

Redes Neuronales y Redes Profundas

Enero 2023

① Introducción

- LAS REDES NEURONALES SON APROX. UNIVERSALES:
UNA RED CON DOS CAPAS, CON LA CAPA DE SALIDA LINEAL, PARA MUCHAS FUNCIONES DE ACTIVACION APROX. CUALQUIERA FUNCION CONTINUA EN UN COMPACTO SIEMPRE Y CUANDO SE TENGAN VARIAS NEURONAS (PAG. 230 DISHOP).

SIN EMBARGO, LAS ACTIVACIONES NO PUEDEN SER
POLINOMIOS

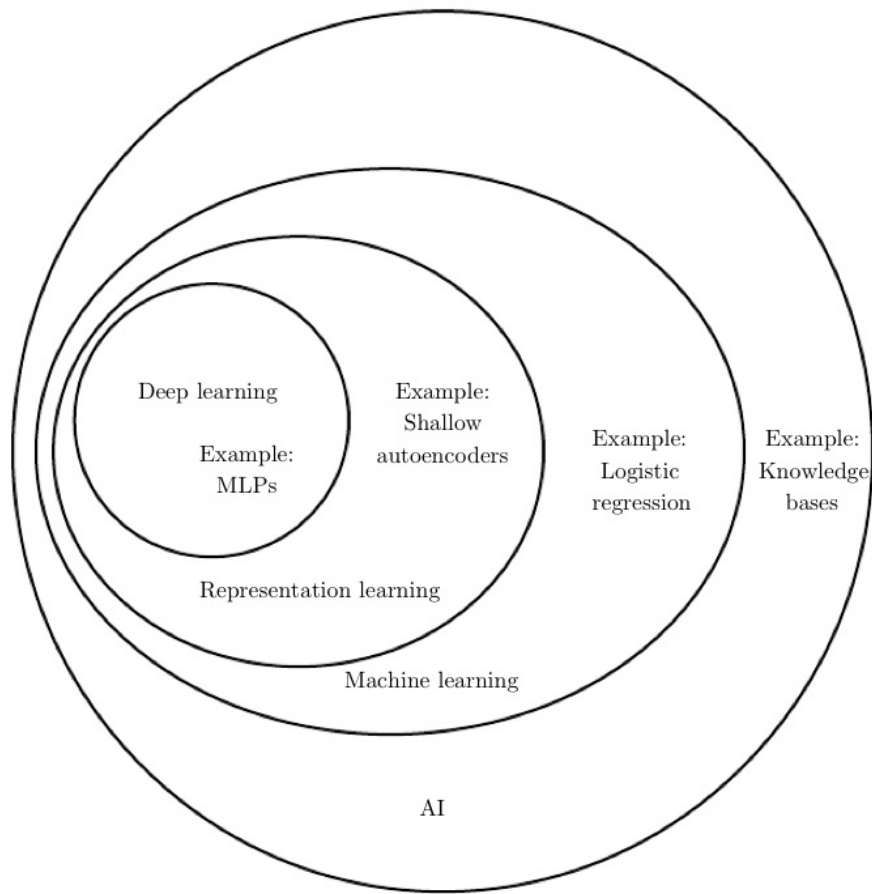


Figure 1.4: A Venn diagram showing how deep learning is a kind of representation learning, which is in turn a kind of machine learning, which is used for many but not all approaches to AI. Each section of the Venn diagram includes an example of an AI technology.

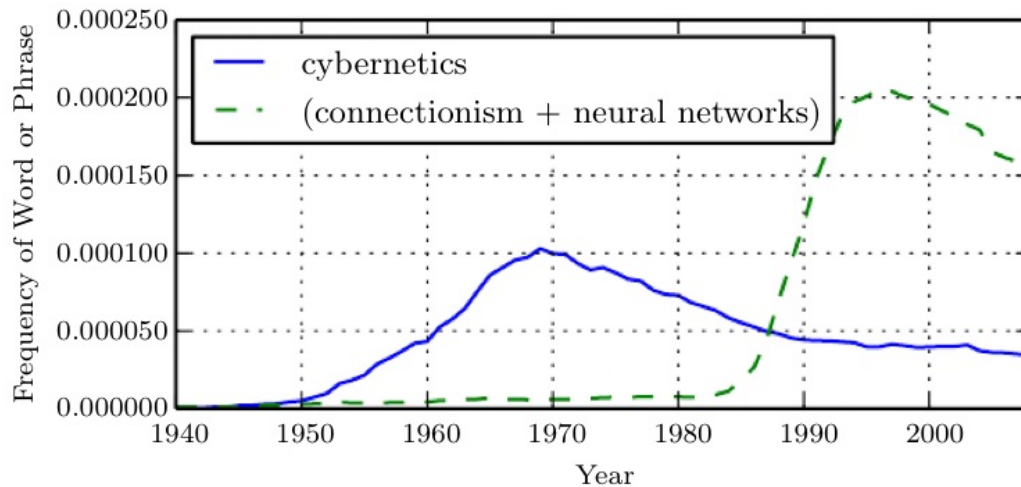


Figure 1.7: The figure shows two of the three historical waves of artificial neural nets research, as measured by the frequency of the phrases “cybernetics” and “connectionism” or “neural networks” according to Google Books (the third wave is too recent to appear). The first wave started with cybernetics in the 1940s–1960s, with the development of theories of biological learning (McCulloch and Pitts, 1943; Hebb, 1949) and implementations of the first models such as the perceptron (Rosenblatt, 1958) allowing the training of a single neuron. The second wave started with the connectionist approach of the 1980–1995 period, with back-propagation (Rumelhart *et al.*, 1986a) to train a neural network with one or two hidden layers. The current and third wave, deep learning, started around 2006 (Hinton *et al.*, 2006; Bengio *et al.*, 2007; Ranzato *et al.*, 2007a), and is just now appearing in book form as of 2016. The other two waves similarly appeared in book form much later than the corresponding scientific activity occurred.

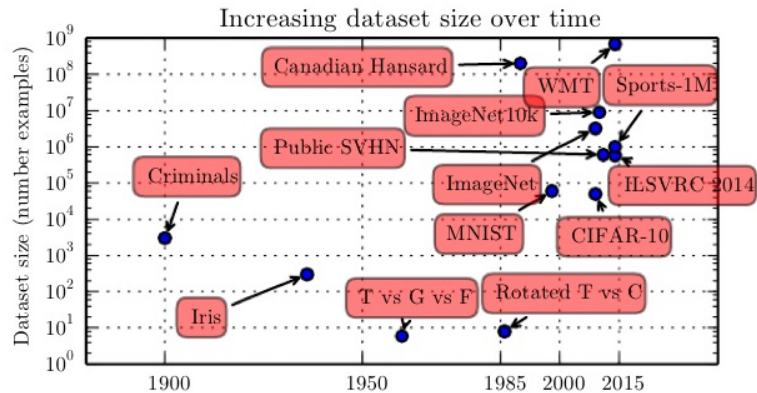


Figure 1.8: Dataset sizes have increased greatly over time. In the early 1900s, statisticians studied datasets using hundreds or thousands of manually compiled measurements (Garson, 1900; Gosset, 1908; Anderson, 1935; Fisher, 1936). In the 1950s through 1980s, the pioneers of biologically inspired machine learning often worked with small, synthetic datasets, such as low-resolution bitmaps of letters, that were designed to incur low computational cost and demonstrate that neural networks were able to learn specific kinds of functions (Widrow and Hoff, 1960; Rumelhart *et al.*, 1986b). In the 1980s and 1990s, machine learning became more statistical in nature and began to leverage larger datasets containing tens of thousands of examples such as the MNIST dataset (shown in Fig. 1.9) of scans of handwritten numbers (LeCun *et al.*, 1998b). In the first decade of the 2000s, more sophisticated datasets of this same size, such as the CIFAR-10 dataset (Krizhevsky and Hinton, 2009) continued to be produced. Toward the end of that decade and throughout the first half of the 2010s, significantly larger datasets, containing hundreds of thousands to tens of millions of examples, completely changed what was possible with deep learning. These datasets included the public Street View House Numbers dataset (Netzer *et al.*, 2011), various versions of the ImageNet dataset (Deng *et al.*, 2009, 2010a; Russakovsky *et al.*, 2014a), and the Sports-1M dataset (Karpathy *et al.*, 2014). At the top of the graph, we see that datasets of translated sentences, such as IBM’s dataset constructed from the Canadian Hansard (Brown *et al.*, 1990) and the WMT 2014 English to French dataset (Schwenk, 2014) are typically far ahead of other dataset sizes.

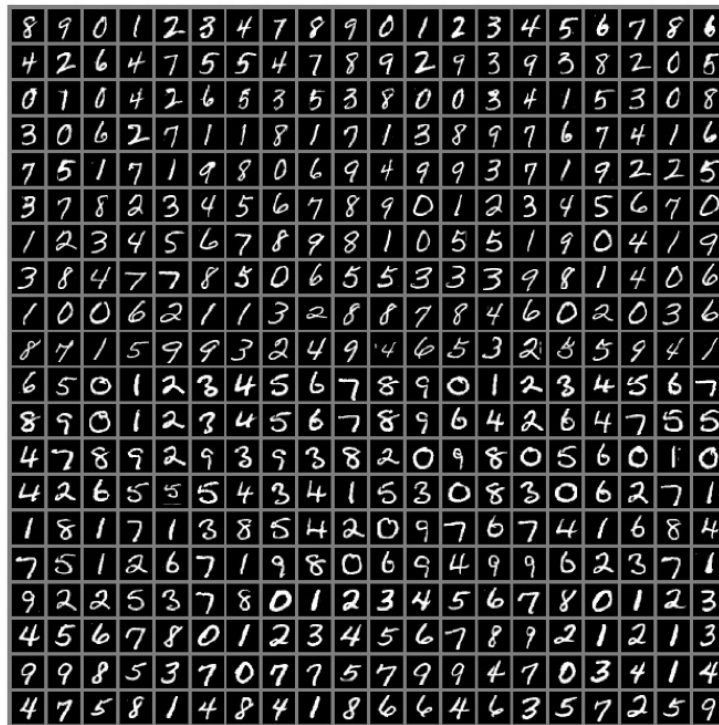


Figure 1.9: Example inputs from the MNIST dataset. The “NIST” stands for National Institute of Standards and Technology, the agency that originally collected this data. The “M” stands for “modified,” since the data has been preprocessed for easier use with machine learning algorithms. The MNIST dataset consists of scans of handwritten digits and associated labels describing which digit 0-9 is contained in each image. This simple classification problem is one of the simplest and most widely used tests in deep learning research. It remains popular despite being quite easy for modern techniques to solve. Geoffrey Hinton has described it as “the *drosophila* of machine learning,” meaning that it allows machine learning researchers to study their algorithms in controlled laboratory conditions, much as biologists often study fruit flies.

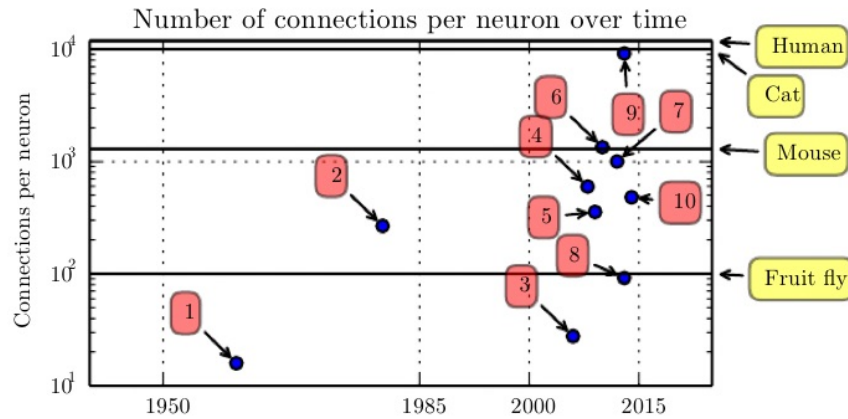


Figure 1.10: Initially, the number of connections between neurons in artificial neural networks was limited by hardware capabilities. Today, the number of connections between neurons is mostly a design consideration. Some artificial neural networks have nearly as many connections per neuron as a cat, and it is quite common for other neural networks to have as many connections per neuron as smaller mammals like mice. Even the human brain does not have an exorbitant amount of connections per neuron. Biological neural network sizes from [Wikipedia \(2015\)](#).

1. Adaptive linear element ([Widrow and Hoff, 1960](#))
2. Neocognitron ([Fukushima, 1980](#))
3. GPU-accelerated convolutional network ([Chellapilla et al., 2006](#))
4. Deep Boltzmann machine ([Salakhutdinov and Hinton, 2009a](#))
5. Unsupervised convolutional network ([Jarrett et al., 2009](#))
6. GPU-accelerated multilayer perceptron ([Ciresan et al., 2010](#))
7. Distributed autoencoder ([Le et al., 2012](#))
8. Multi-GPU convolutional network ([Krizhevsky et al., 2012](#))
9. COTS HPC unsupervised convolutional network ([Coates et al., 2013](#))
10. GoogLeNet ([Szegedy et al., 2014a](#))

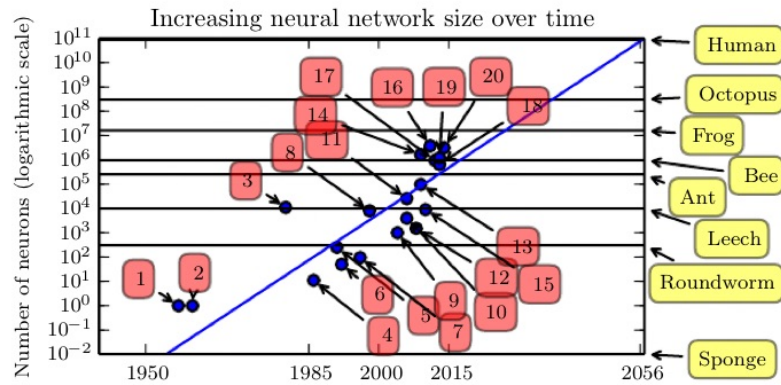


Figure 1.11: Since the introduction of hidden units, artificial neural networks have doubled in size roughly every 2.4 years. Biological neural network sizes from [Wikipedia \(2015\)](#).

1. Perceptron ([Rosenblatt, 1958, 1962](#))
2. Adaptive linear element ([Widrow and Hoff, 1960](#))
3. Neocognitron ([Fukushima, 1980](#))
4. Early back-propagation network ([Rumelhart et al., 1986b](#))
5. Recurrent neural network for speech recognition ([Robinson and Fallside, 1991](#))
6. Multilayer perceptron for speech recognition ([Bengio et al., 1991](#))
7. Mean field sigmoid belief network ([Saul et al., 1996](#))
8. LeNet-5 ([LeCun et al., 1998b](#))
9. Echo state network ([Jaeger and Haas, 2004](#))
10. Deep belief network ([Hinton et al., 2006](#))
11. GPU-accelerated convolutional network ([Chellapilla et al., 2006](#))
12. Deep Boltzmann machine ([Salakhutdinov and Hinton, 2009a](#))
13. GPU-accelerated deep belief network ([Raina et al., 2009](#))
14. Unsupervised convolutional network ([Jarrett et al., 2009](#))
15. GPU-accelerated multilayer perceptron ([Ciresan et al., 2010](#))
16. OMP-1 network ([Coates and Ng, 2011](#))
17. Distributed autoencoder ([Le et al., 2012](#))
18. Multi-GPU convolutional network ([Krizhevsky et al., 2012](#))
19. COTS HPC unsupervised convolutional network ([Coates et al., 2013](#))
20. GoogLeNet ([Szegedy et al., 2014a](#))

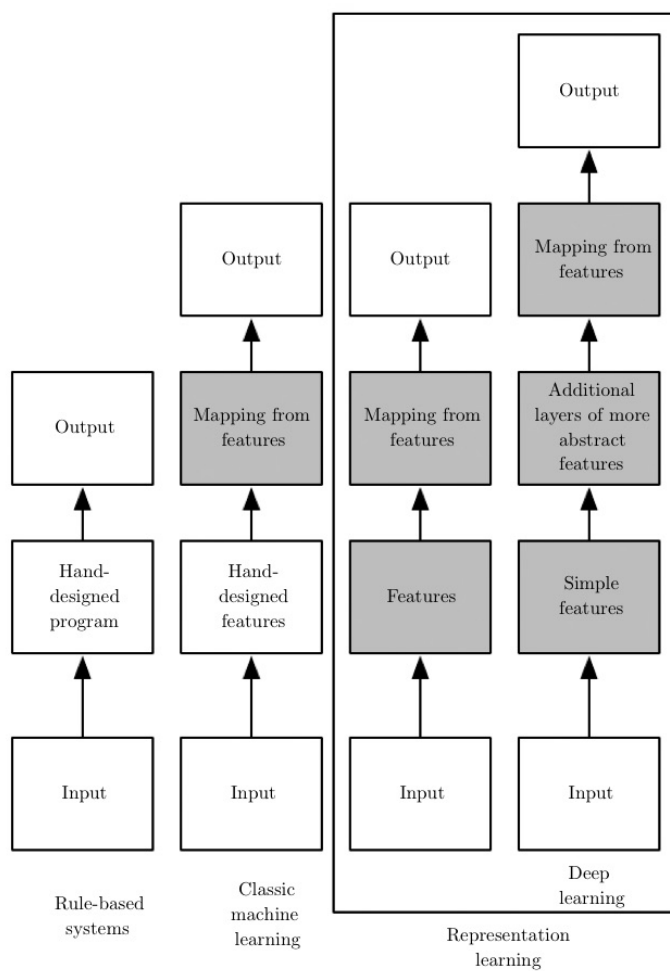


Figure 1.5: Flowcharts showing how the different parts of an AI system relate to each other within different AI disciplines. Shaded boxes indicate components that are able to learn from data.

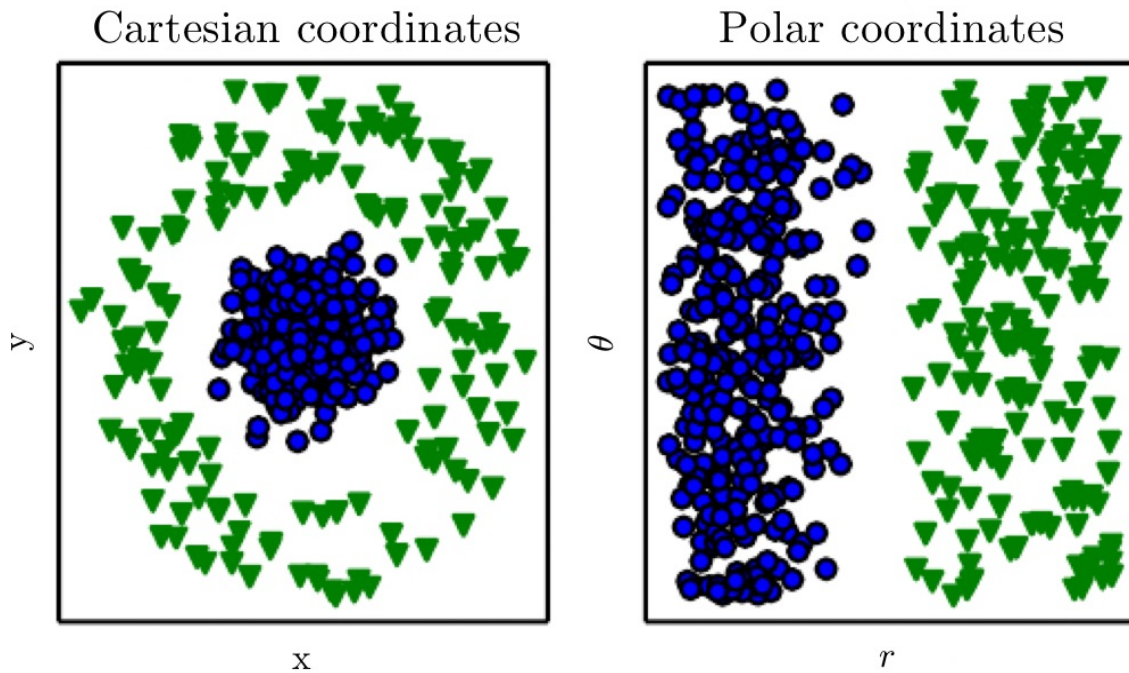


Figure 1.1: Example of different representations: suppose we want to separate two categories of data by drawing a line between them in a scatterplot. In the plot on the left, we represent some data using Cartesian coordinates, and the task is impossible. In the plot on the right, we represent the data with polar coordinates and the task becomes simple to solve with a vertical line. (Figure produced in collaboration with David Warde-Farley)

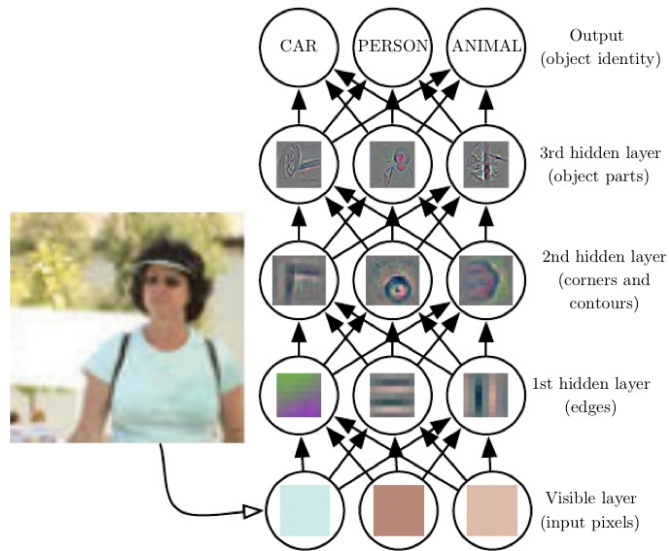


Figure 1.2: Illustration of a deep learning model. It is difficult for a computer to understand the meaning of raw sensory input data, such as this image represented as a collection of pixel values. The function mapping from a set of pixels to an object identity is very complicated. Learning or evaluating this mapping seems insurmountable if tackled directly. Deep learning resolves this difficulty by breaking the desired complicated mapping into a series of nested simple mappings, each described by a different layer of the model. The input is presented at the *visible layer*, so named because it contains the variables that we are able to observe. Then a series of *hidden layers* extracts increasingly abstract features from the image. These layers are called “hidden” because their values are not given in the data; instead the model must determine which concepts are useful for explaining the relationships in the observed data. The images here are visualizations of the kind of feature represented by each hidden unit. Given the pixels, the first layer can easily identify edges, by comparing the brightness of neighboring pixels. Given the first hidden layer’s description of the edges, the second hidden layer can easily search for corners and extended contours, which are recognizable as collections of edges. Given the second hidden layer’s description of the image in terms of corners and contours, the third hidden layer can detect entire parts of specific objects, by finding specific collections of contours and corners. Finally, this description of the image in terms of the object parts it contains can be used to recognize the objects present in the image. Images reproduced with permission from [Zeiler and Fergus \(2014\)](#).

② LAS IDEAS PRINCIPALES

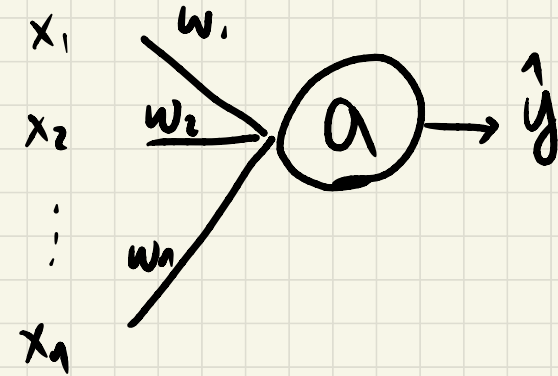
A REDES NEURONALES

B FUNCIONES DE ACTIVACION

C FUNCION DE PERDIDA

D OPTIMIZACION → GRADIENTE DESCENDENTE
↳ "BACK PROPAGATION"

A FUNCIÓN LOGÍSTICA COMO UNA RED



RED UNA
CAPA

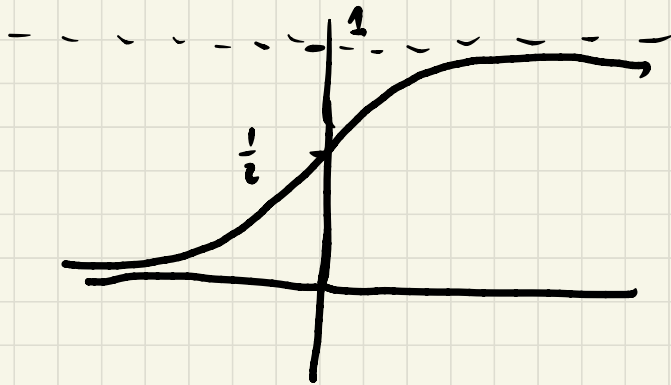
- UNA SOLA CAPA (CAPA SALIDA)
- NO HAY CAPAS OCULTAS

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

$$- \hat{y} = a(z) = \frac{1}{1 + e^{-(w^T x + b)}}$$

B FUNCIÓN DE ACTIVACIÓN

- Función sigmoide: $a(x) = \frac{1}{1 + e^{-x}}$



ε FUNCIÓN DE PÉRDIDA

- ENTROPÍA CRUZADA
- PÉRDIDA CON UN EJEMPLO:

$$\mathcal{L}(\hat{y}^{(i)}, y^{(i)}) = -y^{(i)} \ln(\hat{y}^{(i)}) - (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)})$$

- VARIOS EJEMPLOS:

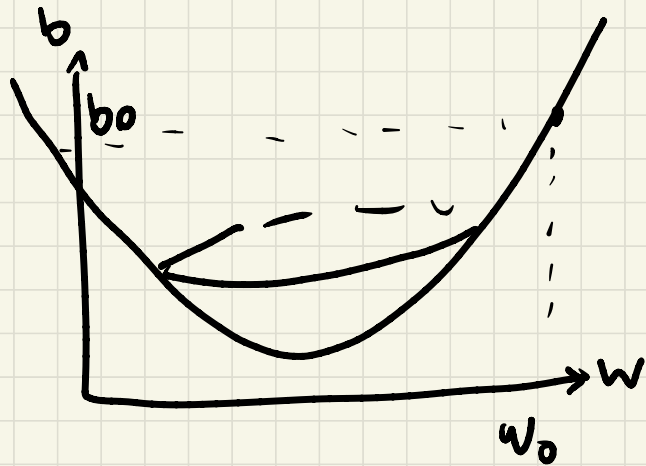
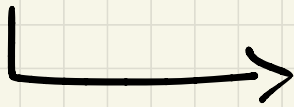
$$\begin{aligned} \mathcal{J}(w, b) &= \frac{1}{m} \sum_{i=1}^m -y^{(i)} \ln(\hat{y}^{(i)}) - (1 - y^{(i)}) \ln(1 - \hat{y}^{(i)}) \\ \text{COSTO} \swarrow & \\ &= \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \end{aligned}$$

0 OPTIMIZACION

- GRADIENTE DESCENDENTE

min
w, b

$J(w, b)$



$$w \leftarrow w - d \underbrace{\nabla J(w_0, b_0)}_{\equiv dw}$$

$$b \leftarrow b - d \underbrace{\nabla J(w_0, b_0)}_{\equiv db}$$

- BATCH LEARNING / ONLINE LEARNING (Bishop 5.2.4)

- BACKPROPAGATION

ES EL PROBLEMA DE CALCULAR:

$\frac{\partial J}{\partial w}$, $\frac{\partial J}{\partial b}$ usando la regla de la cadena

② EJEMPLOS

EJEMPLO FUNCION LOGÍSTICA COMO UNA RED NEURONAL

- DOS VARIABLES
- un solo ejemplo:

$$z = w_1 x_1 + w_2 x_2 + b$$

$$a = \sigma(z)$$

$$L(y, a) = -y \ln(a) - (1-y) \ln(1-a)$$

$$J(w, b) = L(y, a)$$

$$J(w, b) = \sum (y_i - a(z_i))^2, \quad z = w_1 x_1 + w_2 x_2 + b$$

$$\Rightarrow \frac{\partial J(w, b)}{\partial w_1} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial w_1} \equiv \delta w_1$$

$$\frac{\partial y}{\partial a} = \frac{a - y}{a(1-a)}, \quad \frac{\partial a}{\partial z} = a^2 e^{-z} = a(1-a)$$

$$\frac{\partial z}{\partial w_1} = x_1$$

\Rightarrow

$$\frac{\partial J}{\partial b} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b} \equiv \delta b = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot 1$$

en resumen

- con función de activación si modif: g

$$dz \equiv \frac{\partial J}{\partial z} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial z} = \frac{\partial y}{\partial a} \frac{\partial a}{\partial z} = a - y$$

$$dw_i \equiv \frac{\partial J}{\partial w_i} = \frac{\partial J}{\partial a} \cdot \frac{\partial a}{\partial w_i} \cdot \frac{\partial z}{\partial w_i} = (a - y) x_i$$

$$db \equiv \frac{\partial J}{\partial b}$$

- Varios ejemplos

$\equiv dw_i^{(i)}$

$$\frac{\partial J}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial y}{\partial a^{(i)}} \cdot \frac{\partial a^{(i)}}{\partial z^{(i)}} \cdot \frac{\partial z^{(i)}}{\partial w_i}$$

$$\frac{\partial y}{\partial a^{(i)}} = \frac{a^{(i)} - y^{(i)}}{a^{(i)}(1 - a^{(i)})}$$

$$\frac{\partial a^{(i)}}{\partial z^{(i)}} = (a^{(i)})^2 e^{-z^{(i)}} = a^{(i)}(1 - a^{(i)})$$

$$\frac{\partial z^{(i)}}{\partial w_i} = x^{(i)} \quad \frac{\partial y}{\partial b} = \frac{\partial y}{\partial a^{(i)}} \cdot \frac{\partial a^{(i)}}{\partial z^{(i)}}$$

EJEMPLO XOR: APRENDIENDO REPRESENTACIONES

$$x \in \{0,1\} \times \{0,1\}$$

$$f_{\text{XOR}}(x) = \begin{cases} 1 & \text{si } x_1 = 1 \vee x_2 = 1 \\ 0 & \text{CASO CONTRARIO} \end{cases}$$

1	1	0
0	0	1
	0	1

- Si tratamos de aproximar XOR por una función lineal:

$$f_{\text{XOR}}(x) \cong w^T x + b$$

- Si minimizamos el error cuadrático se obtiene $w=0$, $b=1/2$:

$$(f_{\text{XOR}}(0,0) - b)^2 + (f_{\text{XOR}}(0,1) - w_2 - b)^2 + \dots$$

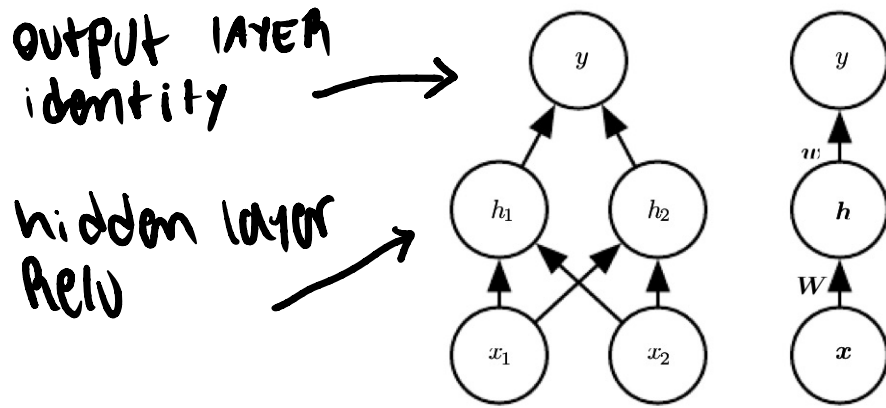


Figure 6.2: An example of a feedforward network, drawn in two different styles. Specifically, this is the feedforward network we use to solve the XOR example. It has a single hidden layer containing two units. *(Left)* In this style, we draw every unit as a node in the graph. This style is very explicit and unambiguous but for networks larger than this example it can consume too much space. *(Right)* In this style, we draw a node in the graph for each entire vector representing a layer's activations. This style is much more compact. Sometimes we annotate the edges in this graph with the name of the parameters that describe the relationship between two layers. Here, we indicate that a matrix \mathbf{W} describes the mapping from \mathbf{x} to \mathbf{h} , and a vector \mathbf{w} describes the mapping from \mathbf{h} to y . We typically omit the intercept parameters associated with each layer when labeling this kind of drawing.

$$- a^{[1]}(x) = \max \left\{ 0, w^{[1]T} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b^{[1]} \\ b^{[1]} \end{bmatrix} \right\} \quad \text{donde}$$

$$w^{[1]} = \begin{bmatrix} w_{11}^{[1]} & w_{12}^{[1]} \\ w_{21}^{[1]} & w_{22}^{[1]} \end{bmatrix}$$

$$- a^{[2]}(x) = w^{[2]T} a^{[1]}(x) + b^{[2]}$$

- Al minimizar el error cuadrático:

$$w^{[1]} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ -1 \end{bmatrix} \quad x \in \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\}$$

$$w^{[2]} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, \quad b^{[2]} = [0] \quad y \in \{ 0, 1, 1, 0 \}$$

- ESTA MED NEURONA APROXIMA PERFECTAMENTE LA
FUNCIÓN XOR:

$$\text{supongamos que } X = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow z^{[1]} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} \Rightarrow a^{[1]} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\Rightarrow z^{[2]} = 1 \Rightarrow a^{[2]} = 1$$

③ NEO DE DOS CAPAS

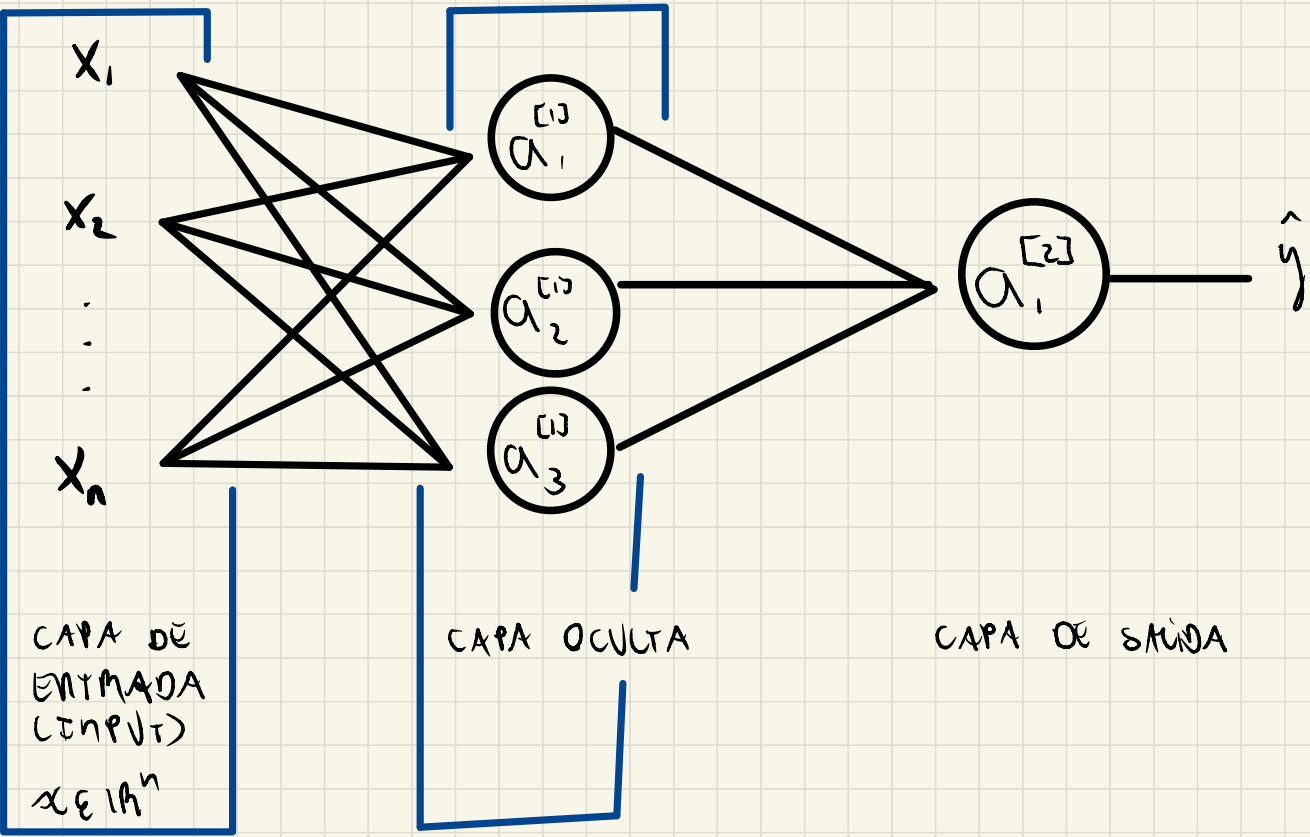
A ARQUITECTURA

B MAS FUNCIONES DE ACTIVACION

C EJEMPLOS DE APOX. FUNCIONES A CLASIFICACION

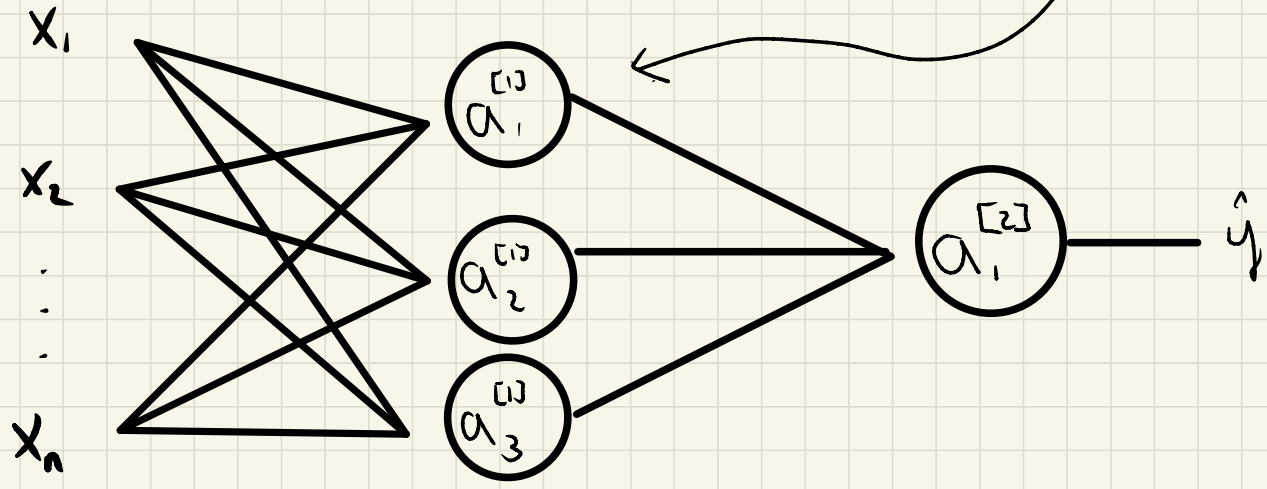
A ARQUITECTURA

DOS CAPAS → ENTRADA
→ UMA CAPA OCULTA



$$z^{[l]} = W^{[l]T} x^{[l]} + b^{[l]}$$

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

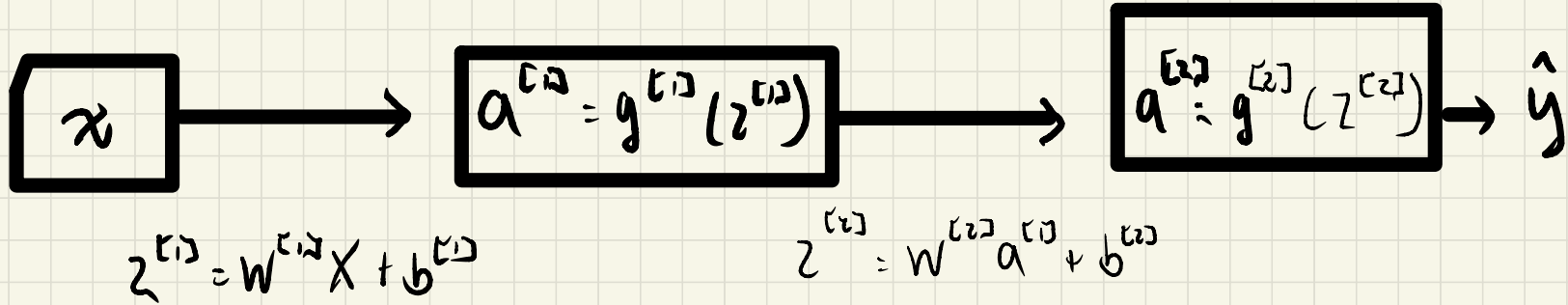


CAPA DE
ENTRADA
(INPUT)
 $x \in \mathbb{R}^n$

CAPA OCULTA

CAPA DE SAIDA

FORWARD PROPAGATION



- FUNCIÓN DE COSTOS UN EJEMPLO:

$$y(w_1, b_1, w_2, b_2) = L(y_i, \hat{y}_i(w_1, b_1, w_2, b_2))$$

- DERIVADAS DE Y CON RESPECTO A w_1, w_2
ES IGUAL QUE EN EL CASO DE LA FUNCIÓN
LOGÍSTICA SOLO QUE AHORA $a^{[i]}$ JUEGA EL PAPEL
DE x .

CASO: OUTPUT BINARIO, PÉRDIDA ENTROPÍA CRUZADA, ACTIVACIÓN SAUDA SIGMOID

CAPA L

$$dz^{[L]} = a^{[L]} - y$$

$$dw^{[L]} = dz^{[L]} \underbrace{a^{[L-1]}}_{\downarrow}$$

en vez de x .

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

$$db^{[L]} = dz^{[L]}$$

CODIGO: EJEMPLO L=2

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

- observarse que una vez se hace el FORWARD PROPAGATION SE PUEDE CALCULAR $dW^{[L]}$, $db^{[L]}$

CAPA L-1: EJEMPLO CON $L=2$, CAPA OCULTA ACTIVACION $a^{[1]} = g^{[1]}(z^{[1]})$

$$y(y, \hat{y}(w_1, b_1, w_2, b_2))$$

$$= y(y, g^{[2]}(w^{[2]} g^{[1]}(w^{[1]} x + b^{[1]} + b^{[2]}))$$

$$\Rightarrow dw^{[1]} = \underbrace{\frac{\partial y}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{[2]}}}_{dz^{[2]} = a^{[2]} - y} \cdot \underbrace{\frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial g^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w^{[1]}}}_{\frac{\partial z^{[2]}}{\partial g^{[1]}} = w_2} \cdot g_1'(z^{[1]}) \cdot x$$

codigo

$$dz^{[1]} = dz^{[2]} w_2 g_1'(z^{[1]})$$

$$dw^{[1]} = dz^{[1]} x$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

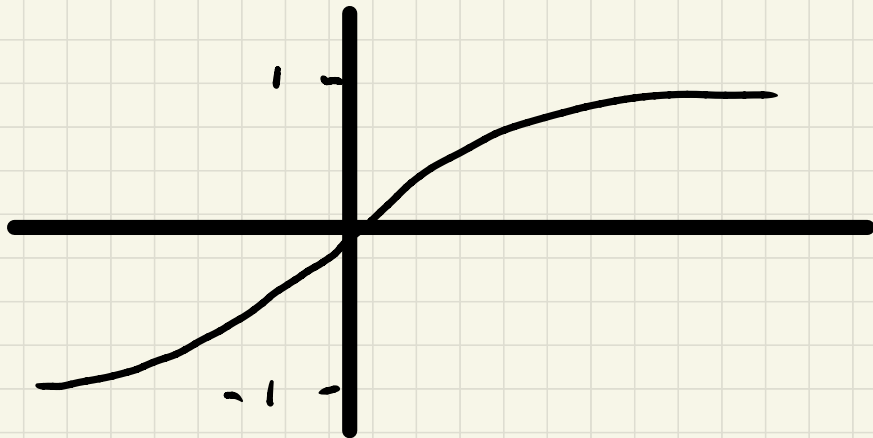
$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

B

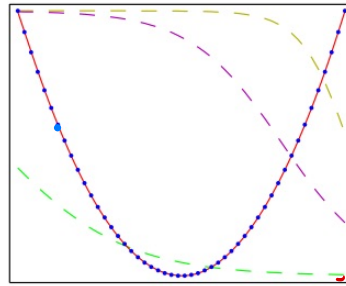
FUNCIONES DE ACTIVACIÓN

- sigmoid (logística)
 - ReLU
 - Tanh
- } común en capas ocultas

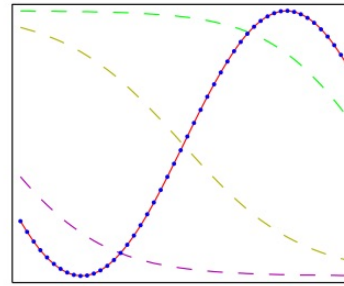


$$\frac{e^z - e^{-z}}{e^z + e^{-z}}$$

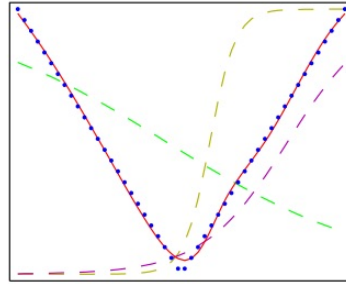
C EJEMPLOS APROXIMACION



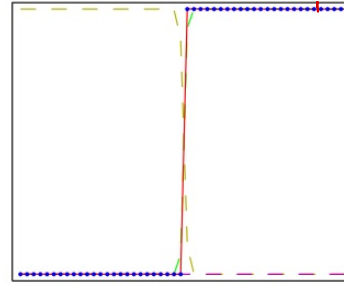
(a)



(b)



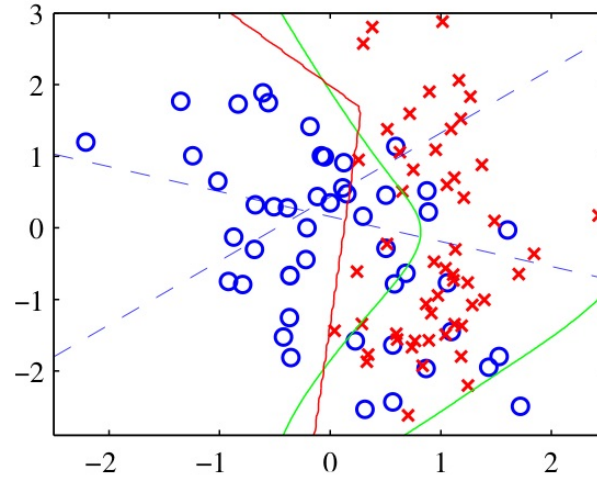
(c)



(d)

- Capacidad de aproximación de una red (e.g., cuadrática, seno, valor absoluto y Heavside). Los datos son los 50 puntos azules. Se entrena una red con dos capas, tres neuronas, función de activación tanh, salidas lineales. Las salidas de las tres neuronas ocultas se muestran con líneas punteadas.

EJEMPLO CLASIFICACION



- Capacidad de aproximación de una red: las líneas punteadas son las salidas de cada una de las dos neuronas (hipersuperficies). Funciones de activación tanh y salida logística sigmoid.
- La línea verde es el clasificador Bayesiano. La roja el clasificador de la red.

④ INTERPRETACION PROBABILISTICA DE LAS SALIDAS DE UNA RED

- nos permite elegir mejor las funciones de ACTIVACION
- DA OTRA INTERPRETACION DE LA FUNCION DE PERDIDA y COSTO

① REGRESION

en un problema de regresion supongamos que $p(y | x, w) = \pi(y | \hat{y}(x, w), \beta^{-1})$ donde β es la varianza inversa (precisión) de la distribución

GAUSSIANA:

$$P(Y | X, \omega) = \frac{1}{(2\pi \sigma)^{1/2}} e^{-\frac{1}{2\sigma^2} (y - \hat{y}(x, \omega))^2}$$

si $\beta = 1/\sigma^2$

$$\Rightarrow P(Y | X, \omega) = \left(\frac{\beta}{2\pi}\right)^{1/2} e^{-\frac{\beta}{2} (y - \hat{y}(x, \omega))^2}$$

suponiendo independencia de y condicional a x , la log verosimilitud es:

$$\Rightarrow \ell(\omega, \beta) = \frac{1}{2} m \ln\left(\frac{\beta}{2\pi}\right) - \frac{\beta}{2} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)}(x^{(i)}, \omega))^2$$

- MAXIMIZAR VEROSIMILITUD ES EQUIVALENTE A MINIMIZAR PÉRDIDA CROMÁTICA

② CLASIFICACIÓN BINARIA

- si la salida es binaria cero o uno y la red genera una probabilidad con una función sigmoide:

$$P(y | x, w) = \hat{y}(x, w)^y (1 - \hat{y}(x, w))^{1-y}$$

- suponiendo independencia de y condicional a x , la log verosimilitud es:

$$J(w) = \frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \ln(\hat{y}(x^{(i)}, w)) + (1 - y^{(i)}) \ln(1 - \hat{y}(x^{(i)}, w)) \right]$$

③ CLASSIFICATION MULTICLASO

$$y \in \{1, \dots, K\}$$

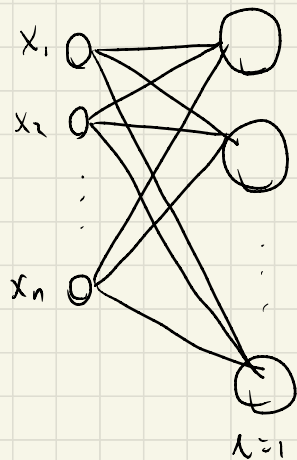
suponiendo la salida de la red es

$p(y_k = 1 | w)$ por ejemplo cuando

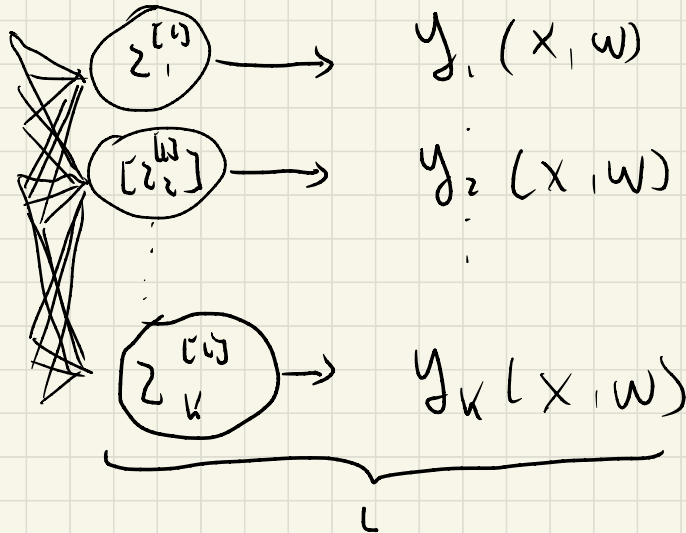
la función ACTIVACION es:

SOFTMAX OUTPUT ACTIVATION :

$$y_k(x, w) = p(y_k = 1 | x) = \frac{e^{z_k^{(L)}}}{\sum_{i=1}^K e^{z_i^{(L)}}}$$



...



Observese que LA ULTIMA CAPA TIENE UN FORMER ESPECIAL.
NO ES LA MISMA FUNCION DE ACTIVACION EN CADA NEURONA.

se puede interpretar como una neurona final que recibe un vector de input $z^{[L]} = (z_1^{[L]}, \dots, z_k^{[L]})$ y la salida es un vector de probabilidades:

$$y = (y_1, \dots, y_k)$$

- en este caso con las mismas hipótesis anteriores el log de la función de verosimilitud es:

$$l(w) = \sum_{i=1}^m \sum_{k=1}^K y_{ik} \ln(\hat{y}_{ik}(x^{(i)}, w))$$

④ TEORÍA DE LA INFORMACIÓN

- Likely events should have low information content, and in the extreme case, events that are guaranteed to happen should have no information content whatsoever.
- Less likely events should have higher information content.
- Independent events should have additive information. For example, finding out that a tossed coin has come up as heads twice should convey twice as much information as finding out that a tossed coin has come up as heads once.

- MOTIVACIÓN PARA DISTRIBUCIONES DISCRETAS: LA INFORMACIÓN de observar x es $I(x) = -\ln p(x)$. en unidades de nats. un nat es la información que tiene un evento con prob. $1/e$. si usamos \log_2 es unidades de bits. un bit es la información de un evento que sucede con prob. $1/2$. ESTA DEFINICIÓN NO SE EXTIENDE a varb. continuas.

- ENTROPIA DE SHANNON (discretas o continuas)

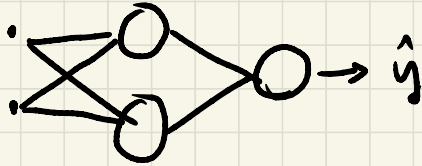
LA ENTROPIA de P es : $H(P) = E_P [-\ln(P(x))]$

- LA DIVERGENCIA KL permite interpretar la maximización de verosimilitud como la minimización de la divergencia de KL entre la verdadera distribución y la distribución empírica.

⑤ INICIALIZACION PARAMETROS

- EN ESTA PARTE DEMOSTRAREMOS PORQUE INICIALIZAR CON CEROS NO ES una buena idea para redes con $L \geq 2$.

- considere esta red:



$$W^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, W^{[2]} = [0 \ 0]$$

$$b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, b^{[2]} = [0]$$

} INICIALIZACION

$$\Rightarrow z^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, a^{[1]} = \begin{bmatrix} 1/2 \\ 1/2 \end{bmatrix}$$

$$z^{[2]} = [0], a^{[2]} = [1/2]$$

} FORMULARIO PROP.

caso logístico
funciona bien
pero no en
general

BACKWARD

$$\begin{aligned} db^{[2]} &= dz^{[2]} = \frac{1}{2} - y & , & \quad dW^2 = \left[\frac{1}{2}(\frac{1}{2} - y) \quad \frac{1}{2}(\frac{1}{2} - y) \right]^T \\ db^{[1]} &= dz^{[1]} = \begin{bmatrix} 0 & 0 \end{bmatrix}^T & , & \quad dW^1 = \begin{bmatrix} 0 & 0 \end{bmatrix} \end{aligned}$$

Observe que :

$$dz_1^{[1]} = dz_2^{[1]}$$

$$dW^1 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \text{---} \\ \text{---} \end{bmatrix} = \begin{bmatrix} \text{N} \rightarrow \\ \text{N} \rightarrow \end{bmatrix} \text{filas iguais}$$

$$W^1 = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} \checkmark \\ \checkmark \end{bmatrix} \rightsquigarrow \boxed{\text{filas iguais}}$$

BACKWARD :

$$dz^{[2]} = a^{[2]} - y$$

$$dW^{[2]} = dz^{[2]} a^{[1]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]})$$

$$dW^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

- PRIMEIRA SEMPLIFICACAO DE PARAMETROS

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} - d \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix} - d \frac{(1/2 - 1)}{2} \begin{bmatrix} 1 & 1 \end{bmatrix}, \quad b^{[2]} = \begin{bmatrix} -\frac{d}{2} (1/2 - 1) \end{bmatrix}$$

- observe-se que :

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$w^{[2]} = \begin{bmatrix} 1 & 1 \end{bmatrix}, \quad b^{[2]} = \begin{bmatrix} 1 \end{bmatrix}$$

- FORWARD PROP (segunda iteración)

$$z^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, a^{[1]} = \begin{bmatrix} \frac{1}{2} \\ \frac{1}{2} \end{bmatrix}$$

$$z^{[2]} = \frac{1}{2} N' + \frac{1}{2} N' = N', a^{[2]} = \frac{1}{1+e^{-N'}}$$

- BACKWARD PROP (segunda iteración)

$$dz_1^{[1]} = dz_1^{[2]}$$

$$dw^{[1]} = \begin{bmatrix} dz_1^{[1]} \\ dz_2^{[1]} \end{bmatrix} [x_1, x_2] = \begin{bmatrix} dz_1^{[1]} x_1 & dz_1^{[1]} x_2 \\ dz_2^{[1]} x_1 & dz_2^{[1]} x_2 \end{bmatrix} = \begin{bmatrix} w'' \\ w'' \end{bmatrix} \text{ filas iguales}$$

- SEGUNDA ACTUALIZACIÓN DE PARÁMETROS

$$w^{[1]} = \begin{bmatrix} - \\ - \end{bmatrix} \text{ filas iguales}$$

- Si continuáramos iterando (más generalmente usando inducción) obtenemos que:

$$W^{[i]} = \begin{bmatrix} r \rightarrow \\ r \rightarrow \end{bmatrix} \rightsquigarrow \boxed{\text{Filas iguales}}$$

Luego:

$$Z^{[i]} = \begin{bmatrix} r \rightarrow \\ r \rightarrow \end{bmatrix} = \begin{bmatrix} r \cdot [x_1, x_2] \\ r \cdot [x_1, x_2] \end{bmatrix} \Rightarrow Z_1^{[i]} = Z_2^{[i]}$$

$\Rightarrow A^{[i]}(Z_1^{[i]}) = A^{[i]}(Z_2^{[i]})$ en todas las iteraciones: las dos unidades están haciendo lo mismo (simétricas). Las unidades no hacen nada

⑥ RED NEURONAL GENERAL

- FORWARD PROPAGATION

$$z^{[k]}(i) = W^{[k]} a^{[k-1]}(i) + b^{[k]}$$

$$a^{[k]}(i) = g(z^{[k]}(i))$$

define:

$$a^{[0]}(i) = x^i$$

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

- BACKWARD PROPAGATION : Implementación recursiva de la REGLA DE LA CADENA

- Cálculos corresponden a capa de salida y activación sigmoid

→ RED NEURONAL

$$\frac{\partial J}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial y}{\partial a^{(i)}} \cdot \frac{\partial a^{(i)}}{\partial z^{(i)}} \cdot \frac{\partial z^{(i)}}{\partial w_i}$$

$$\frac{\partial y}{\partial a^{(i)}} = \frac{a^{(i)} - y^{(i)}}{a^{(i)}(1 - a^{(i)})}$$

$$\frac{\partial a^{(i)}}{\partial z^{(i)}} = (a^{(i)})^2 e^{-2z^{(i)}} = a^{(i)}(1 - a^{(i)})$$

$$\frac{\partial z^{(i)}}{\partial w_i} = x^{(i)}$$

$$\frac{\partial y}{\partial b} = \frac{\partial y}{\partial a^{(i)}} \cdot \frac{\partial a^{(i)}}{\partial z^{(i)}}$$

CAPA L

$$dz^{[L]} = \frac{1}{m} \sum_{i=1}^m (a^{[L]}(i) - y^{(i)})$$

$$dw^{[L]} = \frac{1}{m} \sum_{i=1}^m dz^{[L]}(i) \underbrace{a^{[L-1]}(i)}_{\downarrow}$$

$$z^{[L]}(i) = w^{[L]} a^{[L-1]}(i) + b^{[L]}$$

$$db^{[L]} = \frac{1}{m} \sum_{i=1}^m dz^{[L]}(i)$$

- CAPA $1 \leq \lambda \leq L-1$, $\lambda = k$

$$z^{[k]}(i) = W^{[k]} a^{[k-1]}(i) + b^{[k]}$$

$$a^{[k]}(i) = g^{[k]}(z^{[k]}(i))$$

$$\Rightarrow \frac{\partial a^{[k]}}{\partial w_j^{[k]}} = \frac{\partial g^{[k]}(z^{[k]})}{\partial z} a_j^{[k-1]} = dw_j^{[k]}$$

$$\Rightarrow \frac{\partial a^{[k]}}{\partial b_j} = \underbrace{\frac{\partial g^{[k]}(z^{[k]})}{\partial z}}_{dz^{[k]}}$$