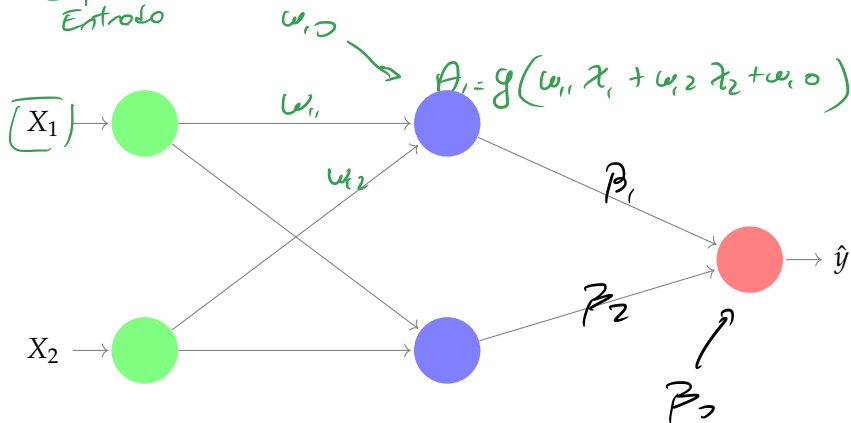
- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

# Single Layer Neural Networks

2 INPUTS, Variables  
Capa de Entrada

Capa Oculta

Capa de Saída



# Single Layer Neural Networks

- ▶ NN are made of **linear combinations** of inputs that are passed through **nonlinear activation functions**

- ▶ The NN model has the form

$$f(X) = f \left[ \beta_0 + \sum_{k=1}^K \beta_k A_k \right]$$

$$= f \left[ \beta_0 + \sum_{k=1}^K \beta_k g \left( \underbrace{w_{k0} + \sum_{j=1}^p w_{kj} X_j}_{\text{Hiper parâmetros}} \right) \right] \quad (2)$$

→ # Nodos o Neuronas

→ # g función de Activación

→ # capa de salida (1)

- ▶ where

- ▶  $g(\cdot)$  is a activation function, the nonlinearity of  $g(\cdot)$  is key
- ▶  $f$  is the output layer of the network

- ▶ both are prespecified

H.w.  $g(z) = z$

# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

## Worked Example: The "Exclusive OR (XOR)" Function

- ▶ The exclusive disjunction of a pair of propositions,  $(p, q)$ , is supposed to mean that  $p$  is true or  $q$  is true, but not both
- ▶ It's truth table is:

<u>q</u>	<u>p</u>	<u>q v p</u>
<u>0</u>	<u>0</u>	<u>0</u>
<u>0</u>	<u>1</u>	<u>1</u>
<u>1</u>	<u>0</u>	<u>1</u>
<u>1</u>	<u>1</u>	<u>0</u>

- ▶ When exactly one of these binary values is equal to 1, the XOR function returns 1. Otherwise, it returns 0

# Worked Example: The "Exclusive OR (XOR)" Function

- Let's use a linear model

$$X = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$

$$\underline{\beta_0} = \frac{1}{2}$$

$$\beta_1 = \beta_2 = 0$$

$$y = \underbrace{\beta_0 + \beta_1 q + \beta_2 p}_{f(q,p)} + u$$

$$y = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

$$\hat{y} = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$$

$$\frac{1}{N} \sum (y - x\beta)^2$$

$$MSE = 1$$



# Worked Example: The "Exclusive OR (XOR)" Function

- ▶ Let's use Single Layer NN containing two hidden units
- ▶ Activation Function: ReLU:  $g(z) = \max\{0, z\}$
- ▶ NN



$$f(X) = \beta_0 + \sum_{k=1}^2 \beta_k g \left( w_{k0} + \sum_{j=1}^{\textcircled{2}} w_{kj} X_j \right) \quad (4)$$

# Worked Example: The "Exclusive OR (XOR)" Function

- Suppose this is the solution to the XOR problem

$$f(x) = \max\{0, xW + W_0\} \beta + \beta_0$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$W_0 = \begin{pmatrix} 0 & -1 \end{pmatrix}$$

$$\beta = \begin{pmatrix} 1 & -2 \end{pmatrix}$$

$$\beta_0 = 0$$

# Worked Example: The "Exclusive OR (XOR)" Function

Lets work out the example step by step

$$\begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 \\ -2 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix}$$

$$\boxed{MSE = 0}$$

# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

# Neural Networks: Activation Functions

HW: como es  
la derivada  
de la logistica.

►  $\text{Sigmoid}(x) = \frac{1}{1+\exp(-x)}$

Logistica  $\mathbb{R} \rightarrow [0, 1]$



►  $\text{ReLU}(x) = \max\{x, 0\}$

$\mathbb{R} \rightarrow [0, \infty)$



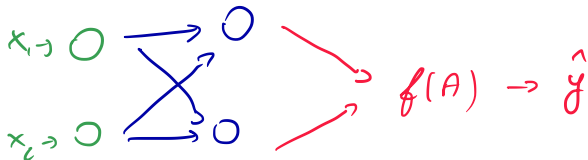
► Among others (see more here)

► Hidden unit design remains an active area of research, and many useful hidden unit types remain to be discovered

# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

# Output Functions



- The choice of output unit is related to the problem at hand

- Regression

$$z = \beta_0 + \sum_{k=1}^K \beta_k A_k$$

- Classification

- Binary

→ logistic

$$\mathbb{R} \rightarrow [0, 1]?$$

$$\frac{1}{1+e^{-z}} \in [0, 1]$$

- Multiclass

→ logistic multinomial  
(softmax) ] →

# Agenda

## 1 Recap: SNN

- Example: XOR
- Activation Functions
- Output Functions

## 2 Training the network

## 3 Architecture Design

- Deep Neural Networks

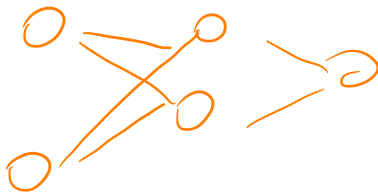
## 4 When to Use Deep Learning?

- # Nodes (Final)

- Activation ✓

- Saliency ✓

$W$ ,  $B$  ✓





# Training the network

- El objetivo es

$$\hat{f} = \operatorname{argmin} \left\{ \sum_{i=1}^n L(y, f(X; \Theta)) \right\} \quad (5)$$

*mse*

*f* → *w, β*

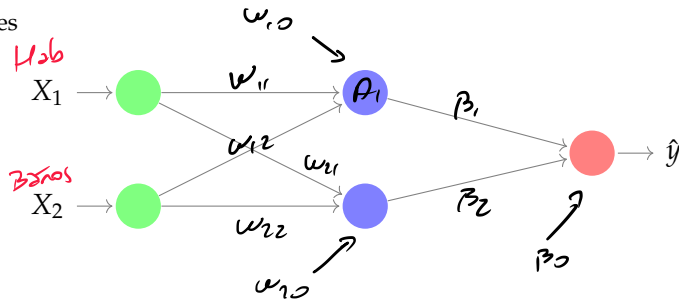
- SNN

$$f(X, \beta, w) = f \left[ \beta_0 + \sum_{k=1}^K \beta_k g \left( w_{k0} + \sum_{j=1}^p w_{kj} X_j \right) \right] \quad (6)$$

*β<sub>0</sub>, β<sub>1</sub>, β<sub>2</sub>*  
*w<sub>k</sub>*

# Training the network

Example: House Prices



## Equations

- Hidden Layer sigmoid (logistic):

- $A_1 = \sigma(w_{11} \cdot X_1 + w_{12} \cdot X_2 + w_{10})$

- $A_2 = \sigma(w_{21} \cdot X_1 + w_{22} \cdot X_2 + w_{20})$

3 param

3 param.

- Output Layer, identity output function:

- $\hat{y}_i = \beta_0 + \beta_1 \cdot A_1 + \beta_2 \cdot A_2$

→ 3 param

9 parameters

- Loss Function  $\Rightarrow$  MSE:  $\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2$

# Training the network

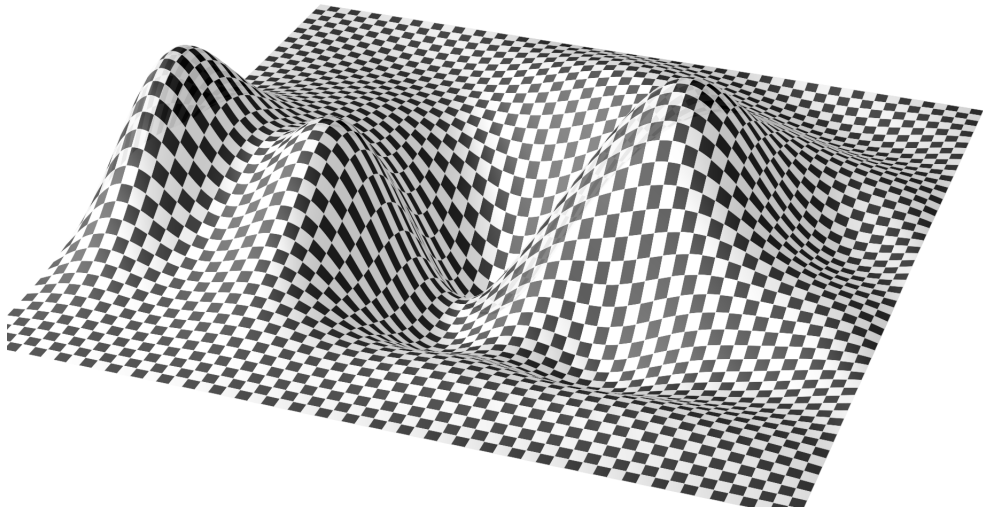
- El objetivo es

$$\hat{f} = \underset{w, \beta}{\operatorname{argmin}} \left\{ \sum_{i=1}^n \ell(y_i, f[\beta_0 + \sum_{k=1}^K \beta_k g(w_{k0} + \sum_{j=1}^J w_{kj} X_j)]) \right\} \quad (7)$$

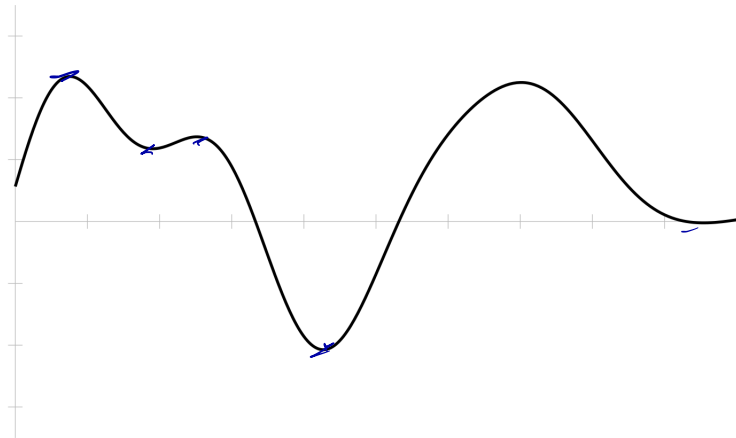
$$w, \beta \Rightarrow \frac{1}{N} \sum_{i=1}^n \left( y_i - \beta_0 - \sum_{k=1}^K \beta_k g\left(w_{k0} + \sum_{j=1}^J w_{kj} x_{ij}\right) \right)^2$$



# Training the network



# Training the network

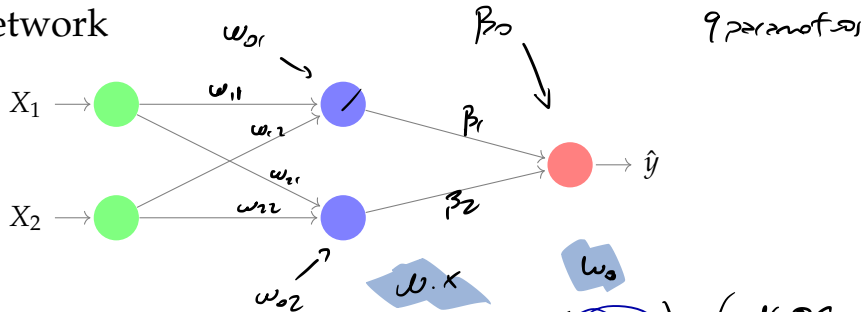


# Training the network

Example: House Prices



# Training the network

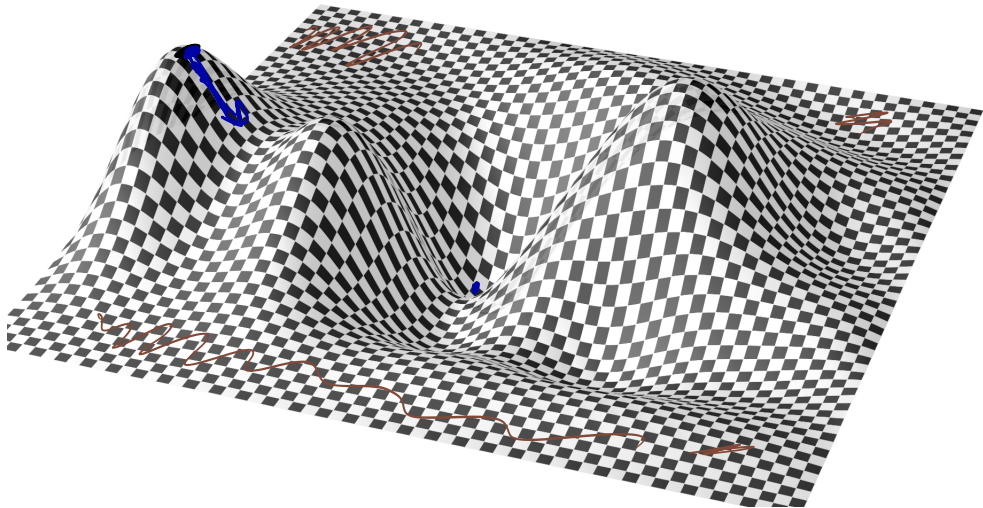


$$W = \begin{pmatrix} \textcircled{2,17} & \textcircled{2,57} \\ \textcircled{\cancel{2,89}} & \textcircled{4,34} \end{pmatrix} \begin{pmatrix} \textcolor{red}{1} \\ \textcolor{red}{3} \end{pmatrix} = \begin{pmatrix} 9,98 \\ 15,88 \end{pmatrix} + \begin{pmatrix} \textcircled{1,87} \\ \textcircled{\cancel{2,09}} \end{pmatrix} = \begin{pmatrix} 11,85 \\ \underline{11,97} \end{pmatrix}$$

$X \rightarrow$  previous obs

$$\text{sig} (W \cdot X + w_0) = \begin{pmatrix} 0,98 \\ 1,99 \end{pmatrix} \cdot \begin{pmatrix} \textcircled{28,33} \\ \textcircled{60,17} \end{pmatrix} = \textcircled{+59} = \frac{\textcolor{red}{142,905}}{\textcolor{blue}{215}}$$

# Training the network





# Training the network

## Backpropagation

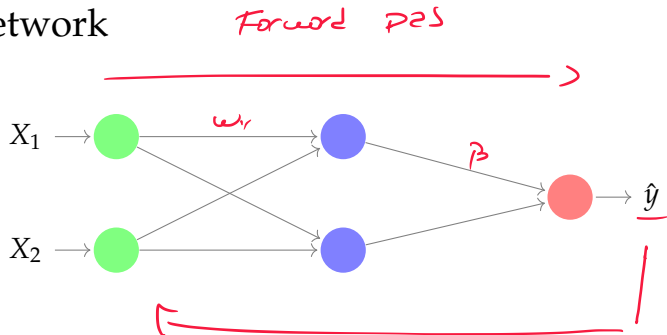
$$\theta' = \theta \ominus \underbrace{\epsilon \cdot \nabla_{\theta} \mathcal{L}(\theta)}$$



Funcion de Perda  
depende de copas

# Training the network

Example: House Prices



## Equations

- Hidden Layer sigmoid (logistic):

- $A_1 = \sigma(w_{11} X_1 + w_{12} X_2 + w_{10})$

- $A_2 = \sigma(w_{21} X_1 + w_{22} X_2 + w_{20})$

- Output Layer, identity output function:

- $\hat{y}_i = \beta_0 + \beta_1 A_1 + \beta_2 A_2$

- Loss Function  $\Rightarrow$  MSE:  $\frac{1}{n} \sum_i^n (y_i - \hat{y}_i)^2 = \mathcal{J}$

# How does backpropagation work?

Updating a single weight

- For simplicity let's focus on updating  $w_{11}$

$$\frac{\partial A_i}{\partial w_{11}} = \sigma'(\cdot) \cdot x_i$$
$$= \sigma(\cdot) \cdot [1 - \sigma(\cdot)] \cdot x_i$$

$$w_{11}' = w_{11} - \epsilon \frac{\partial \mathcal{L}}{\partial w_{11}}$$

$$\frac{\partial \mathcal{L}}{\partial w_{11}} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial A_1} \cdot \frac{\partial A_1}{\partial w_{11}} \cdot \frac{\partial}{\partial}$$

$$\frac{\partial \hat{y}}{\partial A_1} = \hat{\beta}_1$$

H.W.  
derivada  
de la  
relu. ?

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

215-147

# Training the network

Updating a weight in the output layer

- Now let's update one of the weights from the hidden layer to the output layer,  $\beta_1$

$$\beta_1' = \beta_1 - \epsilon \frac{\partial \mathcal{L}}{\partial \beta_1}$$

$$\frac{\partial \mathcal{L}}{\partial \beta_1} = \frac{\partial \mathcal{L}}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \beta_1} = -2(y - \hat{y}) A_1$$

# Training the network

## Batch Gradient Descent

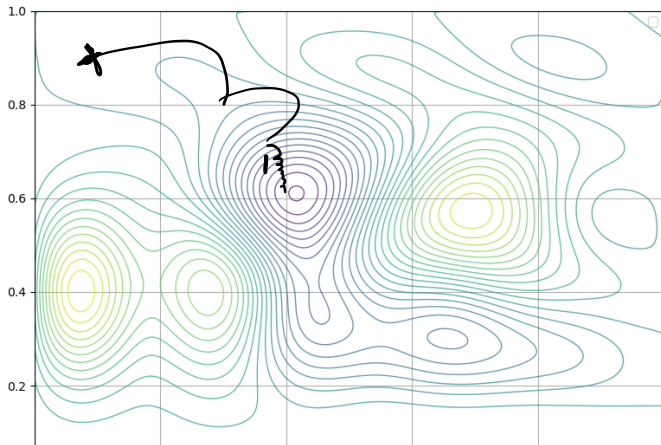
- ▶ Notice that this formula involves calculations over the full data set, at each Gradient Descent step!
- ▶ This is why the algorithm is called **Batch Gradient Descent**: it uses the whole batch of training data at every step.
- ▶ As a result it is terribly slow on very large data sets.

# Training the network

## Stochastic Gradient-Based Optimization

- Stochastic Gradient Descent just picks a random observation at every step and computes the gradients based only on that single observation.

$\epsilon$   
Gradient



# Training the network

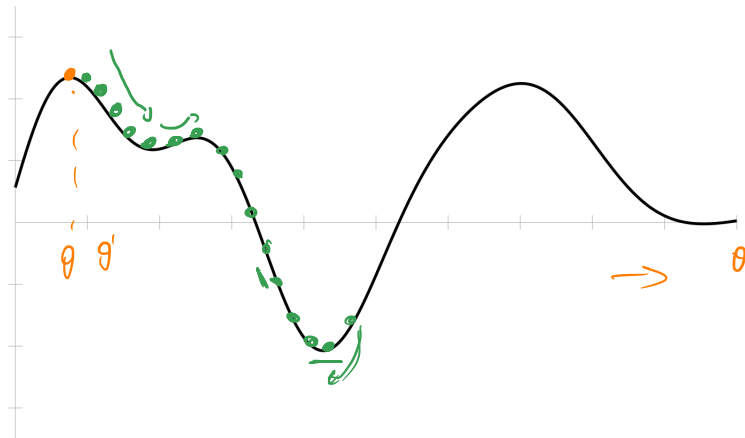
## Mini-batch Gradient Descent

- ▶ Batch Gradient Descent involves calculations over the full data set
- ▶ Stochastic Gradient Descent just picks one random observation at every step
- ▶ At each step, mini-batch GD computes the gradients on small random sets of observations called mini-batches.

# Training the network

$$\theta' = \theta - \underbrace{(\epsilon)}_{(-)} \underbrace{\nabla_{\theta} \mathcal{L}(\theta)}_{(-)}$$

sgd (mini batch)





# Training the network

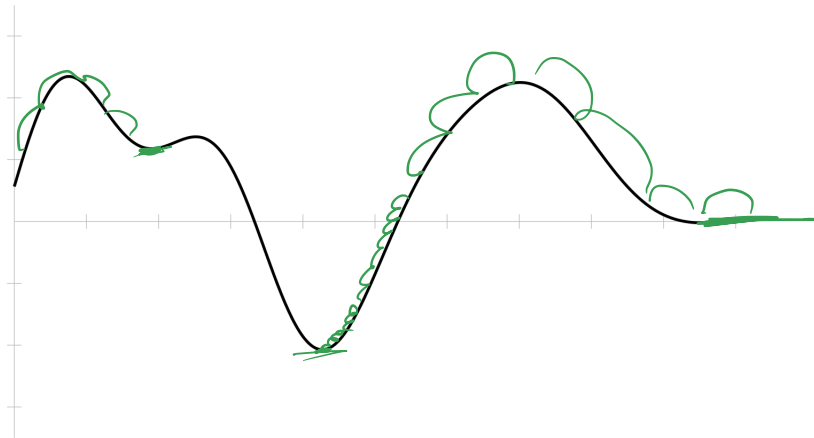
SGD + Momentum, Polyak, 1964

$$v' = \gamma v - \epsilon \cdot \nabla_{\theta} \mathcal{L}(\theta)$$
$$\theta' = \theta + v'$$

Handwritten green annotations: an arrow points from  $Q^a$  to the  $\gamma$  term, and the entire first equation is circled.

- Agrega una "inercia" al gradiente.
- Permite acumular dirección → suaviza la trayectoria.
- Ayuda a escapar de mínimos poco profundos.

# Training the network



# Training the network

RMSProp, Tieleman and Hinton, 2012

$$\frac{\epsilon}{\lfloor \rfloor \uparrow \downarrow}$$

$$E[\nabla_{\theta} \mathcal{L}(\theta)^2] = \rho E[\nabla_{\theta} \mathcal{L}(\theta)^2]_{iter-1} + (1 - \rho) \nabla_{\theta} \mathcal{L}(\theta)^2$$
$$\theta' = \theta - \frac{\epsilon}{\sqrt{E[\nabla_{\theta} \mathcal{L}(\theta)^2] + \eta}} \cdot \nabla_{\theta} \mathcal{L}(\theta)_t$$

$\downarrow$   
constante

# Training the network

Adam, Kingma and Ba, 2014

$$\begin{aligned} m &= \beta_1 m_{iter-1} + (1 - \beta_1) \nabla_{\theta} \mathcal{L}(\theta) \\ v &= \beta_2 v_{iter-1} + (1 - \beta_2) \nabla_{\theta} \mathcal{L}(\theta)^2 \\ \hat{m} &= \frac{m}{1 - \beta_1^{iter}}, \quad \hat{v} = \frac{v}{1 - \beta_2^{iter}} \\ \theta' &= \theta - \epsilon \cdot \frac{\hat{m}}{\sqrt{\hat{v} + \eta}} \end{aligned}$$

- Combina Momentum + RMSProp.
- Corrige sesgo en los primeros pasos.
- Muy popular por su robustez y rapidez.

# Training the network

Which one to use? Zhou et al., 2020

- ▶ Aunque Adam converge más rápido, SGD suele encontrar soluciones que generalizan mejor.
- ▶ Zhou et al., 2020 modela ambos métodos
- ▶ SGD: Puede
  - ▶ Escapar más fácilmente de mínimos locales.
  - ▶ Alcanzar soluciones que generalizan mejor.
- ▶ Adam: quedar atrapado en mínimos subóptimos.

**Conclusión:** El ruido inherente al SGD actúa como una forma de regularización que favorece la generalización.

# Training the network

Example: MNIST



# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

# Architecture Design

- ▶ A key design consideration for neural networks is determining the architecture.
- ▶ The word architecture refers to the overall structure of the network: how many units it should have and how these units should be connected to each other.
- ▶ The universal approximation theorem (Hornik et al., 1989; Cybenko, 1989) guarantees that regardless of what function we are trying to learn, a sufficiently large MLP will be able to represent this function.



# Architecture Design

- ▶ A key design consideration for neural networks is determining the architecture.
- ▶ The word architecture refers to the overall structure of the network: how many units it should have and how these units should be connected to each other.
- ▶ The universal approximation theorem (Hornik et al., 1989; Cybenko, 1989) guarantees that regardless of what function we are trying to learn, a sufficiently large MLP will be able to represent this function.
- ▶ However, learning can fail for two different reasons.
  - 1 The optimization algorithm used for training may not be able to find the value of the parameters that corresponds to the desired function.
  - 2 The training algorithm might choose the wrong function as a result of overfitting

# Architecture Design

- ▶ Using deeper models can reduce the number of units required to represent the desired function and can reduce the amount of generalization error.
- ▶ The ideal network architecture for a task must be found via experimentation guided by monitoring the validation set error

# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?

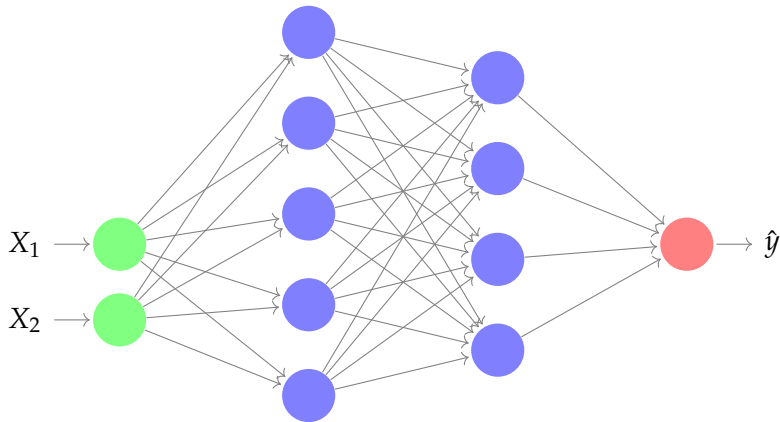
# Architecture Design

## Multilayer Neural Networks

- ▶ Modern neural networks typically have more than one hidden layer, and often many units per layer.
- ▶ In theory a single hidden layer with a large number of units has the ability to approximate most functions.
- ▶ However, the learning task of discovering a good solution is made much easier with multiple layers each of modest size.

# Architecture Design

## Multilayer Neural Networks



# Architecture Design

## Network Tuning

- ▶ Training networks requires a number of choices that all have an effect on the performance:
  - ▶ The number of hidden layers,
  - ▶ The number of units per layer
  - ▶ Details of stochastic gradient descent.
    - SGD
    - RMS Prop
    - Adam
  - ▶ Regularization of parameters
    - Funcion de Activacion
- ▶ This is an active research area that involves a lot of trial and error, and overfitting is a latent danger at each step.
  - Dropout rate
  - L1 → 2
  - L2 → 1

} Hyperparameters  
(tuneable)

# Training the network

Example: MNIST



# Agenda

- 1 Recap: SNN
  - Example: XOR
  - Activation Functions
  - Output Functions
- 2 Training the network
- 3 Architecture Design
  - Deep Neural Networks
- 4 When to Use Deep Learning?



# When to Use Deep Learning?

- ▶ The performance of deep learning usually is very impressive.
- ▶ The question that then begs an answer is: should we discard all our older tools, and use deep learning on every problem with data?

# When to Use Deep Learning?

