Back propagation

El modo de entrenar una red, que lo cambio todo

biases

pesos

5

0

0

-5

+1

-1

2

-2

$x$

$f(x)$

Vimos que una red es capaz de
aproximar cualquier funcion razonable

Rossenblatt ya habia propuesto que para entrenar a una red, se podia trabajar iterativamente, computando la diferencia entre lo que plantea la red y lo esperado para el ejemplo, y modificando a la red de tal modo de achicar esa diferencia

(data, resultado)=(X, Y)

Rossenblatt ya habia propuesto que para entrenar a una red, se podia trabajar iterativamente, computando la diferencia entre lo que plantea la red y lo esperado para el ejemplo, y modificando a la red de tal modo de achicar esa diferencia
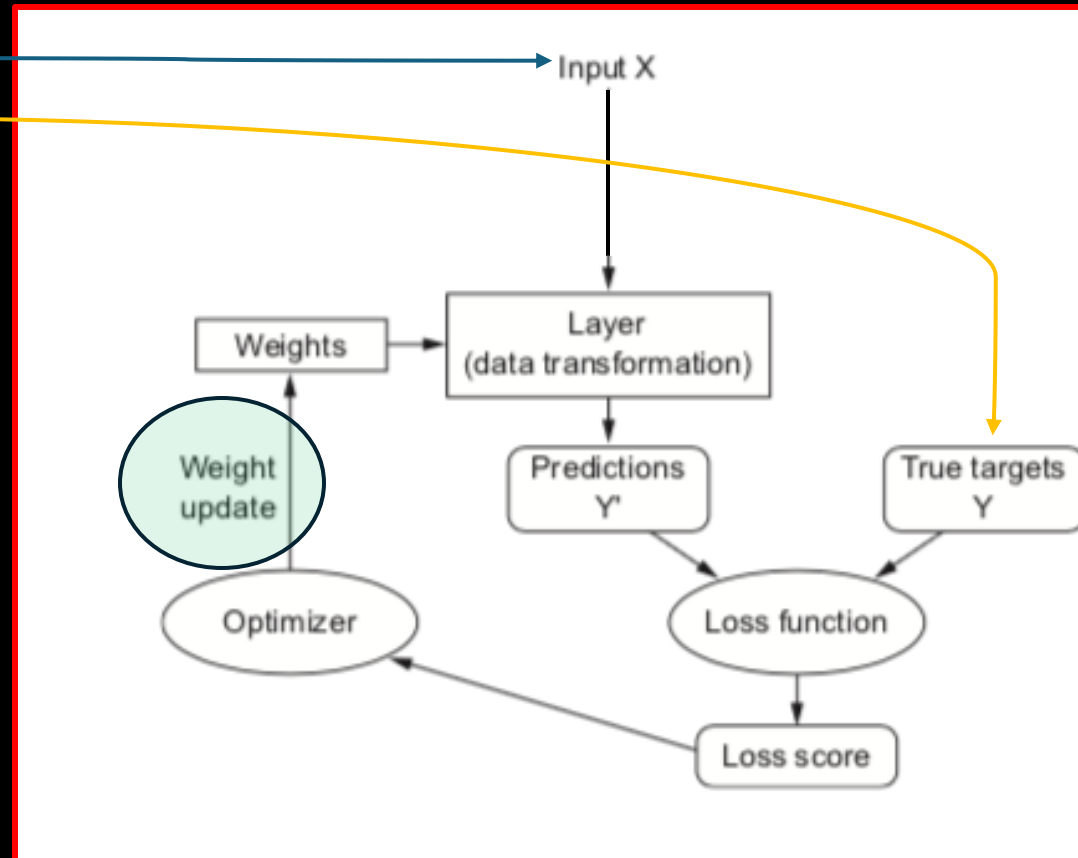
(data, resultado)=(X, Y)

Como hacer el update?



Input X

Weights → Layer (data transformation)

Weight update

Optimizer

Predictions Y'

True targets Y

Loss function

Loss score

# Redes de propagación hacia delante y su entrenamiento
## (Rumelhart, Hinton, Williams, Nature 1986)

# a. classification

1. Traduce the input image into a matrix

$x_j$

2. For a given set of parameters:

inputs

Probability of being the sonogram of the song of a finch

Probability of being the sonogram of the song of a canary

input layer          hidden layer          output layer

**Input, un vector de "features"**

$$x_j = S\left(\sum_{i=1}^{n} W_{ji}\, x_i\right), \qquad s(x) = \frac{1}{1 + e^{-x}}$$

$$x_j = S\left(\sum_{i=1}^{n} W_{ji}\, x_i\right)$$

$$s(x) = \frac{1}{1 + e^{-x}}$$

inputs

Se traduce el diagrama en numeros sobre una matriz

$x_j$

If with the present parameters, we get an error

We correct the weights

0.7 1

0.3 0

input layer     hidden layer     output layer

$$x_j = S(\sum_{i=1}^{n} W_{ji} \, x_i)$$

$$s(x) = \frac{1}{1 + e^{-x}}$$

**Como hacer para que esa correccion no sea solo un "ir a otro lado", sino "ir a uno mejor"?**

$x_j$

If with the present parameters, we get an error

inputs

We correct the weights

0.7 | 1

0.3 | 0

Se traduce el diagrama en numeros sobre una matriz

input layer          hidden layer          output layer

$$x_j = S(\sum_{i=1}^{n} W_{ji} \, x_i)$$

$$s(x) = \frac{1}{1 + e^{-x}}$$

# How does it actually work?
## (let us start with one feedforward architecture)

$O_i$

Output layer

$W_{ii}$   weights

The class the image belongs to

Input layer

The pixels of an image

# How does it actually work?
## (let us start with one feedforward architecture)

How to choose them?

$O_i$

Output layer

The class the image belongs to

$W_{ii}$   weights

Input layer

The pixels of an image

# How does it actually work?

Suposse the unit *j is* in the last layer,
and we touch a weight
in the immediately previous layer

Expected value in unit j

Actual output in unit j

$y_j$    $o_j$

$W_{ii}$

$o_i$

$E = y_j - o_j$

(*E* stands for "error")

# How does it actually work?

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

We want to **correct the weight** $W_{ij}$ in order to decrease the error

$O_j$

$W_{ii}$

$O_i$

# How does it actually work?

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}}.$$

$O_j$

$W_{ii}$

$O_i$

We use the **chain rule**:

the error changes because changing the weight $W_{ij}$
1. It changes the net activity arriving at j
2. If the net activity changes, the output of j changes
3. If the output of j changes, the error changes

# How does it actually work?

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

$O_j$

$W_{ii}$

$O_i$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \boxed{\frac{\partial net_j}{\partial W_{ij}}}$$

$$\boxed{\frac{\partial net_j}{\partial W_{ij}}} = O_i. \quad \text{since } net_j = \sum W_{ij} O_i$$

# How does it actually work?

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}}.$$

$O_j$

$W_{ii}$

$O_i$

$$\frac{\partial net_j}{\partial W_{ij}} = O_i. \quad \text{Since } net_j = \sum W_{ij} O_i$$

As $O_j = S(net_j)$ and $S(net_j) = 1/(1 + e^{-net_j})$,

then $$\frac{\partial O_j}{\partial net_j} = O_j(1 - O_j).$$

# How does it actually work?

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}}.$$

$$\frac{\partial net_j}{\partial W_{ij}} = O_i. \quad \text{Since } net_j = \sum W_{ij} O_i$$

as $O_j = S(net_j)$ and $S(net_j) = 1/(1 + e^{-net_j})$,

then $\quad \dfrac{\partial O_j}{\partial net_j} = O_j(1 - O_j).$

$$\frac{\partial E}{\partial O_j} = (y^j_{elemento} - O_j).$$

# Example

o1        o2        o3        o4

○         ○         ○         ○

w61

o5    \   o6        o7

○         ○         ○

o8        o9        o10       o11

○         ○         ○         ○

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j}\frac{\partial O_j}{\partial net_j}\frac{\partial net_j}{\partial W_{ij}}.$$
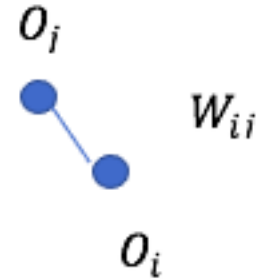
# Example

o1     o2     o3     o4

w61

o5     o6     o7

o8     o9     o10     o11

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j}\frac{\partial O_j}{\partial net_j}\frac{\partial net_j}{\partial W_{ij}}.$$

$$\frac{\partial E}{\partial W61} = \left(o1 - t1\right)\left(o1(1 - o1)\right)o6$$

# Example

o1    o2    o3    o4

o5    o6    o7

w61

o8    o9    o10    o11

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}} .$$
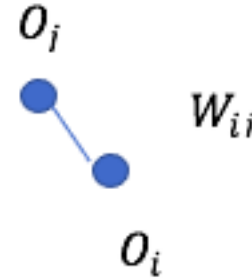
$$\frac{\partial E}{\partial W61} = (o1 - t1)(o1(1 - o1))o6$$

# Example

o1  o2  o3  o4



$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}} \cdot$$

w61

o5  o6  o7

$$\frac{\partial E}{\partial W61} = (o1 - t1)(o1(1 - o1))o6$$

o8  o9  o10  o11

# Example

o1      o2      o3      o4



w61

o5      o6      o7

$$\frac{\partial E}{\partial W_{ij}} = \frac{\partial E}{\partial O_j} \frac{\partial O_j}{\partial net_j} \frac{\partial net_j}{\partial W_{ij}} .$$
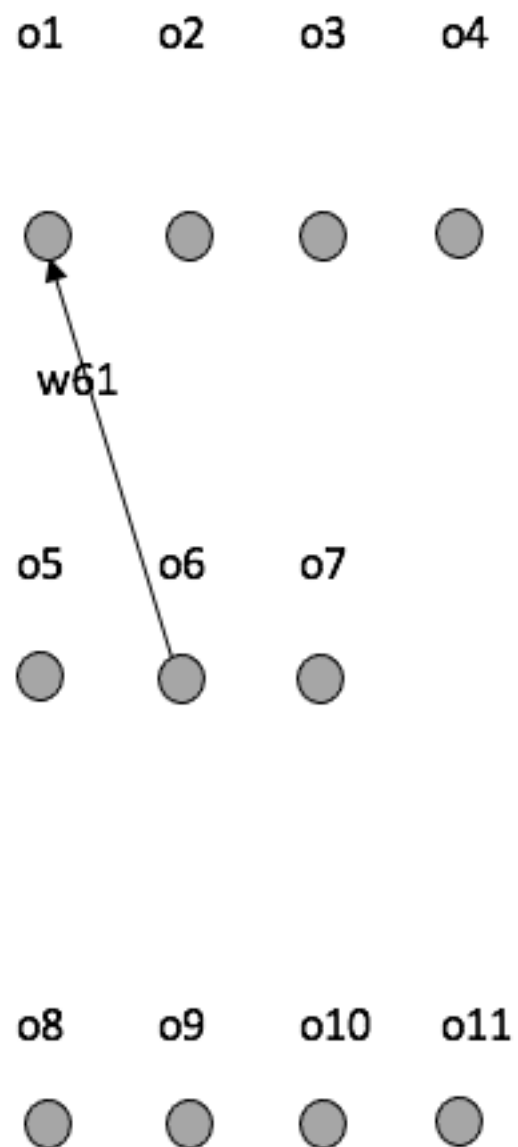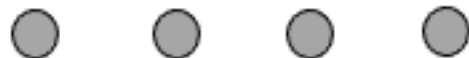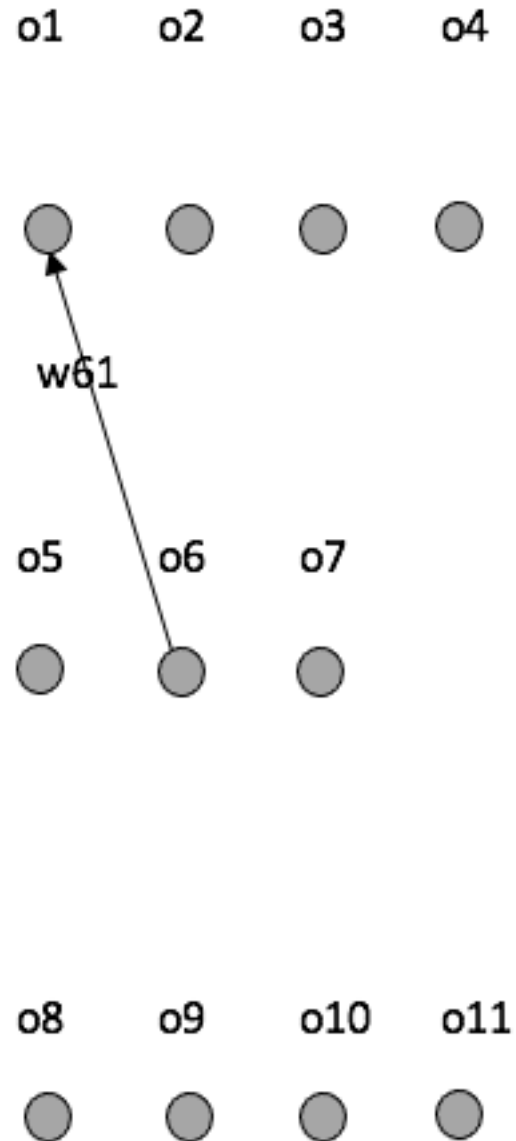
$$\frac{\partial E}{\partial W61} = (o1 - t1)\big(o1(1 - o1)\big)o6$$

The interesting thing is that
it only depends on the outputs O
that I got with the original $W$ values

o8      o9      o10     o11

# Example 2

o1    o2    o3    o4

⬤   ⬤   ⬤   ⬤

o5    o6    o7

⬤   ⬤   ⬤   ⬤

w86

o8    o9   o10   o11

⬤   ⬤   ⬤   ⬤

If the weight being modified is in the **previous to the last layer**, the error changes due to changes in all the units of the last layer

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\frac{\partial O6}{\partial net6}\frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\left(o6(1-o6)\right)o8$$

# Example 2

o1     o2     o3     o4

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\frac{\partial O6}{\partial net6}\frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\big(o6(1-o6)\big)o8$$

o5     o6     o7

$$\frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} = \frac{\partial E}{\partial O1}\frac{\partial O1}{\partial net1}\frac{\partial net1}{\partial O6} + \frac{\partial E}{\partial O2}\frac{\partial O2}{\partial net2}\frac{\partial net2}{\partial O6} + \frac{\partial E}{\partial O3}\frac{\partial O3}{\partial net3}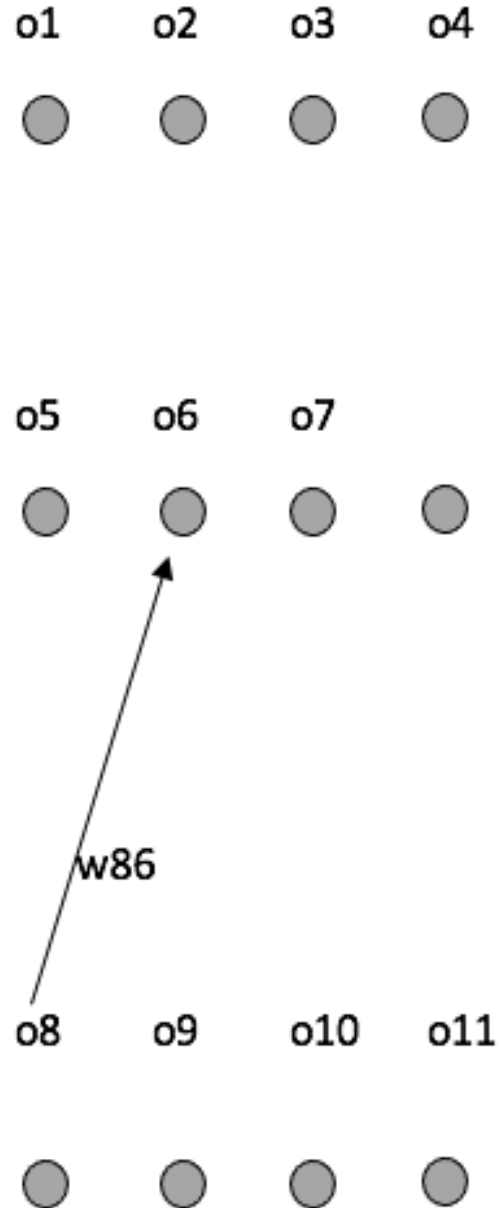\frac{\partial net3}{\partial O6} + \frac{\partial E}{\partial O4}\frac{\partial O4}{\partial net4}\frac{\partial net4}{\partial O6}$$

w86

o8     o9     o10     o11

Example 2

o1  o2  o3  o4

o5  o6  o7

w86

o8  o9  o10  o11

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} \frac{\partial O6}{\partial net6} \frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} \big(o6(1-o6)\big)o8$$

$$\frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} = \frac{\partial E}{\partial O1} \frac{\partial O1}{\partial net1} \frac{\partial net1}{\partial O6} + \frac{\partial E}{\partial O2} \frac{\partial O2}{\partial net2} \frac{\partial net2}{\partial O6} + \frac{\partial E}{\partial O3} \frac{\partial O3}{\partial net3} \frac{\partial net3}{\partial O6} + \frac{\partial E}{\partial O4} \frac{\partial O4}{\partial net4} \frac{\partial net4}{\partial O6}$$

$$\frac{\partial E}{\partial O6} = (O1-t1)\big(O1(1-O1)\big)W16$$
$$+(O2-t2)\big(O2(1-O2)\big)W26$$
$$+(O3-t3)\big(O3(1-O3)\big)W36$$
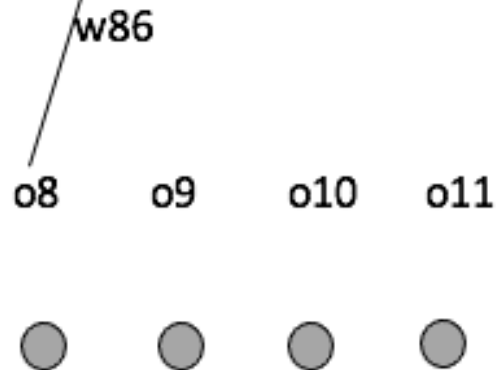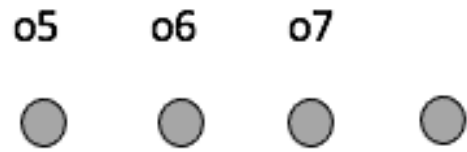$$+(O4-t4)\big(O4(1-O4)\big)W46$$

## Example 2

o1    o2    o3    o4

○    ○    ○    ○

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\frac{\partial O6}{\partial net6}\frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\big(o6(1-o6)\big)o8$$

o5    o6    o7

○    ○    ○    ○

$$\frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} = \frac{\partial E}{\partial O1}\frac{\partial O1}{\partial net1}\frac{\partial net1}{\partial O6} + \frac{\partial E}{\partial O2}\frac{\partial O2}{\partial net2}\frac{\partial net2}{\partial O6} + \frac{\partial E}{\partial O3}\frac{\partial O3}{\partial net3}\frac{\partial net3}{\partial O6} + \frac{\partial E}{\partial O4}\frac{\partial O4}{\partial net4}\frac{\partial net4}{\partial O6}$$

w86

$$\frac{\partial E}{\partial O6} = (O1 - t1)\big(O1(1-O1)\big)W16$$
$$+(O2 - t2)\big(O2(1-O2)\big)W26$$
$$+(O3 - t3)\big(O3(1-O3)\big)W36$$
$$+(O4 - t4)\big(O4(1-O4)\big)W46$$

o8    o9    o10    o11

○    ○    ○    ○

Example 2

o1   o2   o3   o4

○    ○    ○    ○

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} \frac{\partial O6}{\partial net6} \frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} \big(o6(1-o6)\big)o8$$

o5   o6   o7

○    ○    ○    ○

$$\frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} = \frac{\partial E}{\partial O1} \frac{\partial O1}{\partial net1} \frac{\partial net1}{\partial O6} + \frac{\partial E}{\partial O2} \frac{\partial O2}{\partial net2} \frac{\partial net2}{\partial O6} + \frac{\partial E}{\partial O3} \frac{\partial O3}{\partial net3} \frac{\partial net3}{\partial O6} + \frac{\partial E}{\partial O4} \frac{\partial O4}{\partial net4} \frac{\partial net4}{\partial O6}$$

w86

$$\frac{\partial E}{\partial O6} = (O1 - t1)\big(O1(1-O1)\big)W16$$
$$+(O2 - t2)\big(O2(1-O2)\big)W26$$
$$+(O3 - t3)\big(O3(1-O3)\big)W36$$
$$+(O4 - t4)\big(O4(1-O4)\big)W46$$

o8   o9   o10   o11

○    ○    ○    ○

## Example 2

o1    o2    o3    o4

⬤    ⬤    ⬤    ⬤

$$\frac{\partial E}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\frac{\partial O6}{\partial net6}\frac{\partial net6}{\partial W86} = \frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6}\big(o6(1-o6)\big)o8$$

o5    o6    o7

⬤    ⬤    ⬤    ⬤

$$\frac{\partial \boldsymbol{E}}{\partial \boldsymbol{O}6} = \frac{\partial E}{\partial O1}\frac{\partial O1}{\partial net1}\frac{\partial net1}{\partial O6} + \frac{\partial E}{\partial O2}\frac{\partial O2}{\partial net2}\frac{\partial net2}{\partial O6} + \frac{\partial E}{\partial O3}\frac{\partial O3}{\partial net3}\frac{\partial net3}{\partial O6} + \frac{\partial E}{\partial O4}\frac{\partial O4}{\partial net4}\frac{\partial net4}{\partial O6}$$

w86

o8    o9    o10    o11

⬤    ⬤    ⬤    ⬤

$$\frac{\partial E}{\partial O6} = (O1 - t1)\big(O1(1-O1)\big)W16$$
$$+(O2-t2)\big(O2(1-O2)\big)W26$$
$$+(O3-t3)\big(O3(1-O3)\big)W36$$
$$+(O4-t4)\big(O4(1-O4)\big)W46$$

Tedious, but only
*O* y *W* are involved

inputs

we traduce
the image
into numbers

$x_j$

And request, for
example, this output

1

0

input layer    hidden layer    output layer

$$x_j = S(\sum_{i=1}^{n} W_{ji}\, x_i)$$

$$s(x) = \frac{1}{1 + e^{-x}}$$

$x_j$

We correct the weights

inputs

0.7 | 1

0.3 | 0

Se traduce el diagrama en numeros sobre una matriz

input layer

hidden layer

output layer

$$x_j = S(\sum_{i=1}^{n} W_{ji}\, x_i)$$

$$s(x) = \frac{1}{1 + e^{-x}}$$

Does this mean that the scary AI that will leave us out of jobs is



= the chain rule?

Does this mean that the scary AI that will leave us out of jobs is



= the chain rule?

Yeap.

a

Input Data          Encoded Data          Reconstructed Data

Autoencoder

b

Max-Pooling    Convolution+Relu    Max-Pooling    Fully-Connected    Fully-Connected

Convolutional neural network

c

o

Unfold

$o_{t-1}$          $o_t$          $o_{t+1}$

W          W          W          W

V

h          $h_{t-1}$          $h_t$          $h_{t+1}$

V          V          V          V

U          U          U          U

x          $x_{t-1}$          $x_t$          $x_{t+1}$

Recurrent neural network

# Autoencoder
# (dimensional reduction)



Input Data          Encoded Data          Reconstructed Data

Convolution+Relu  Max-Pooling  Convolution+Relu  Max-Pooling  Fully-Connected  Fully-Connected

P(dog)=0.8
P(wolf)=0.1
P(cow)=0.001

P(bird)=0.0001

Convolutional neural network
classify

# Recurrent neural network
## (temporal predictions)

a

# En todas, lo esencial es buscar minimizar el error entre lo predicho y lo deseado segun los ejemplos

prediccion

El dato de test

Error E (loss)

Y lo emplea para calcular los cambios en la red:

$$\Delta W_{ij} = -\eta \frac{\partial E}{\partial W_{ij}},$$

Con la idea de encontrar un mínimo en el cual $\Delta W_{ij} = 0$, resulta que

Datos



| Datos de entrenamiento | Datos de test |
|---|---|

Over-fitting

(a)

Error

test

entrenamiento

epocas

Under-fitting

(b)

Error

epocas

Datos



| Datos de entrenamiento | Datos de test |
|---|---|



Over-fitting

(a)

Error

test

entrenamiento

epocas

Under-fitting

(b)

Error

epocas

La explicacion mas clara del overfitting

Training and validation loss

**Modos de evitar esto:**

1. **Hacer chico al modelo**. Un modelo con muchos parametros tiene lugar para aprender cada ejemplo sin necesidad de generalizar. Hay que empezar con uno chico, que no ajuste al training set, y empezar a agrandar. Nunca al revés.
2. **Hay que regularizar**. Los modelos, para que sean sencillos y no ajusten caprichosamente a los datos, deben tener pocos pesos, y de tamaños comparables. Una red con un parámetro desproporcionadamente mas grande que el resto probablemente este haciendo contorsiones para ajustar algún conjunto particular de datos. Para esto se penaliza tener muchos parámetros distintos de cero.
3. Hay trucos, como el **dropout** (setear a cero algunos parámetros al azar)

# Librerías para redes

# Librerías para redes:   K Keras

Pasos para entrenar una red
        neuronal como esta:

# Librerías para redes: K Keras

Pasos para entrenar una red neuronal como esta:



Inputs Hidden Output(s)

1) **Crear el Objeto**

2) **Definir la arquitectura (agregar capas)**

3) **Compilar**

4) **Entrenar**

# Librerías para redes: K Keras

Pasos para entrenar una red neuronal como esta:



1)
```
from keras.models import Sequential
model = Sequential()
```

2) **Definir la arquitectura (agregar capas)**

3) **Compilar**

4) **Entrenar**

# Librerías para redes: K Keras

Pasos para entrenar una red neuronal como esta:



Inputs    Hidden    Output(s)

1)
```
from keras.models import Sequential
model = Sequential()
```

2)
```
model.add(Dense(5), activation = 'sigmoid')
model.add(Dense(1),activation = linear)
```

3) **Compilar**

4) **Entrenar**

# Librerías para redes: K Keras

Pasos para entrenar una red neuronal como esta:



Inputs

Hidden

Output(s)

1)
```
from keras.models import Sequential
model = Sequential()
```

2)
```
model.add(Dense(5), activation = 'sigmoid')
model.add(Dense(1),activation = linear)
```

3)
```
model.compile(loss = 'mse', optimizer='adam')
```

4) **Entrenar**

# Librerías para redes: **K** Keras

Pasos para entrenar una red neuronal como esta:



Inputs
Hidden
Output(s)

1)
```
from keras.models import Sequential
model = Sequential()
```

2)
```
model.add(Dense(5), activation = 'sigmoid')
model.add(Dense(1),activation = linear)
```

3)
```
model.compile(loss = 'mse', optimizer='adam')
```
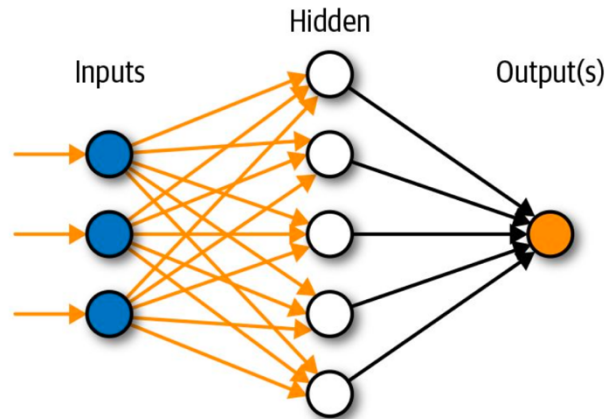
4)
```
model.fit(X,y,batch_size=32, epochs=100)
```

# Función de Costo: Definición

Es aquello que buscamos optimizar:

- Costo de una instancia:

$$J_i(\hat{y}_i, y_i) = J_i(f_\Theta(x_i), y_i)$$

# Función de Costo: Definición

Es aquello que buscamos optimizar.

- Costo de una instancia:

$$J_i(\hat{y}_i, y_i) = J_i(f_\Theta(x_i), y_i)$$

Predicción
del modelo

Valor Real

Depende de los
parámetros
(pesos de la red)

# Función de Costo: Selección

¿Cómo elijo la función de costo J? Hay muchas disponibles.

### Clasificación

- Binary Cross Entropy
- Categorical Cross-entropy
- Poisson Loss
- Custom

### Regresión

- Mean Squared Error (MSE)
- Mean Absolute Error (MAE)
- MAPE
- Custom

https://neptune.ai/blog/keras-loss-functions                    https://keras.io/api/losses/

# Función de Costo: Selección

¿Cómo elijo la función de costo J? Hay muchas disponibles.

**Vamos a ver 3 escenarios:**

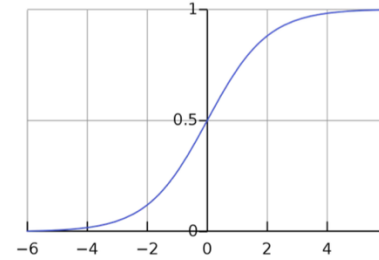- **Escenario 1:**  **Clasificación Binaria**

- **Escenario 2:**  **Clasificación Multiclase**

- **Escenario 3:**  **Regresión**

# Escenario 1: Clasificación Binaria

- **Activacion:** Sigmoide

$$h(t) = \frac{1}{1+e^{-t}}$$



- **Función de costo:** Binary Cross-Entropy

$$J(\Theta) = \frac{1}{N} \sum_i -y_i \cdot log\left(f_\Theta(x_i)\right) - (1 - y_i) \cdot log\left(1 - f_\Theta(x_i)\right)$$

# Escenario 1: Clasificación Binaria

- **Activacion:** **Sigmoide**

```
model.add(layers.Dense(1, activation='sigmoid'))
```

- **Función de costo:** **Binary Cross-Entropy**

```
model.compile(optimizer="Adam",
              loss=tf.keras.losses.BinaryCrossentropy())
```

# Escenario 2: Clasificación Multiclase

- **Activacion: SoftMax**

$$h(t)_j = \frac{e^{t_j}}{\sum_{k=1}^{K} e^{t_j}}$$

Generalización de la función logística (Normalizada)

- **Función de costo: Categorical Cross-Entropy**

$$J(\Theta) = -\frac{1}{N} \sum_i \sum_k y_i \cdot log\left(f_\Theta(x_i)\right)$$

Generalización de la binaria (k clases)

# Escenario 2: Clasificación Multiclase

- **Activacion:** SoftMax

```
model.add(layers.Dense(num_clases, activation='softmax'))
```

- **Función de costo:** Categorical Cross-Entropy

```
model.compile(optimizer="Adam",
              loss=tf.keras.losses.categorical_crossentropy())
```
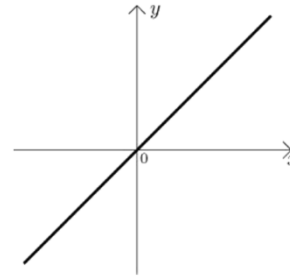
# Escenario 3: Regresión

- **Activacion:** Lineal

$$h(t) = t$$



- **Función de costo:** Mean Squared Error

$$J(\Theta) = \frac{1}{N} \sum_i (f_\Theta(x_i) - y_i)^2$$

# Escenario 3: Regresión

- **Activacion:** **Lineal**

```
model.add(layers.Dense(1, activation='linear'))
```

- **Función de costo:** **Mean Squared Error**

```
model.compile(optimizer="Adam", loss='mse')
```