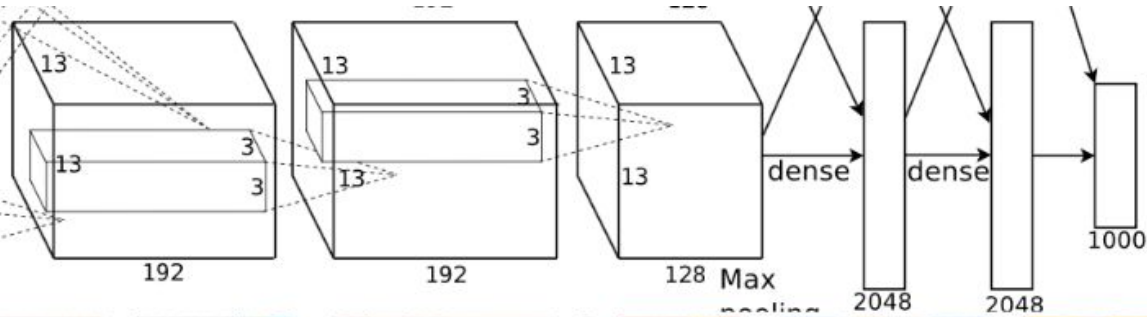


Técnicas Avanzadas de Data Mining y Sistemas Inteligentes

Maestría en Informática
Escuela de Posgrado
Pontificia Universidad Católica del Perú

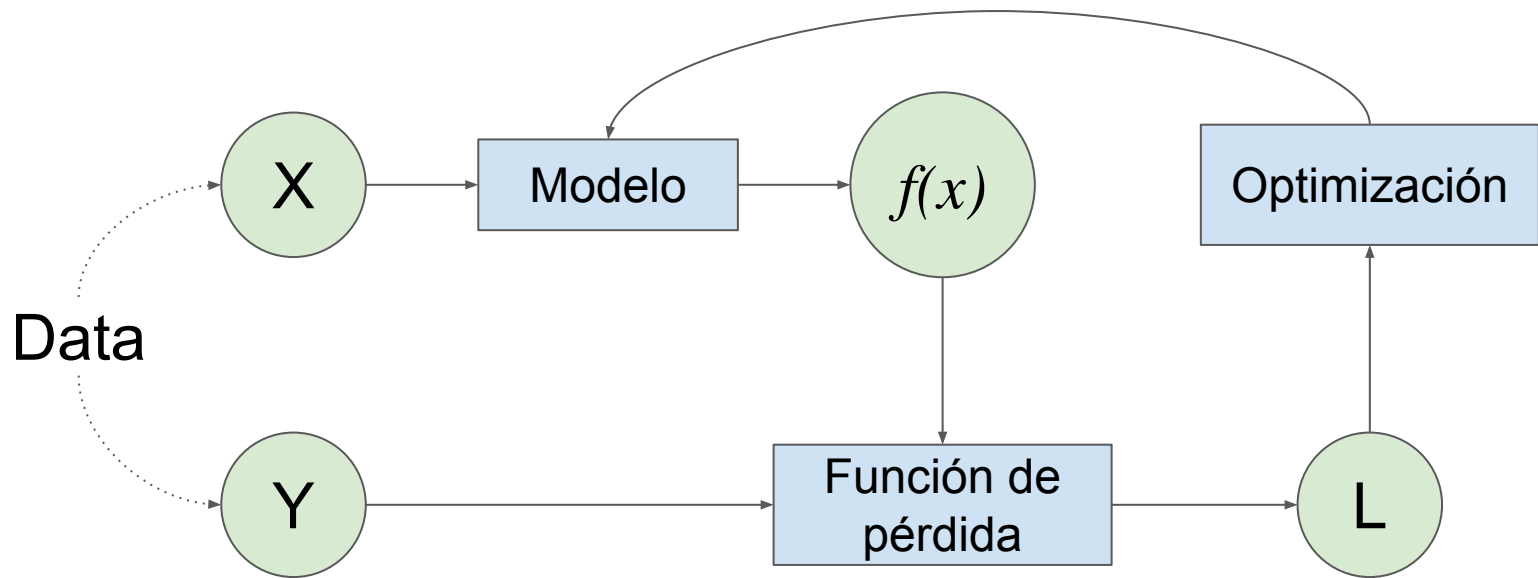
2018-2

Algorithms

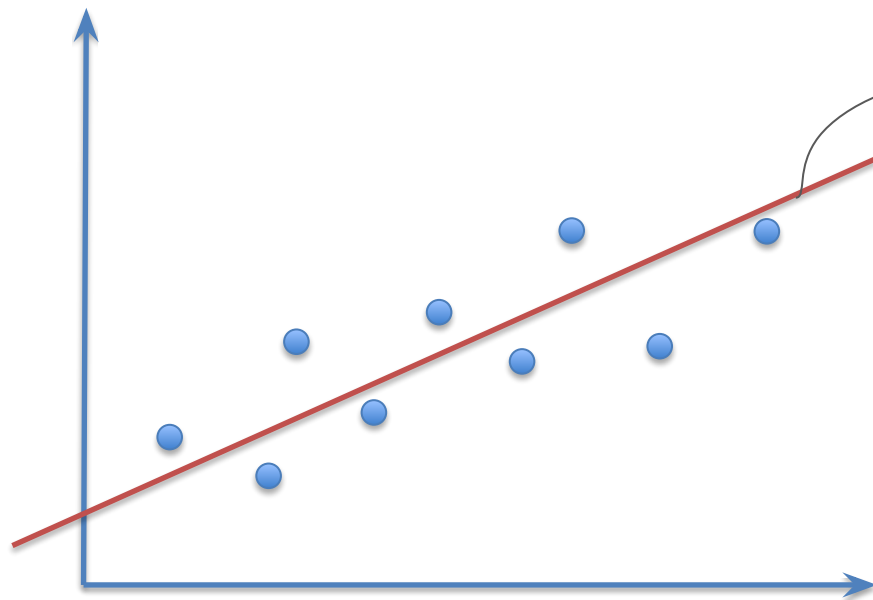


Data

Computation



Regresión lineal simple



$$f(x) = ax + b$$

Donde:

a = pendiente

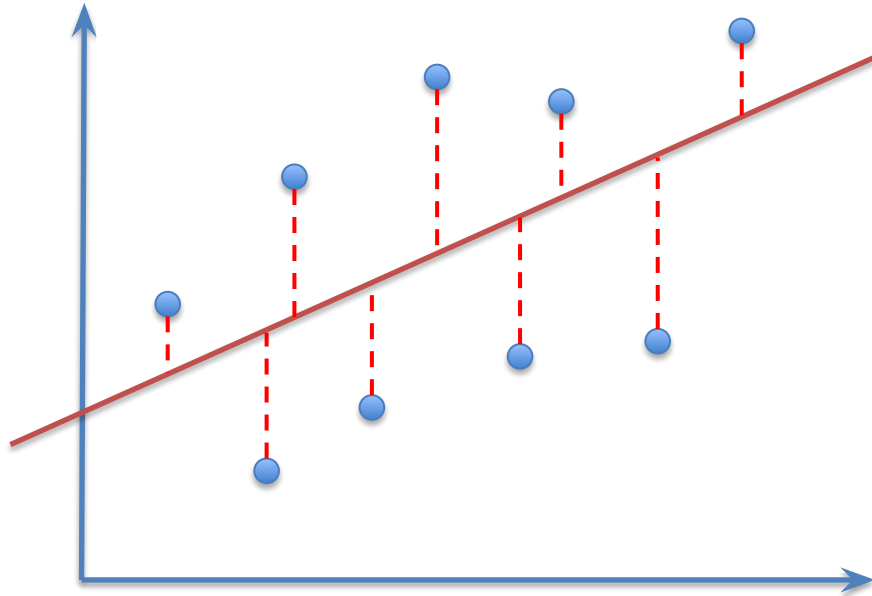
b = término de intercepción (bias)

⬢ Data

— Modelo




$$\begin{array}{ccc}
 \mathbf{x} & & \mathbf{w} \\
 \begin{bmatrix} 1 & 1 \\ 0 & 1 \\ 2 & 1 \\ 1 & 1 \end{bmatrix} & \times & \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 1 \\ 5 \\ 3 \end{bmatrix} \\
 \mathbf{4 \times 2} & & \mathbf{2 \times 1} \quad \mathbf{4 \times 1}
 \end{array}$$

Función de pérdida: MSE (Mean squared error)



Mean squared error:

$$L = \frac{1}{2m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

-  Data x
-  Modelo $f(x)$
-  Error $f(x) - y$

Optimización: Gradient descent

$$f(x_i) = ax_i + b$$

$$L = \frac{1}{2m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

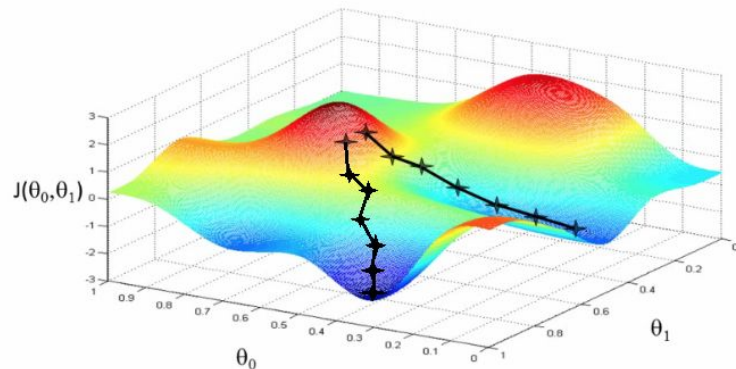
$$\frac{dL}{da} = \sum_{i=1}^m (f(x_i) - y_i)(x_i)$$

$$\frac{dL}{db} = \sum_{i=1}^m (f(x_i) - y_i)$$

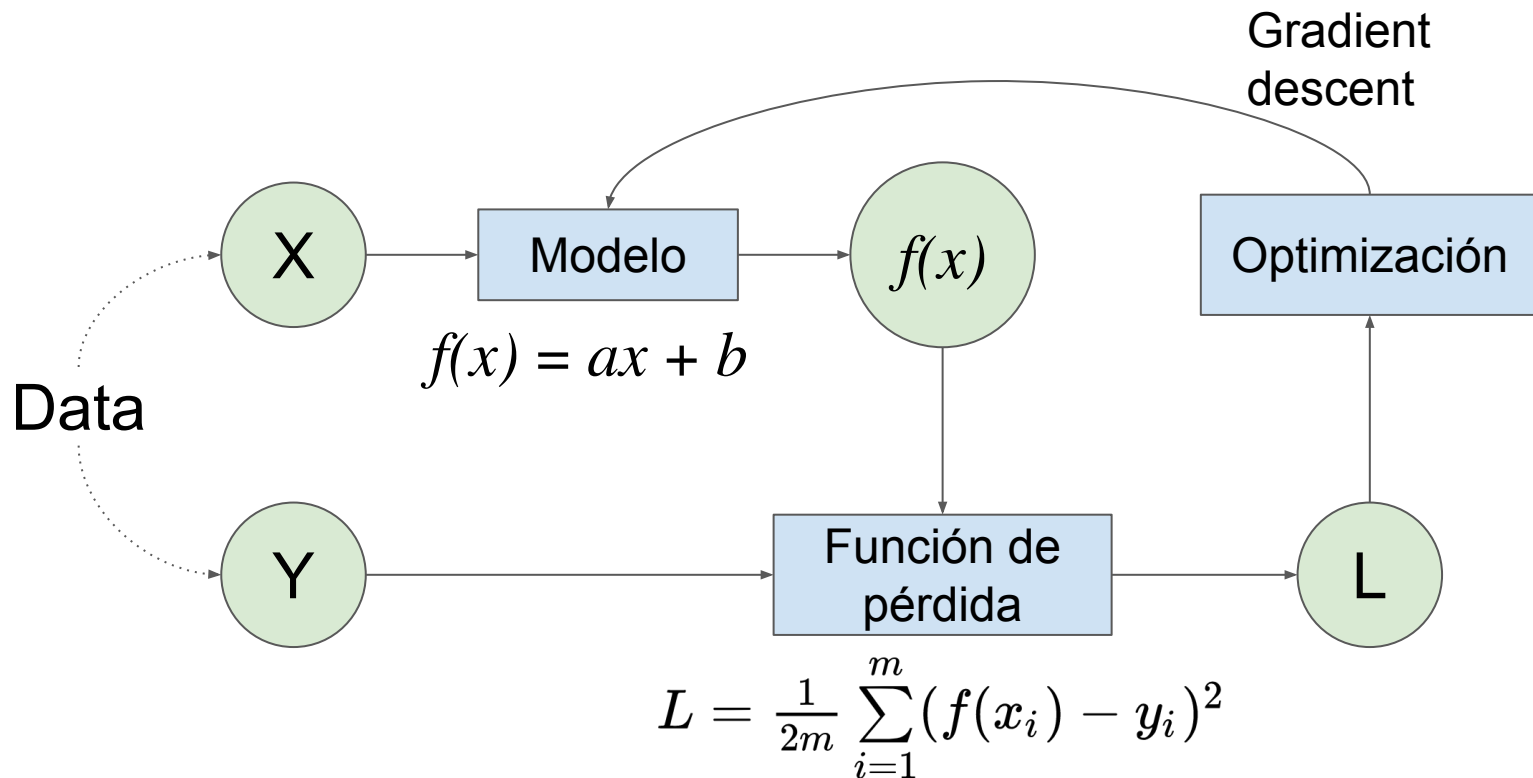
iterar :

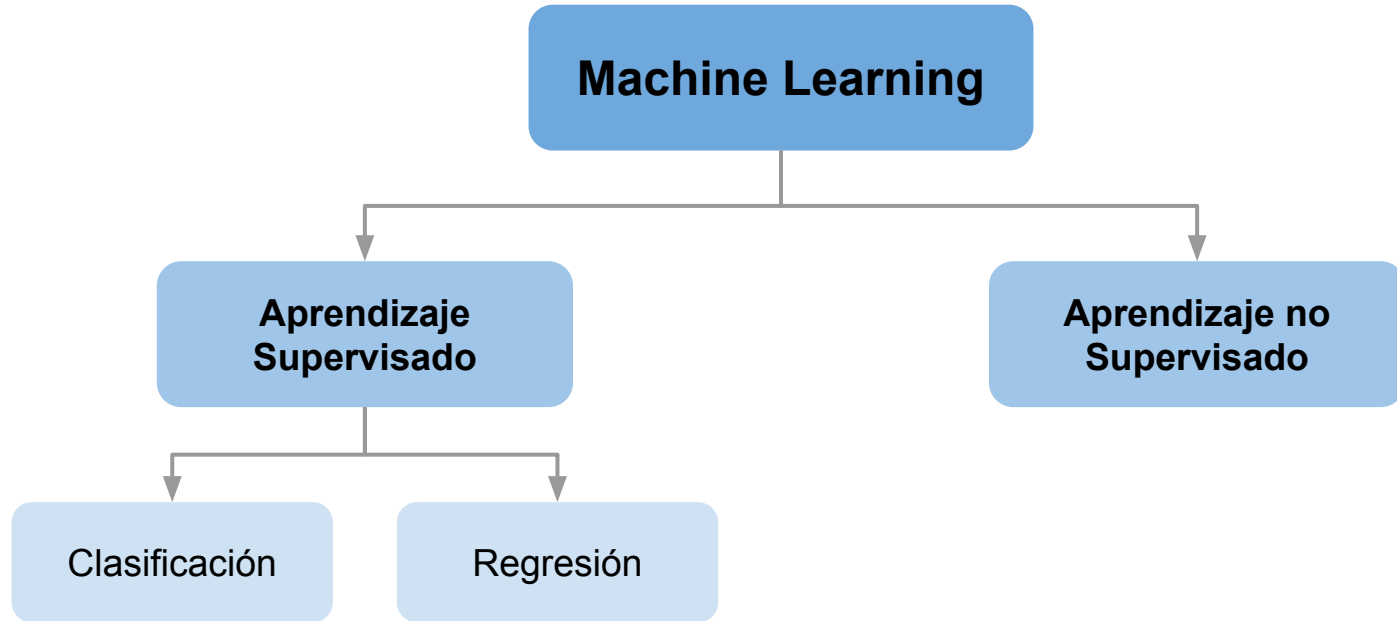
$$new\ a = a - lr * \left(\frac{dL}{da} \right)$$

$$new\ b = b - lr * \left(\frac{dL}{db} \right)$$

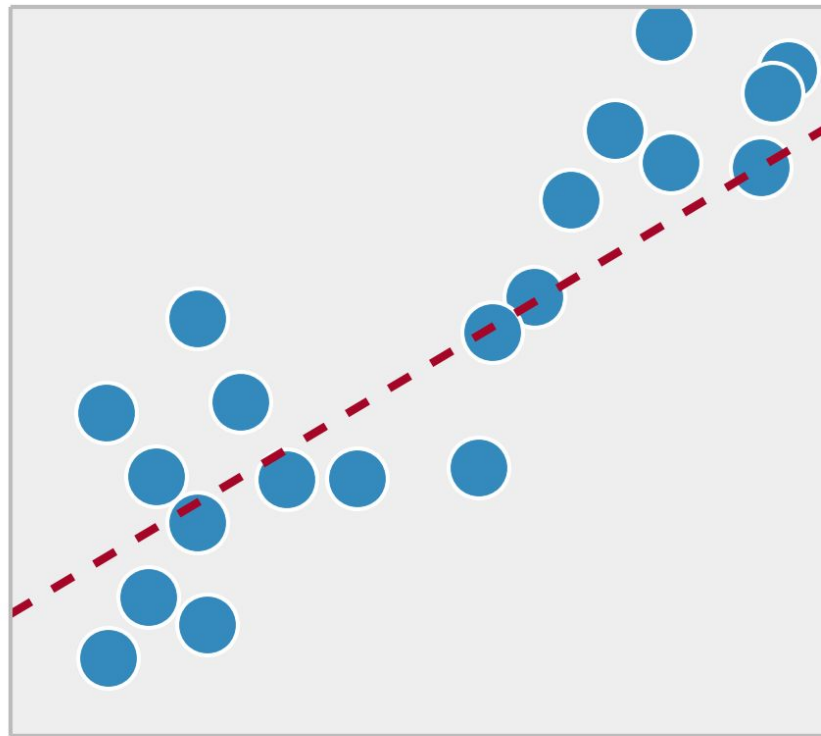
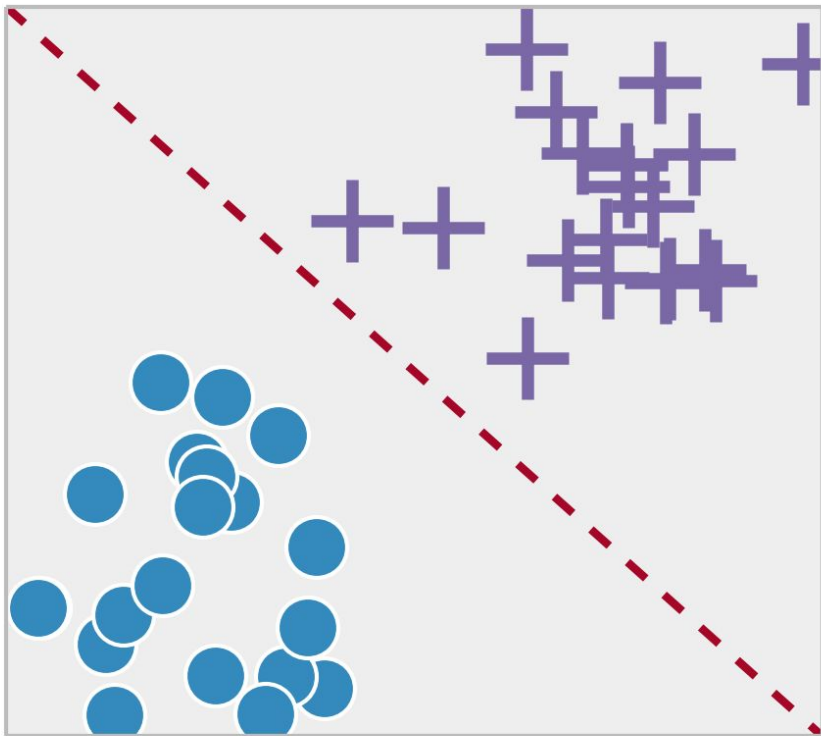


Regresión lineal simple



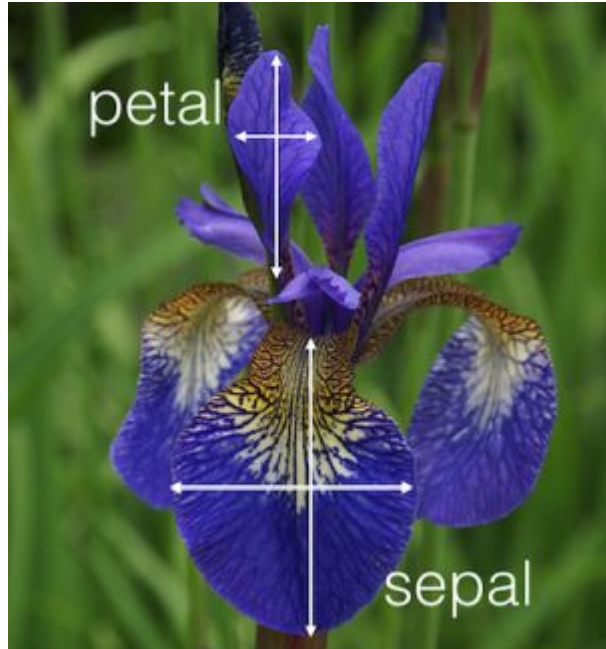


Clasificación vs Regresión



Regresión Logística

Iris dataset



sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2
5.7	4.4	1.5	0.4
6.7	3.1	4.4	1.4
4.8	3.4	1.6	0.2
4.4	3.2	1.3	0.2

Target: class of iris plant

Iris dataset



Iris setosa



Iris versicolor



Iris virginica

Objetivo

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

???

setosa!
~~versicolor~~
~~virginica~~

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

???

Variables
independientes:
[4x1]

setosa!

3 resultados,
indicando la
probabilidad de
cada clase:

[1, 0, 0]

Setosa Versicolor Virginica
(one hot encoding)

[3x1]

Definición del modelo

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

$f(\mathbf{x}, \mathbf{W})$

Variables
independientes:
[4x1]

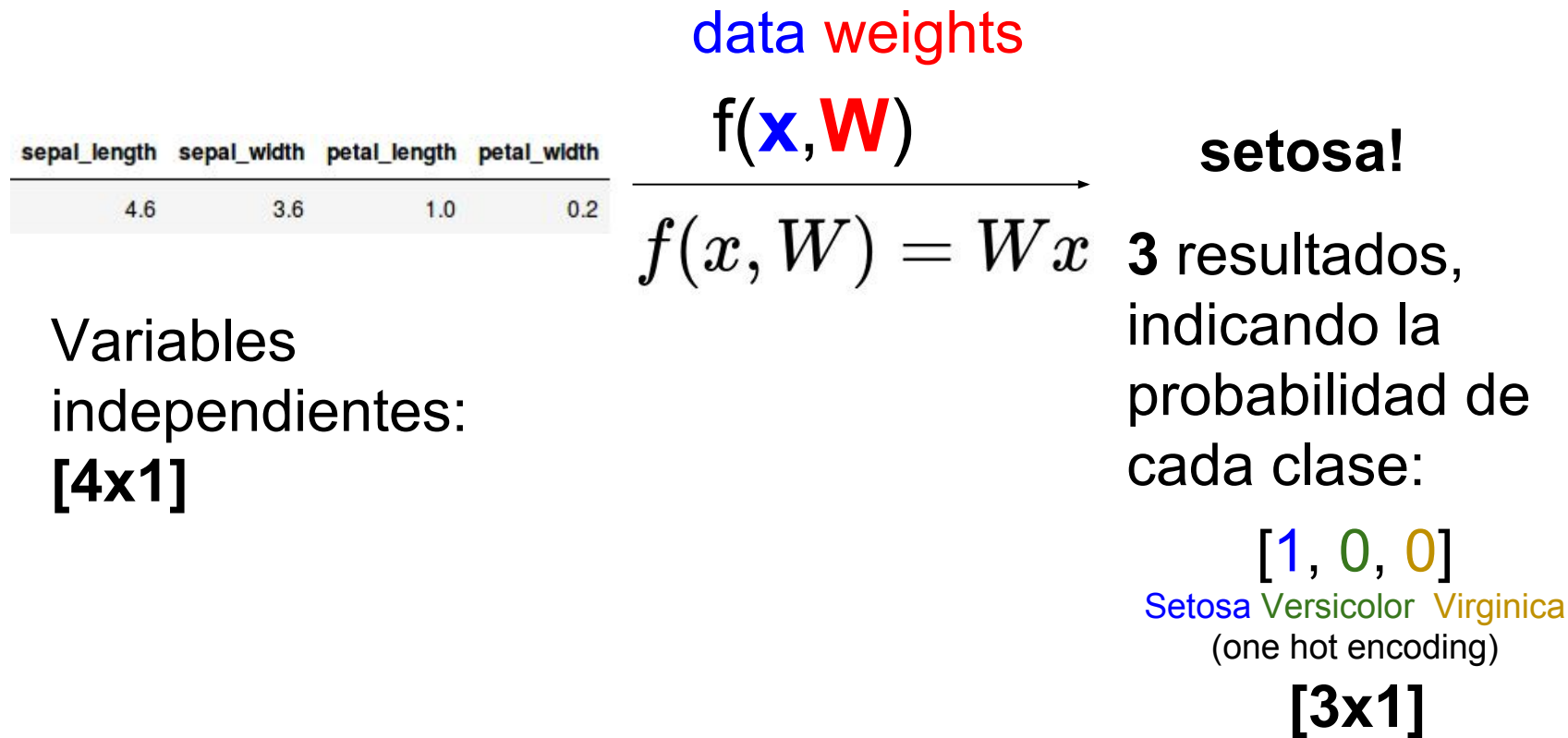
setosa!

3 resultados,
indicando la
probabilidad de
cada clase:

[1, 0, 0]
Setosa Versicolor Virginica
(one hot encoding)

[3x1]

Definición del modelo: Clasificador Lineal



Definición del modelo: Clasificador Lineal

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

setosa!

$$\boxed{f(x, W)} = \boxed{W} \boxed{x} \quad \begin{matrix} 4 \times N \\ 3 \times N \end{matrix}$$

?

parameters, or “weights”

Definición del modelo: Clasificador Lineal

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

setosa!

$$\boxed{f(x, W)} = \boxed{W} \boxed{x}$$

$3 \times N$ $[3 \times 4]$ $4 \times N$

parameters, or “weights”

Definición del modelo: Clasificador Lineal

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

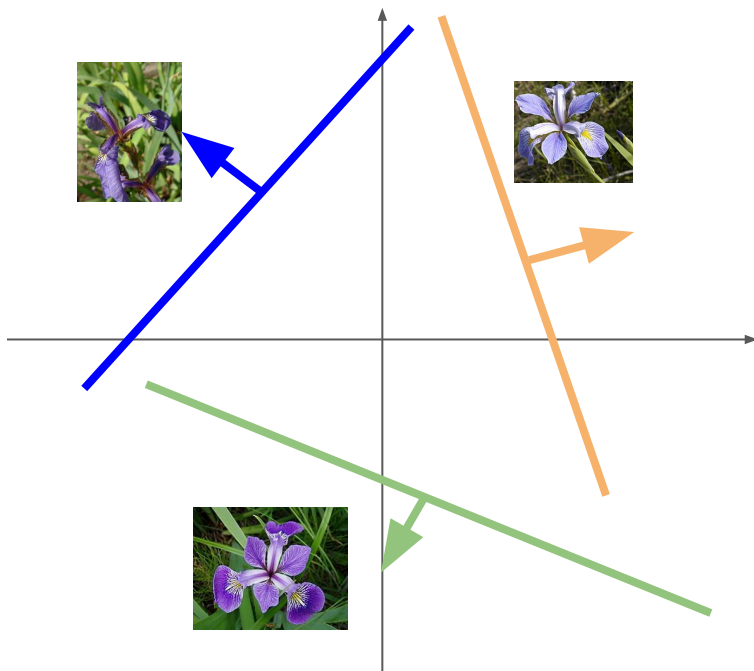
setosa!

$$\boxed{f(x, W)}_{3 \times N} = \boxed{W}_{[3 \times 4]} \boxed{x}_{4 \times N} \boxed{(+b)}_{3 \times 1} \text{ bias}$$

parameters, or “weights”

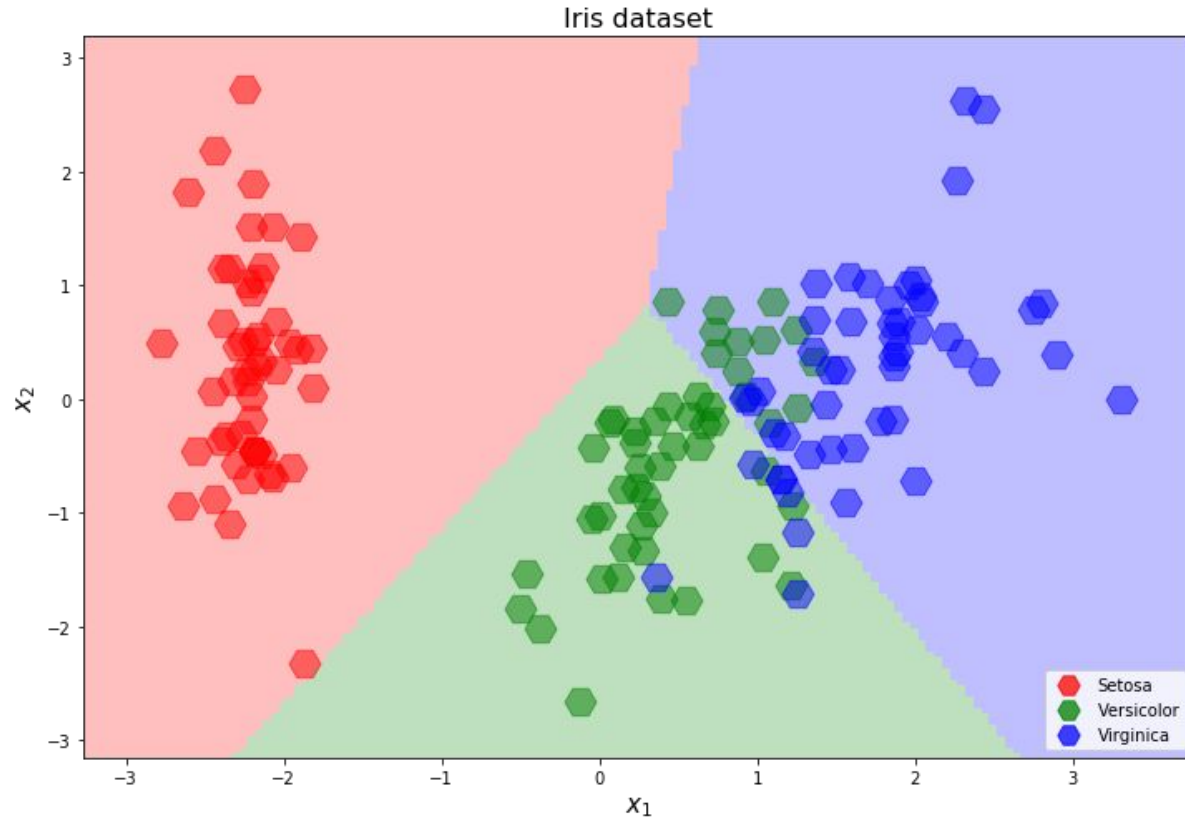
Interpretando el clasificador lineal

$$f(x_i) = Wx_i + b$$



Setosa
Versicolor
Virginica

Interpretando el clasificador lineal



Ej:

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

setosa!

[1, 0, 0]

Setosa Versicolor

Virginica

4.6
3.6
1.0
0.2

X_i

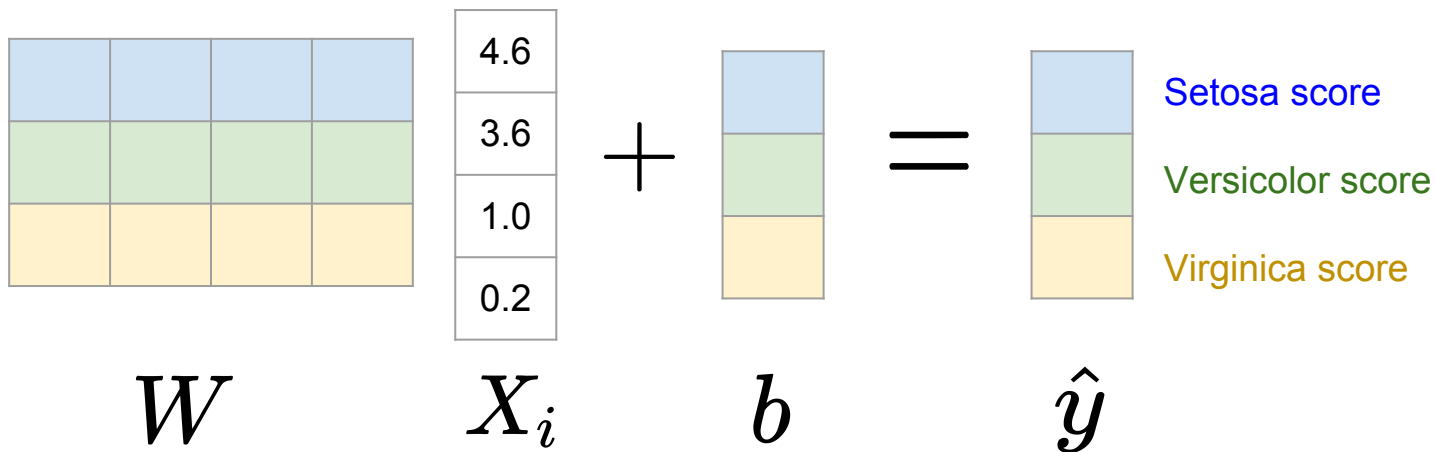
1
0
0

y_i

Ej:

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

$$\hat{y} = f(x_i) = Wx_i + b$$



Ej:

sepal_length	sepal_width	petal_length	petal_width
4.6	3.6	1.0	0.2

Random
numbers!

$$\hat{y} = f(x_i) = Wx_i + b$$

-0.4	-0.1	-2.1	1.6
-1.8	-0.8	0.5	-1.2
-1.1	-0.9	0.6	2.3

 W

4.6
3.6
1.0
0.2

 X_i
 $+$

1.8
0.4
0.1

 b
 $=$

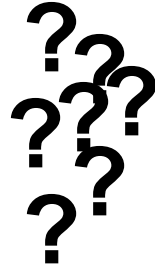
-2.18
-10.5
-7.14

 \hat{y}

Setosa score
Versicolor score
Virginica score

¿Cómo interpretamos estos scores?

-2.18	Setosa score
-10.5	Versicolor score
-7.14	Virginica score



Softmax Classifier (Multinomial Logistic Regression)

-2.18	Setosa score
-10.5	Versicolor score
-7.14	Virginica score

$$\text{softmax}(\hat{y}_i) = \frac{e^{\hat{y}_i}}{\sum_j e^{\hat{y}_j}}$$

Softmax Classifier (Multinomial Logistic Regression)

$$\text{softmax}(\hat{y}_i) = \frac{e^{\hat{y}_i}}{\sum_j e^{\hat{y}_j}}$$

Setosa score

-2.18

Versicolor score

-10.5

Virginica score

-7.14

\hat{y}_i

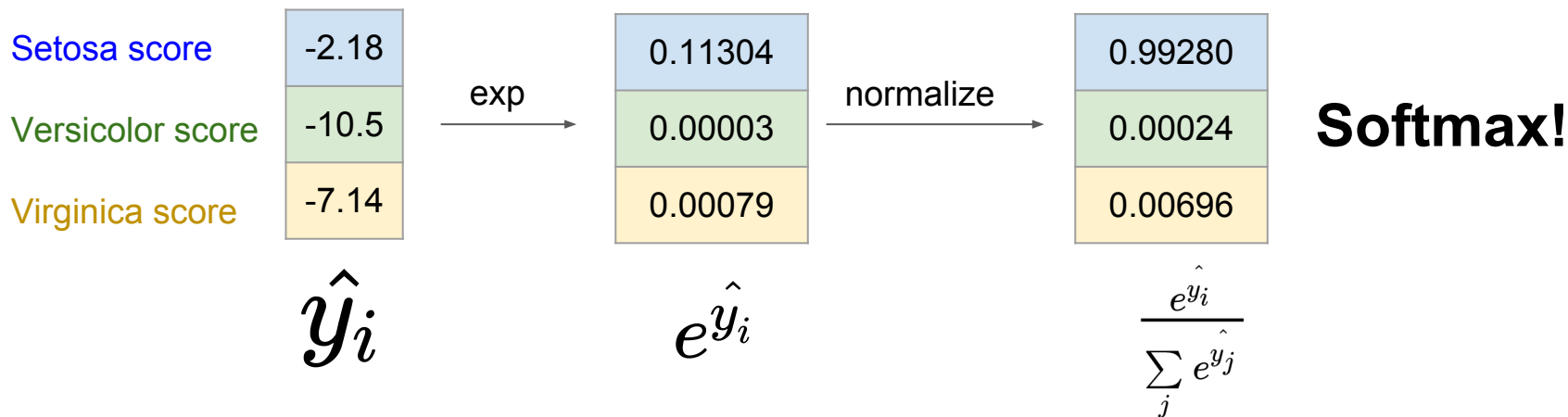
Softmax Classifier (Multinomial Logistic Regression)

$$\text{softmax}(\hat{y}_i) = \frac{e^{\hat{y}_i}}{\sum_j e^{\hat{y}_j}}$$

Setosa score	-2.18	exp	0.11304
Versicolor score	-10.5		0.00003
Virginica score	-7.14		0.00079
	\hat{y}_i		$e^{\hat{y}_i}$

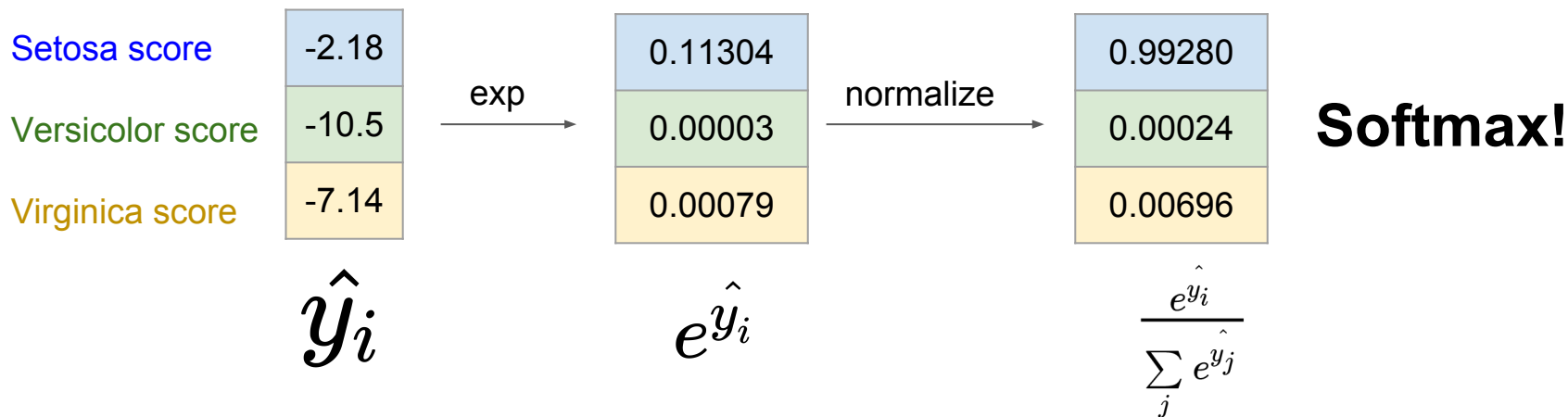
Softmax Classifier (Multinomial Logistic Regression)

$$\text{softmax}(\hat{y}_i) = \frac{e^{\hat{y}_i}}{\sum_j e^{\hat{y}_j}}$$

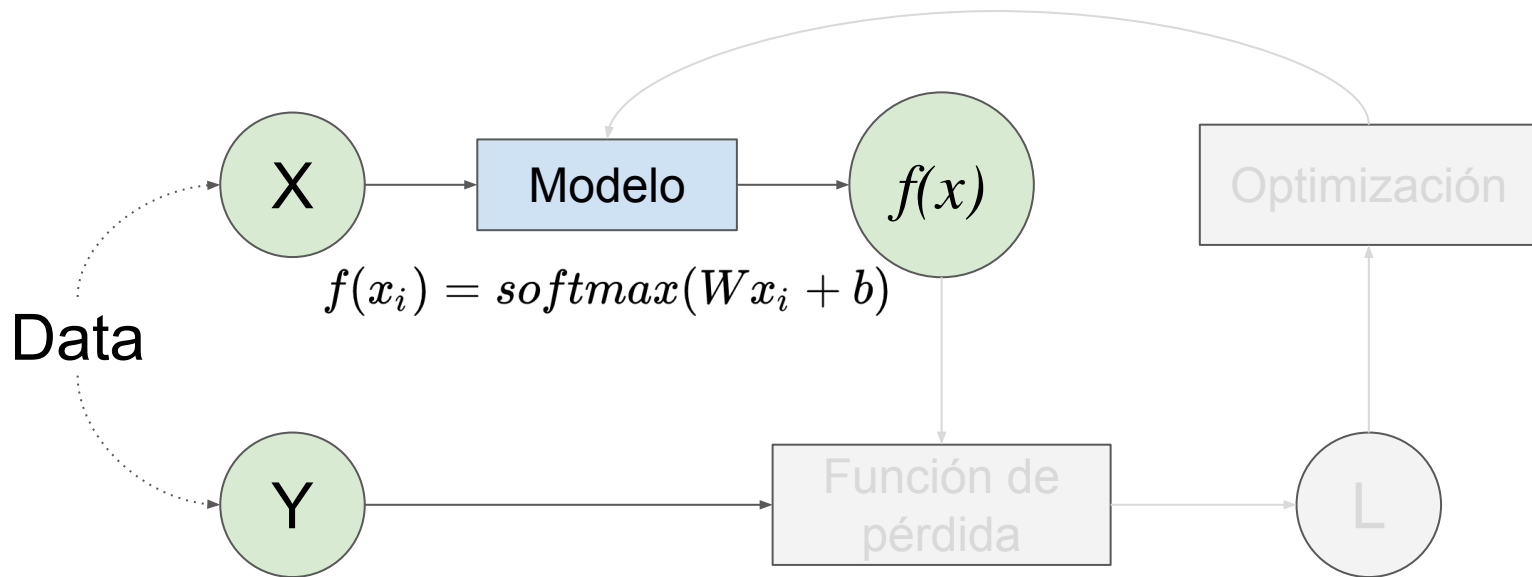


Definición final del modelo

$$f(x_i) = \text{softmax}(Wx_i + b)$$



Regresión logística



¿Cómo evaluamos el modelo?

Cross Entropy (Función de pérdida)

$$L = - \sum p(x) \log q(x)$$

$$L = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

Cross Entropy (Función de pérdida)

$$L = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

The diagram illustrates the Cross Entropy loss function with a specific example for the Iris dataset. It shows the ground truth vector y_i and the predicted probability vector \hat{y}_i for a single instance i .

Ground Truth (y_i): A red box highlights the ground truth vector y_i , which is a 3x1 column vector. The values are 1 for Setosa, 0 for Versicolor, and 0 for Virginica.

Predicted Probabilities (\hat{y}_i): A green box highlights the predicted probability vector \hat{y}_i , which is a 3x1 column vector. The values are 0.99280 for Setosa, 0.00024 for Versicolor, and 0.00696 for Virginica.

Legend:

- Setosa (Blue)
- Versicolor (Green)
- Virginica (Yellow)

Arrows indicate the mapping from the ground truth vector y_i to the term y_i in the equation, and from the predicted probability vector \hat{y}_i to the term \hat{y}_i in the equation.

Cross Entropy (Función de pérdida)

$$L = - \sum_{i=1}^m \boxed{y_i} \boxed{\log(\hat{y}_i)}$$

Setosa

Versicolor

Virginica

1
0
0

y_i

-0.00723
-8.32723
-4.96723

$\log(\hat{y}_i)$

Cross Entropy (Función de pérdida)

$$L = - \sum_{i=1}^m \boxed{y_i} \boxed{\log(\hat{y}_i)}$$

Setosa

Versicolor

Virginica

1
0
0

y_i

*

-0.00723
-8.32723
-4.96723

$\log(\hat{y}_i)$

=

-0.00723
0
0

Cross Entropy (Función de pérdida)

$$L = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

$$L = - \sum \begin{array}{|c|} \hline -0.00723 \\ \hline 0 \\ \hline 0 \\ \hline \end{array} = \mathbf{0.00723}$$

Review:

Entropy, Cross-entropy

Ej:

Queremos usar una antena para enviar el estado del clima.



Ej:

Queremos usar una antena para enviar el estado del clima.



code

bits

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
000	001	010	011	100
3	3	3	3	3

cantidad promedio de bits usados = 3

Ej:

Queremos usar una antena para enviar el estado del clima.



Veamos las probabilidades
del clima en una ciudad A

$p(x)$

code

bits

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
000	001	010	011	100
3	3	3	3	3

$$\text{bits promedio} = 3(0.25) + 3(0.45) + 3(0.2) + 3(0.09) + 3(0.01) = 3$$

Ej:

Queremos usar una antena para enviar el estado del clima.



$p(x)$

code

bits

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
1100	0	10	1101	1110
4	1	2	4	4

Usamos códigos más
óptimos para las
probabilidades de la ciudad A

$$\text{bits promedio} = 4(0.25) + 1(0.45) + 2(0.2) + 4(0.09) + 4(0.01) = 2.25$$

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
$p(x)$ 0.25	0.45	0.2	0.09	0.01

¿Cómo obtenemos
mínimo número de bits
que se puede usar?

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
4	2.22	5	11.11	100

¿Cómo obtenemos
mínimo número de bits
que se puede usar?

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
4	2.22	5	11.11	100
2	1.15	2.32	3.47	6.64

¿Cómo obtenemos
mínimo número de bits
que se puede usar?

← $\log_2\left(\frac{1}{p}\right) = -\log_2(p)$

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
2	1.15	2.32	3.47	6.64

¿Cómo obtenemos
mínimo número de bits
que se puede usar?

$\log_2\left(\frac{1}{p}\right) = -\log_2(p)$

bits promedio = $0.25(2) + 0.45(1.15) + 0.2(2.32) + 0.09(3.47) + 0.01(6.64) = 1.86$

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
2	1.15	2.32	3.47	6.64

¿Cómo obtenemos
mínimo número de bits
que se puede usar?

Entropy

$$-\sum p(x) \log p(x)$$



$$\text{bits promedio} = 0.25(2) + 0.45(1.15) + 0.2(2.32) + 0.09(3.47) + 0.01(6.64) = 1.86$$

Ej:

Queremos usar una antena para enviar el estado del clima.

				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
2	1.15	2.32	3.47	6.64

¿Cómo obtenemos
mínimo número de bits
que se puede usar?


Entropy

$$-\sum p(x) \log p(x)$$

La entropía en este caso es 1.86, esta cantidad de bits representa la información que uno obtendría en promedio, al conocer el clima en la ciudad A.

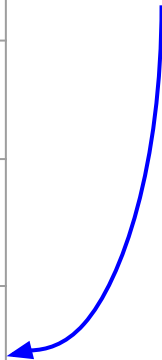
¿Qué pasa si usamos la misma antena para una ciudad B?

¿Qué pasa si usamos la misma antena para una ciudad B?



				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
2	1.15	2.32	3.47	6.64
0.3	0.2	0.2	0.25	0.05

Si la antena usa los bits óptimos para la ciudad A.

¿Cuál sería la cantidad de bits promedio que enviaría la antena en la ciudad B?



¿Qué pasa si usamos la misma antena para una ciudad B?



				
Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
0.25	0.45	0.2	0.09	0.01
2	1.15	2.32	3.47	6.64
0.3	0.2	0.2	0.25	0.05

Si la antena usa los bits óptimos para la ciudad A.

¿Cuál sería la cantidad de bits promedio que enviaría la antena en la ciudad B?

$$\text{bits promedio} = 0.3(2) + 0.2(1.15) + 0.2(2.32) + 0.25(3.47) + 0.05(6.64) = 2.5$$

¿Qué pasa si usamos la misma antena para una ciudad B?

					
	Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
$p(x)$ en A	0.25	0.45	0.2	0.09	0.01
$-\log_2(p_A)$	2	1.15	2.32	3.47	6.64
$p(x)$ en B	0.3	0.2	0.2	0.25	0.05

Si la antena usa los bits óptimos para la ciudad A.

¿Cuál sería la cantidad de bits promedio que enviaría la antena en la ciudad B?

Cross-entropy

$$-\sum p(x) \log q(x)$$

$$\text{bits promedio} = 0.3(2) + 0.2(1.15) + 0.2(2.32) + 0.25(3.47) + 0.05(6.64) = 2.5$$

¿Qué pasa si usamos la misma antena para una ciudad B?

					
	Soleado	Nublado	Nublado parcial	Lluvia	Tormenta
$p(x)$ en A	0.25	0.45	0.2	0.09	0.01
$-\log_2(p_A)$	2	1.15	2.32	3.47	6.64
$p(x)$ en B	0.3	0.2	0.2	0.25	0.05

Si la antena usa los bits óptimos para la ciudad A.

¿Cuál sería la cantidad de bits promedio que enviaría la antena en la ciudad B?

Cross-entropy

$$-\sum p(x) \log q(x)$$

La entropía cruzada en este caso es 2.5, esta cantidad de bits representa la información que uno obtendría en promedio al conocer el clima en la ciudad B usando la codificación óptima para la ciudad A.

Entropy

$$-\sum p(x) \log p(x)$$

Cross-entropy

$$-\sum p(x) \log q(x)$$

Cross Entropy (Función de pérdida)

$$L = - \sum_{i=1}^m y_i \log(\hat{y}_i) = - \sum$$

Setosa

Versicolor

Virginica

1
0
0

y_i

0.99280
0.00024
0.00696

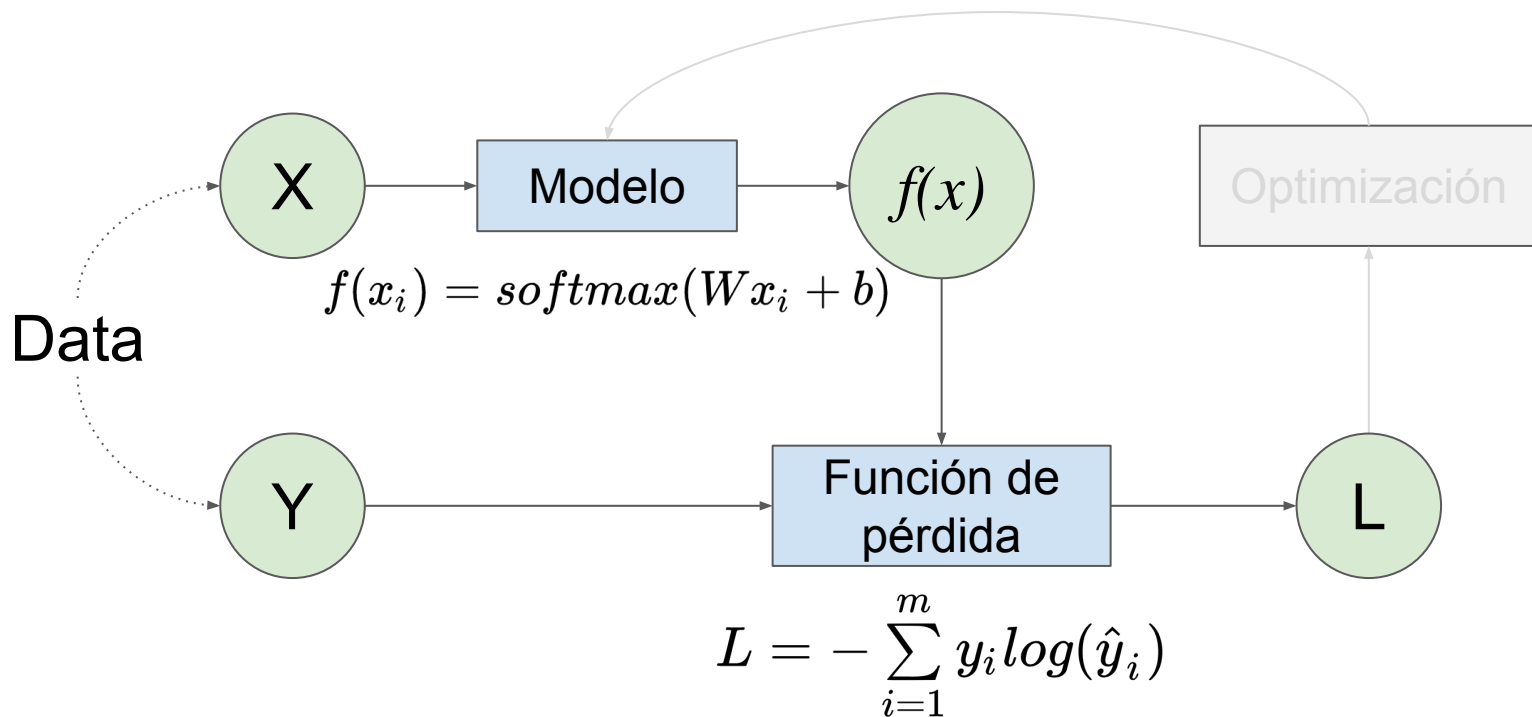
\hat{y}_i

= 0.00723

Si usamos la predicción del modelo, **0.00723** representaría la información que obtendríamos al conocer la data real.

-0.00723
0
0

Regresión logística



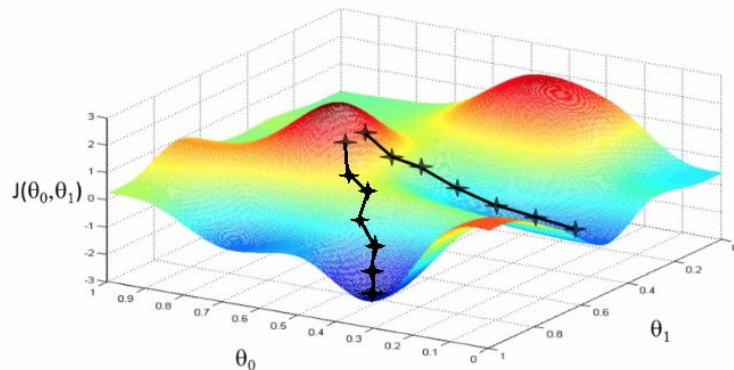
Optimización

Optimización

$$f(x_i) = \textit{softmax}(Wx_i + b)$$

$$L = - \sum_{i=1}^m y_i \log(\hat{y}_i)$$

$$\frac{dL}{dW}, \frac{dL}{db} \quad ?$$

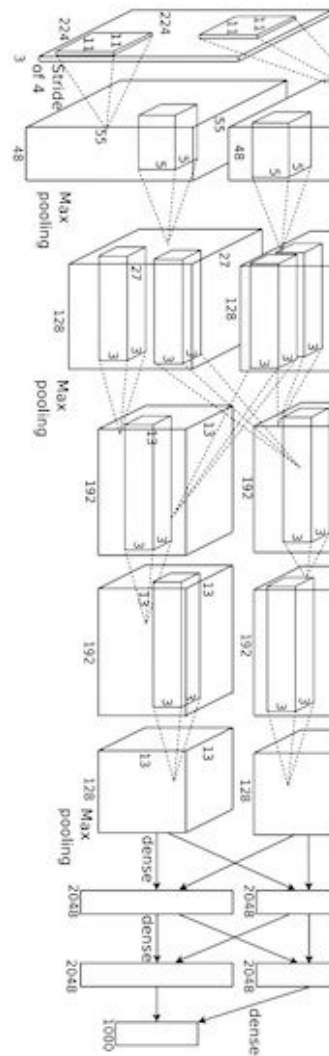


Convolutional Network (AlexNet)

input image
weights

loss

$$\frac{dL}{dW} ?$$

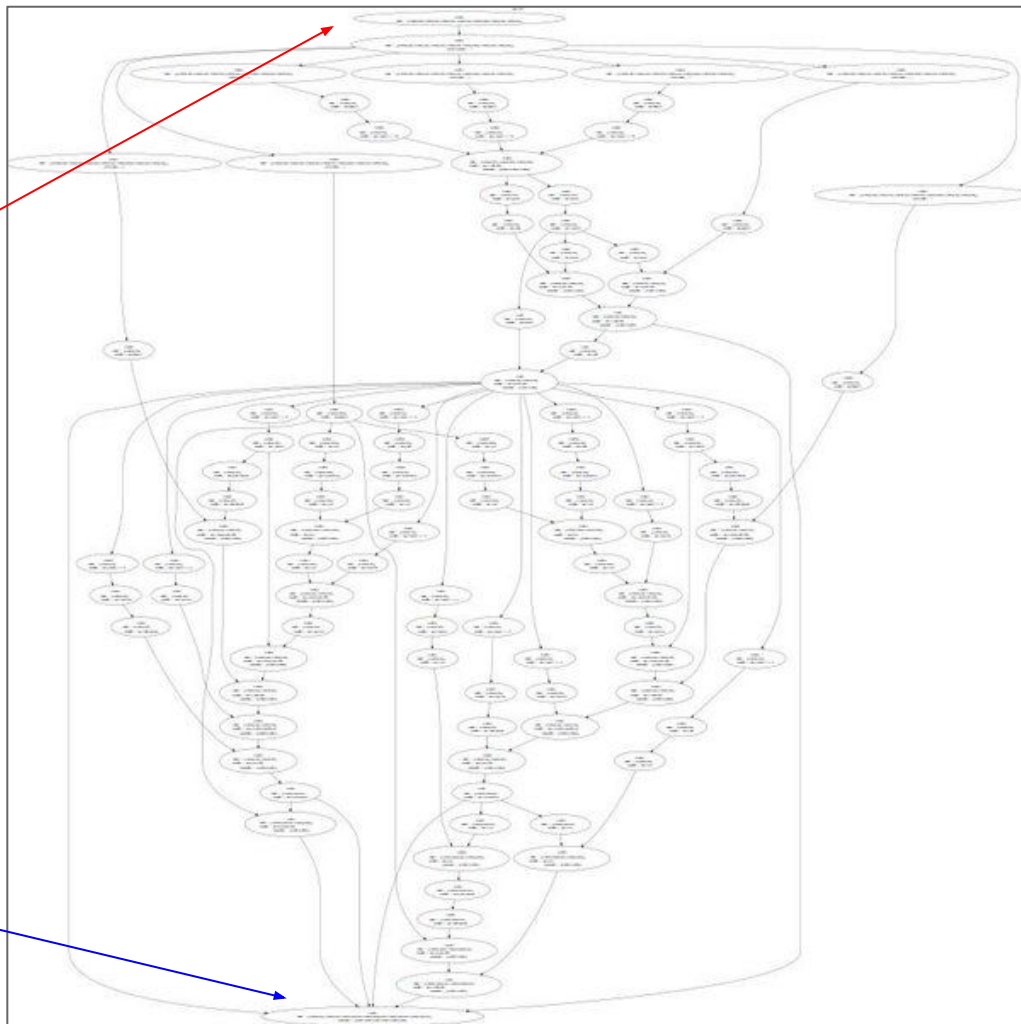


Neural Turing Machine

input tape

loss

$$\frac{dL}{dW} ?$$

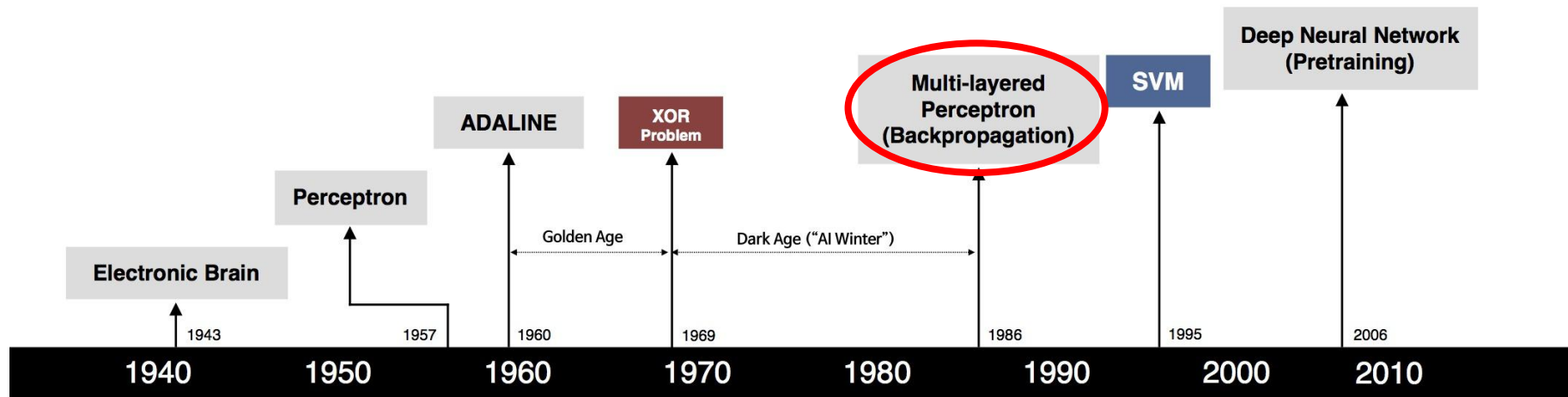


¿Cómo derivamos las gradientes?

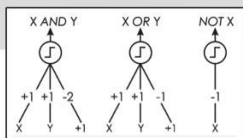


Backpropagation !





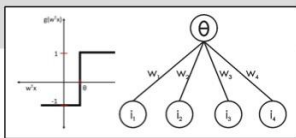
S. McCulloch – W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



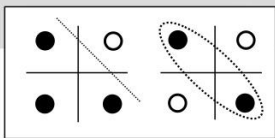
- Learnable Weights and Threshold



B. Widrow – M. Hoff



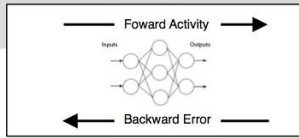
M. Minsky – S. Papert



- XOR Problem



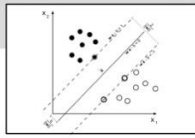
D. Rumelhart – G. Hinton – R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



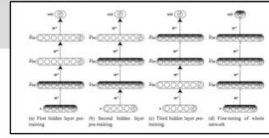
V. Vapnik – C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention



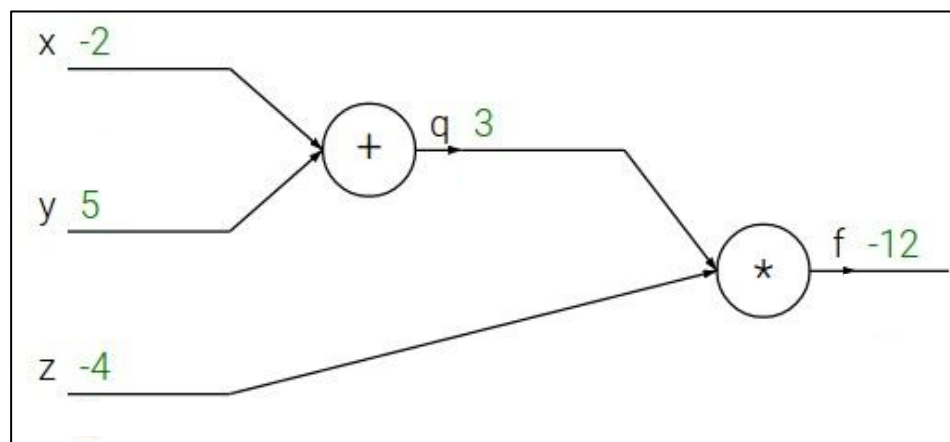
G. Hinton – S. Ruslan



- Hierarchical feature Learning

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

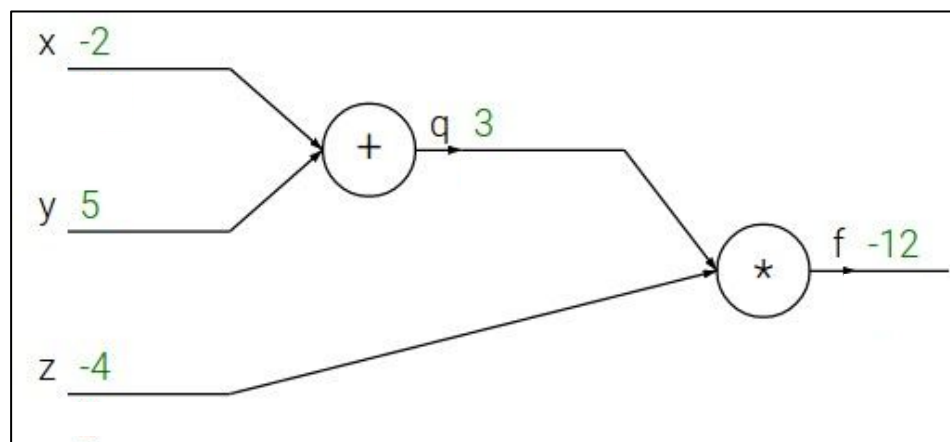


Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$



Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

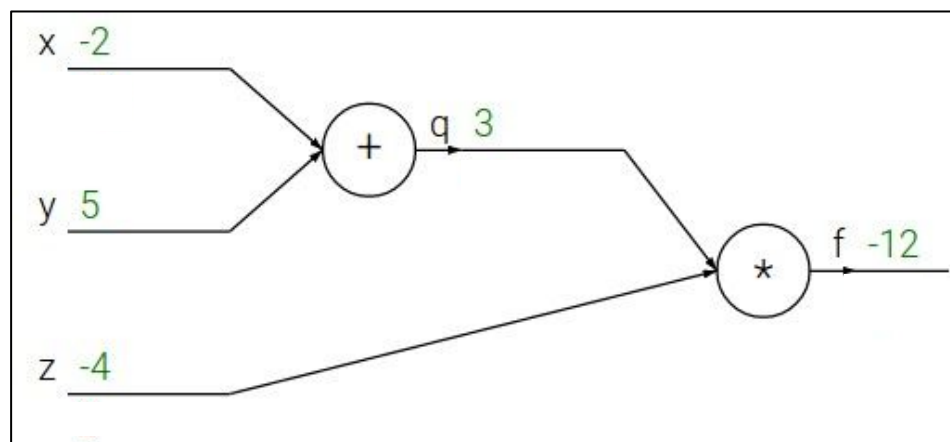
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



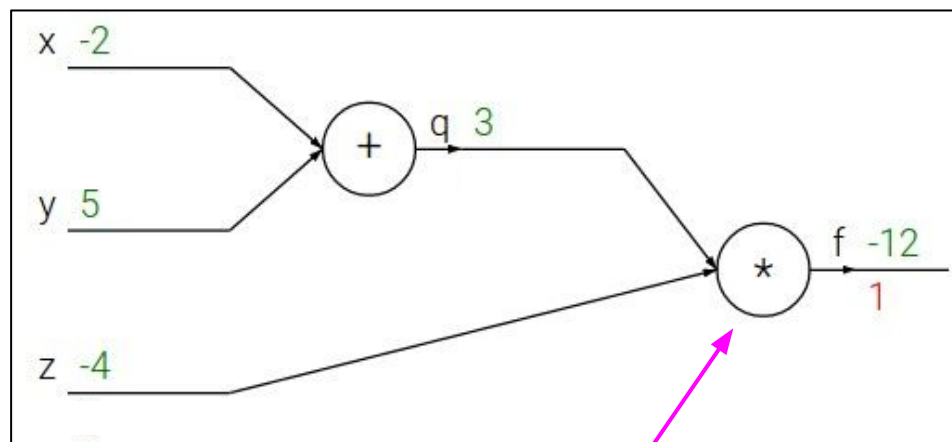
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{df}{df} = 1$$

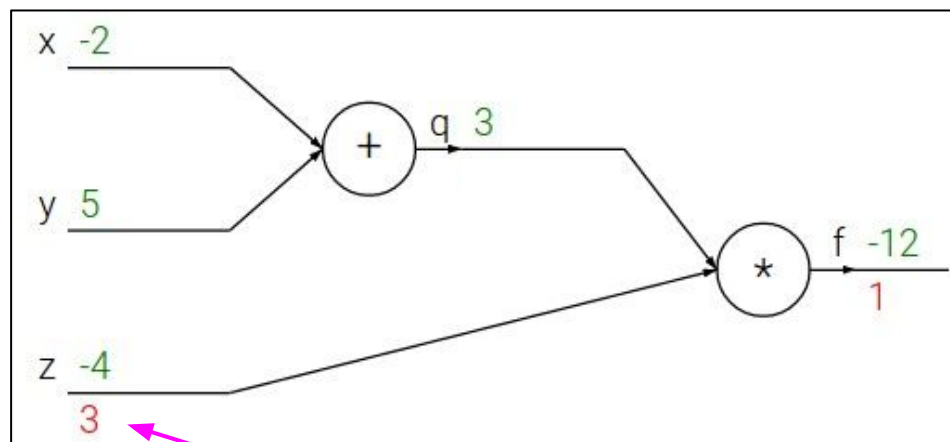
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{df}{dz} = q = 3$$

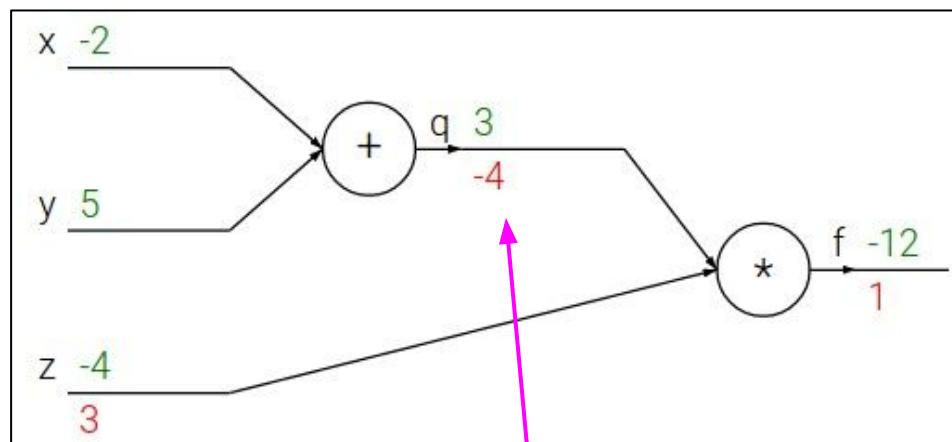
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{df}{dq} = z = -4$$

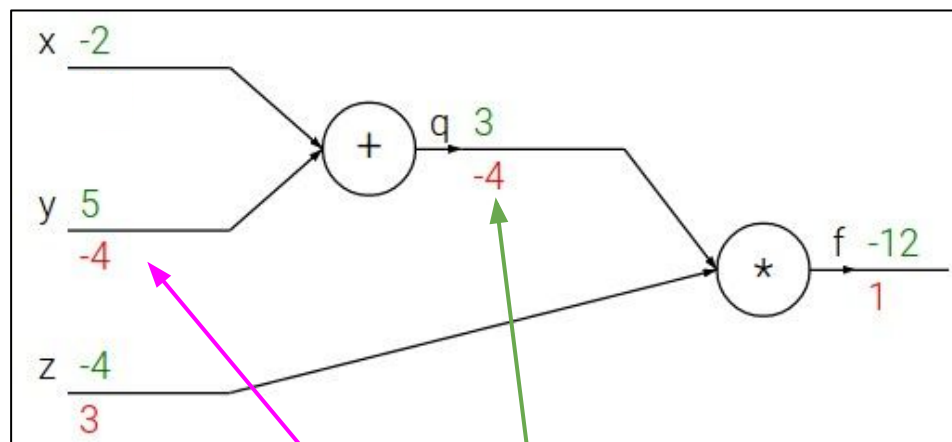
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



Chain rule:

$$\frac{df}{dy} = \frac{df}{dq} \frac{dq}{dy} = -4 * 1 = -4$$

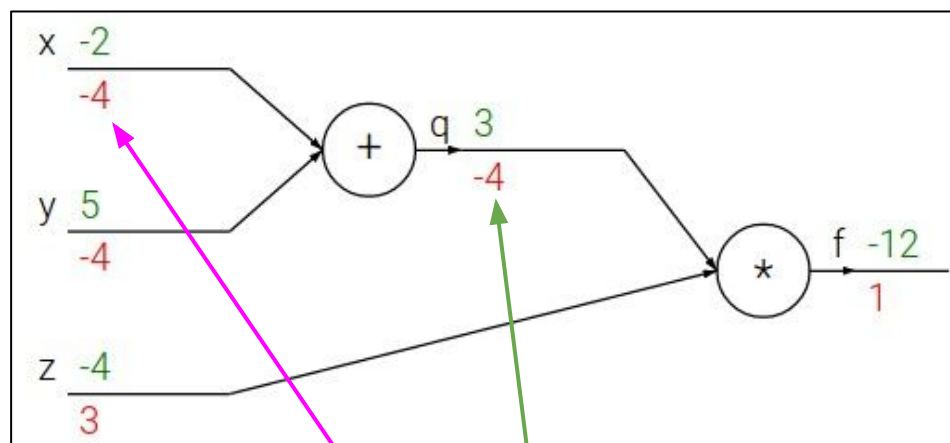
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

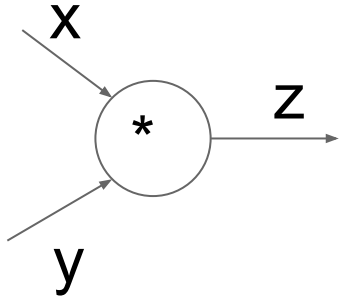
Queremos: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



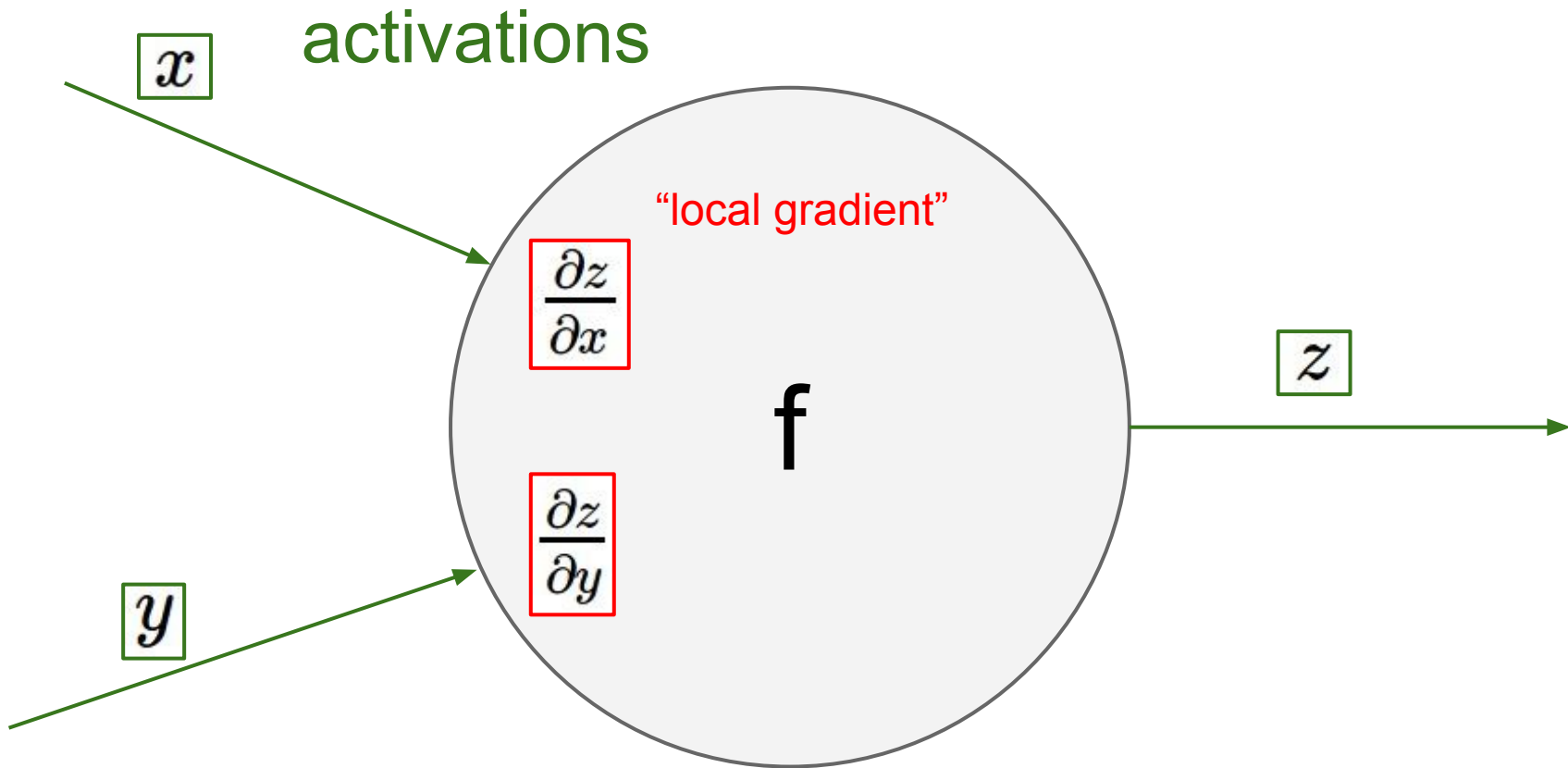
Chain rule:

$$\frac{df}{dx} = \frac{df}{dq} \frac{dq}{dx} = -4 * 1 = -4$$

Implementation: forward/backward API



```
class MultiplyGate(object):  
    def forward(x,y):  
        z = x*y  
        self.x = x # must keep these around!  
        self.y = y  
        return z  
    def backward(dz):  
        dx = self.y * dz # [dz/dx * dL/dz]  
        dy = self.x * dz # [dz/dy * dL/dz]  
        return [dx, dy]
```



Forward pass

activations

x

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial x}$$

"local gradient"

$$\frac{\partial z}{\partial x}$$

$$\frac{\partial z}{\partial y}$$

f

z

$$\frac{\partial L}{\partial z}$$

gradients

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

y

Backward pass

Example: Torch MulConstant

$$f(X) = aX$$

initialization

forward()

backward()

```
1 local MulConstant, parent = torch.class('nn.MulConstant', 'nn.Module')
2
3 function MulConstant:__init(constant_scalar, ip)
4     parent.__init(self)
5     assert(type(constant_scalar) == 'number', 'input is not scalar!')
6     self.constant_scalar = constant_scalar
7
8     -- default for inplace is false
9     self.inplace = ip or false
10    if (ip and type(ip) ~= 'boolean') then
11        error('in-place flag must be boolean')
12    end
13 end
```

```
14
15 function MulConstant:updateOutput(input)
16     if self.inplace then
17         input:mul(self.constant_scalar)
18         self.output = input
19     else
20         self.output:resizeAs(input)
21         self.output:copy(input)
22         self.output:mul(self.constant_scalar)
23     end
24     return self.output
25 end
```

```
26
27 function MulConstant:updateGradInput(input, gradOutput)
28     if self.gradInput then
29         if self.inplace then
30             gradOutput:mul(self.constant_scalar)
31             self.gradInput = gradOutput
32             -- restore previous input value
33             input:div(self.constant_scalar)
34         else
35             self.gradInput:resizeAs(gradOutput)
36             self.gradInput:copy(gradOutput)
37             self.gradInput:mul(self.constant_scalar)
38         end
39         return self.gradInput
40     end
41 end
```

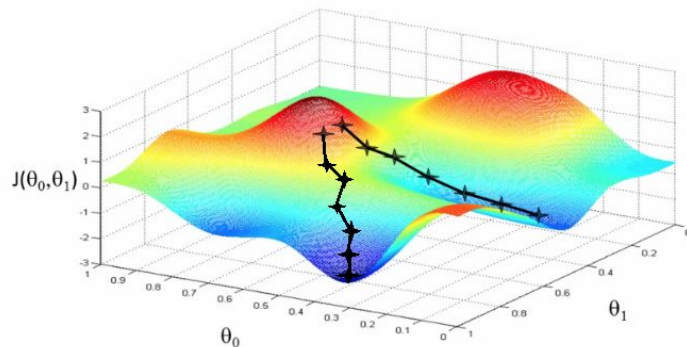
TensorFlow Layers

[illegible][illegible]

Optimización: Gradient descent

Iterar:

1. **Forward:** obtener la pérdida (Loss).
2. **Backprop:** calcular las gradientes.
3. **Actualizar** los parámetros del modelo.

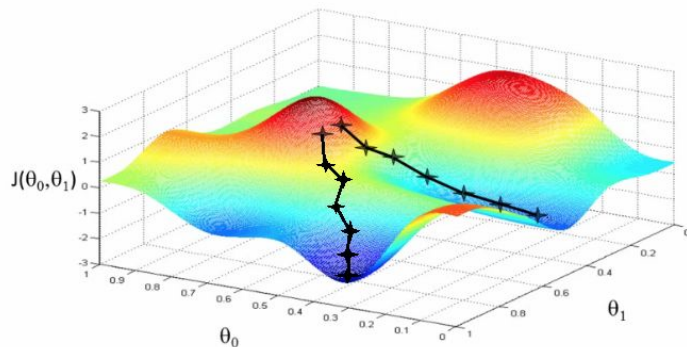


Optimización: Gradient descent

Iterar:

1. **Forward:** obtener la pérdida (Loss).
2. **Backprop:** calcular las gradientes.
3. **Actualizar** los parámetros del modelo.

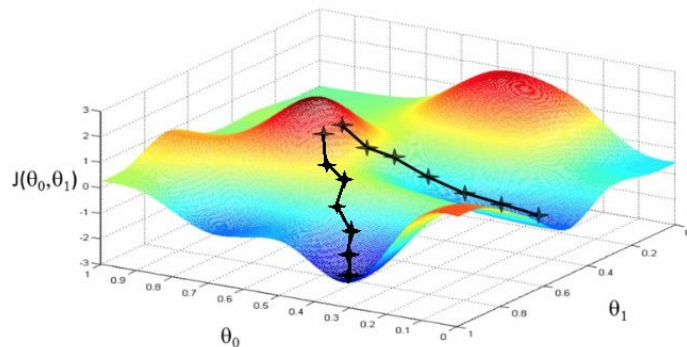
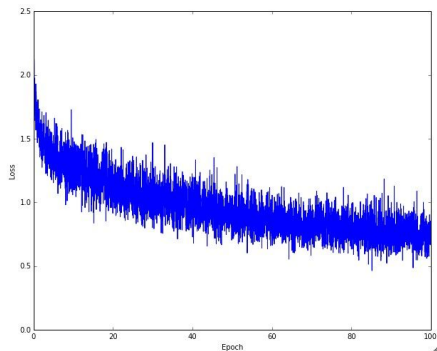
¿Y si tenemos mucha data?



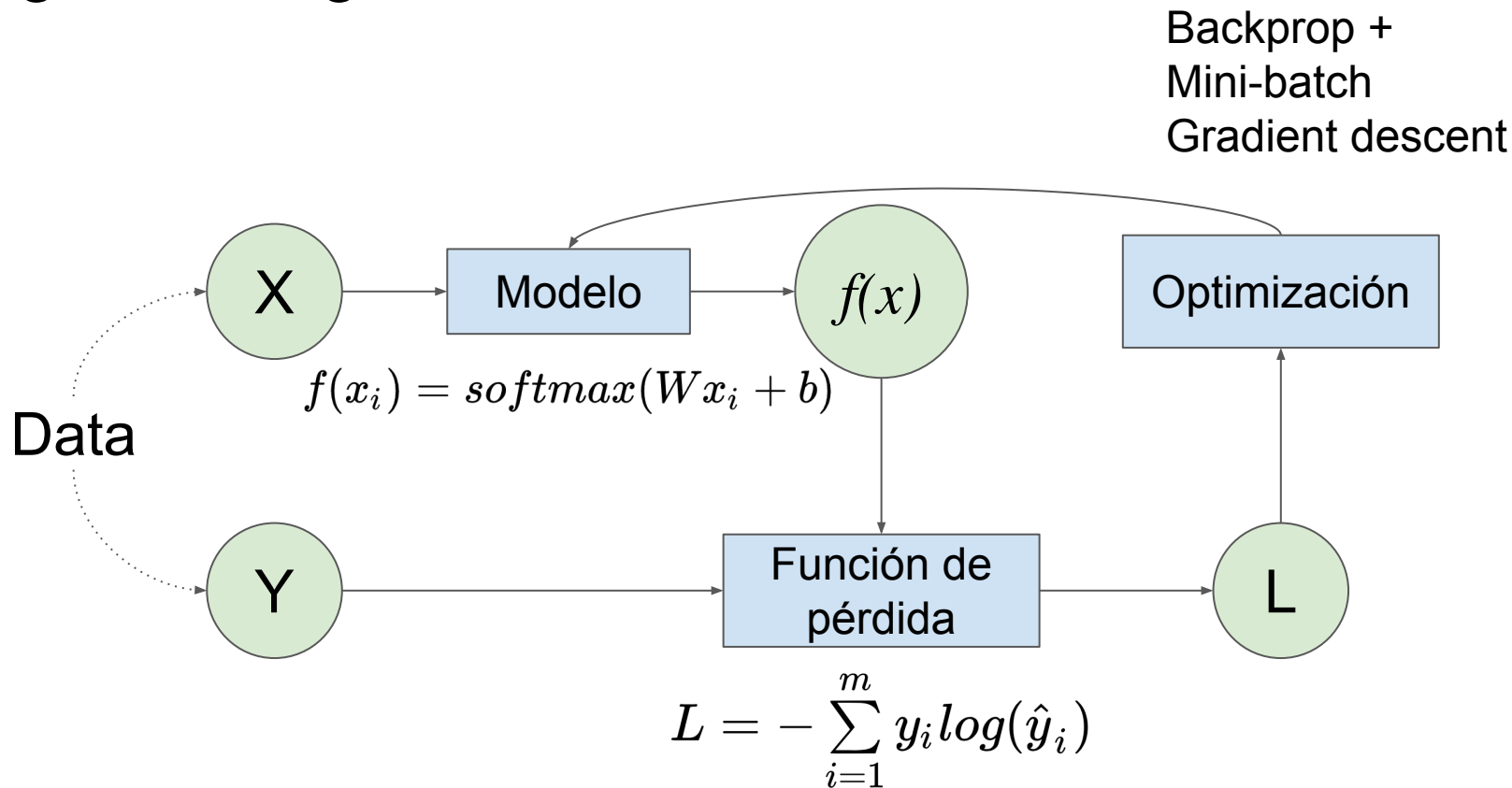
Optimización: Mini-batch Gradient descent

Iterar:

1. **Sample:** obtener una muestra de la data.
2. **Forward:** obtener la pérdida (Loss).
3. **Backprop:** calcular las gradientes.
4. **Actualizar** los parámetros del modelo.



Regresión logística



Redes Neuronales

Redes Neuronales

Hasta ahora teníamos funciones lineales antes de aplicar softmax:

$$f(x_i) = \textit{softmax}(Wx_i + b)$$

¿Cómo incrementamos la complejidad del modelo?

Redes Neuronales

Hasta ahora teníamos funciones lineales antes de aplicar softmax:

$$f(x_i) = \text{softmax}(Wx_i + b)$$

¿Cómo incrementamos la complejidad del modelo?

Podemos añadir operaciones intermedias

Redes Neuronales

$$\boxed{h(x_i)} = \boxed{W_1} \boxed{x_i} + \boxed{b_1}$$

10xN **10x4** **4xN** **10x1**

Redes Neuronales

$$\boxed{h(x_i)} = \boxed{W_1} \boxed{x_i} + \boxed{b_1}$$

$10 \times N$ 10×4 $4 \times N$ 10×1

$$\boxed{f(x_i)} = \text{softmax}(\boxed{W_2} \boxed{h(x_i)} + \boxed{b_2})$$

$3 \times N$ 3×10 $10 \times N$ 3×1

Redes Neuronales

$$h_1(x_i) = W_1 x_i + b_1$$

$$h_2(x_i) = W_2 h_1(x_i) + b_2$$

⋮

$$f(x_i) = \textit{softmax}(W_m h_n(x_i) + b_m)$$

Redes Neuronales

$$h_1(x_i) = W_1 x_i + b_1$$

$$h_2(x_i) = W_2 h_1(x_i) + b_2$$

⋮

$$f(x_i) = \textit{softmax}(W_m h_n(x_i) + b_m)$$

Pero una serie de operaciones lineales se puede reducir a solo 1.

$$y_1 = 2x + 10$$

$$y_2 = 3y_1 - 5$$

$$y_2 = 6x + 25$$

Redes Neuronales

Pero una serie de operaciones lineales se puede reducir a solo 1.

$$h_1(x_i) = \sigma(W_1 x_i + b_1)$$

$$h_2(x_i) = \sigma(W_2 h_1(x_i) + b_2)$$

⋮

$$f(x_i) = \textit{softmax}(W_m h_n(x_i) + b_m)$$

**Añadimos
funciones de
activación**

Redes Neuronales

Pero una serie de operaciones lineales se puede reducir a solo 1.

$$h_1(x_i) = \sigma(W_1 x_i + b_1)$$

$$h_2(x_i) = \sigma(W_2 h_1(x_i) + b_2)$$

⋮

$$f(x_i) = \underline{\text{softmax}}(W_m h_n(x_i) + b_m)$$

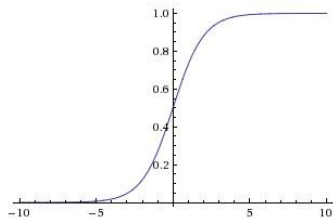
softmax también es una función de activación

**Añadimos
funciones de
activación**

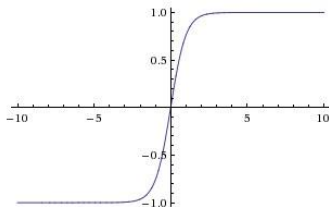
Funciones de Activación

Sigmoid

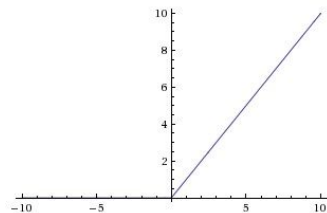
$$\sigma(x) = 1/(1 + e^{-x})$$



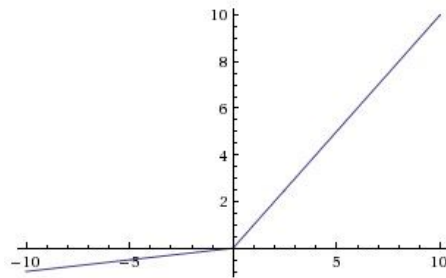
tanh tanh(x)



ReLU max(0,x)

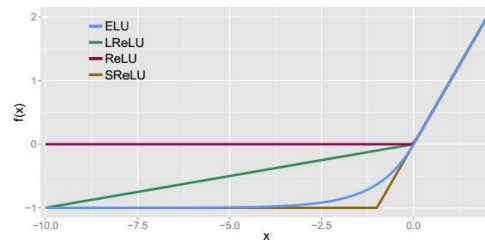


Leaky ReLU max(0.1x, x)



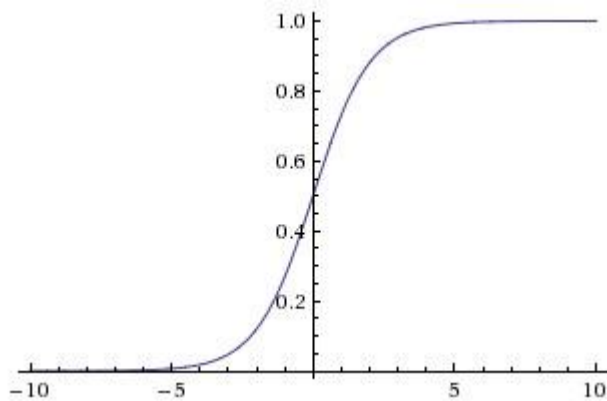
ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



Funciones de Activación

$$\sigma(x) = 1/(1 + e^{-x})$$

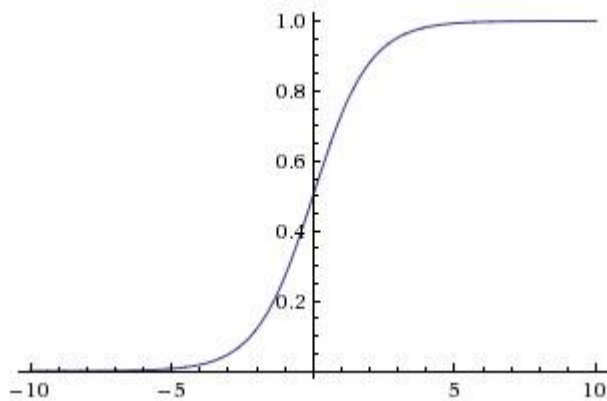


Sigmoid

- Reduce el rango a $[0,1]$
- Históricamente popular, por la interpretación de “activar” una neurona.

Funciones de Activación

$$\sigma(x) = 1/(1 + e^{-x})$$

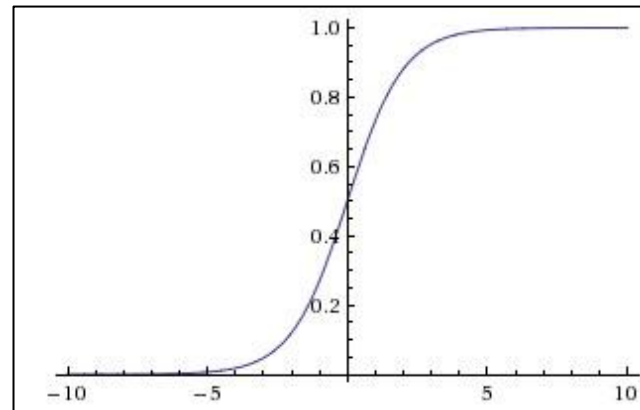
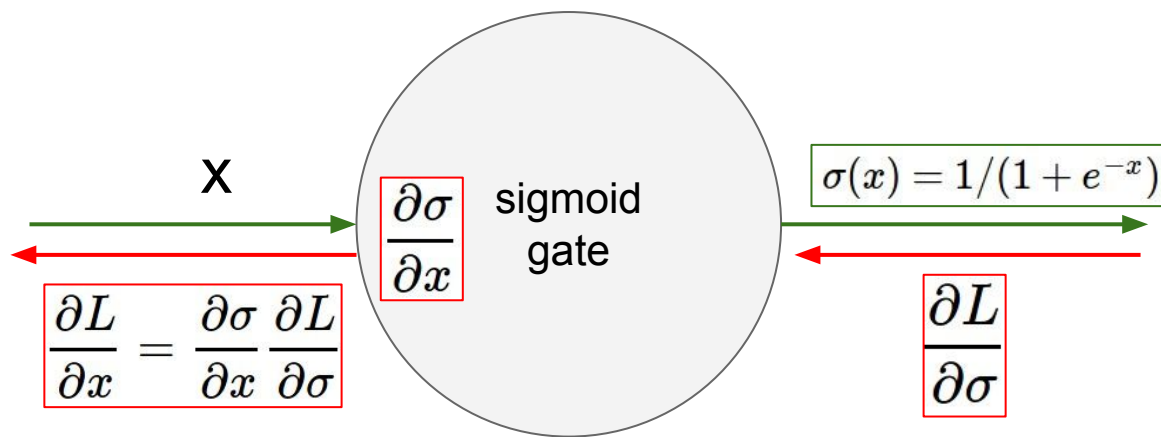


Sigmoid

- Reduce el rango a $[0,1]$
- Históricamente popular, por la interpretación de “activar” una neurona.

3 problemas:

1. Las neuronas saturadas “matan” las gradientes.

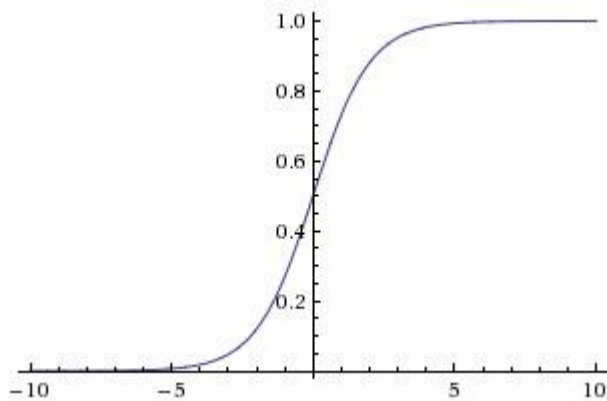


¿Qué pasa cuando $x = -10$?

¿Qué pasa cuando $x = 10$?

Funciones de Activación

$$\sigma(x) = 1/(1 + e^{-x})$$



Sigmoid

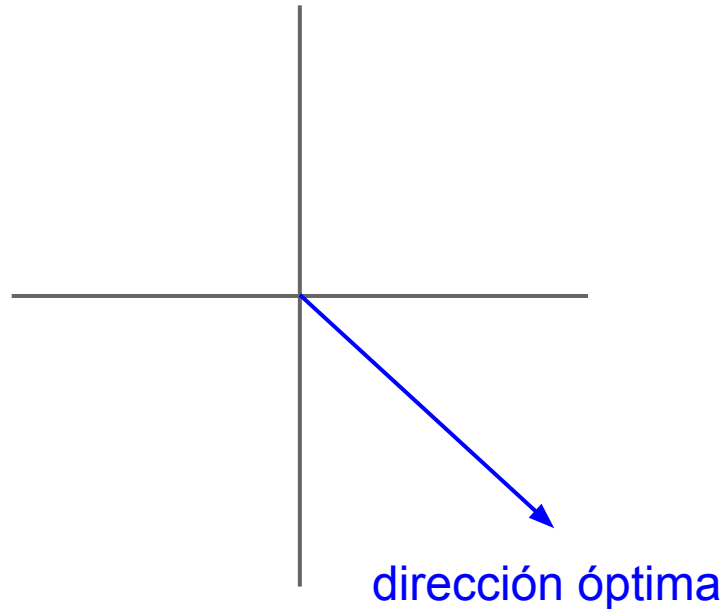
- Reduce el rango a $[0, 1]$
- Históricamente popular, por la interpretación de “activar” una neurona.

3 problemas:

1. Las neuronas saturadas “matan” los gradientes.
2. El outputs no está centrado en cero.

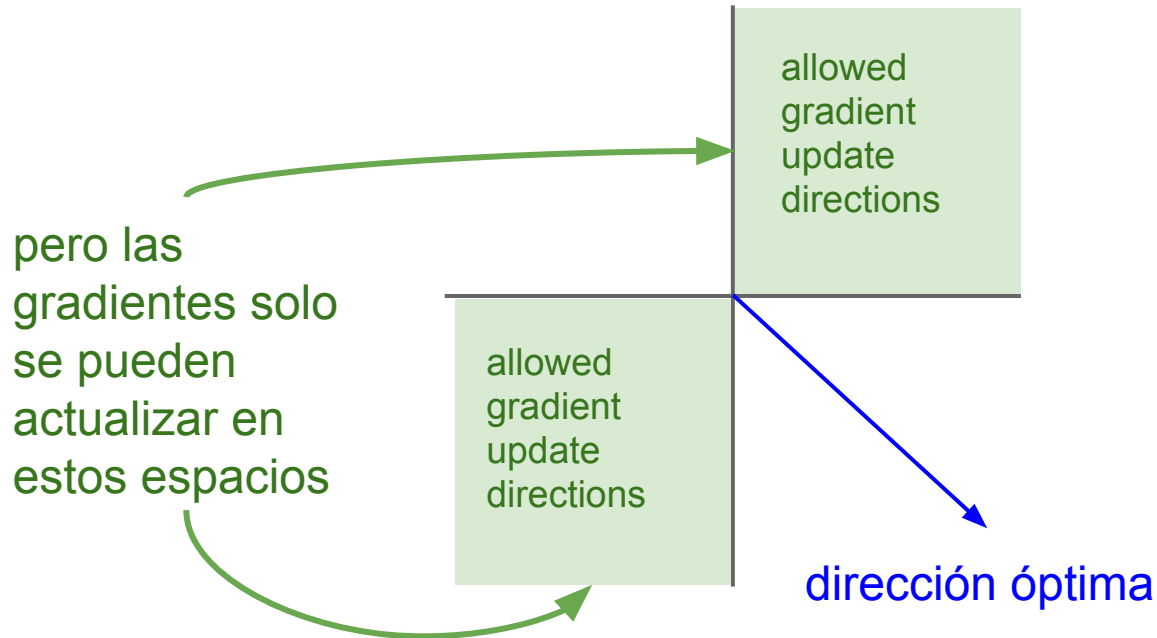
Cuando el input (x) de una neurona siempre es positivo, las gradientes van a ser todas positivas o todas negativas.

Espacio de optimización de la gradiente:



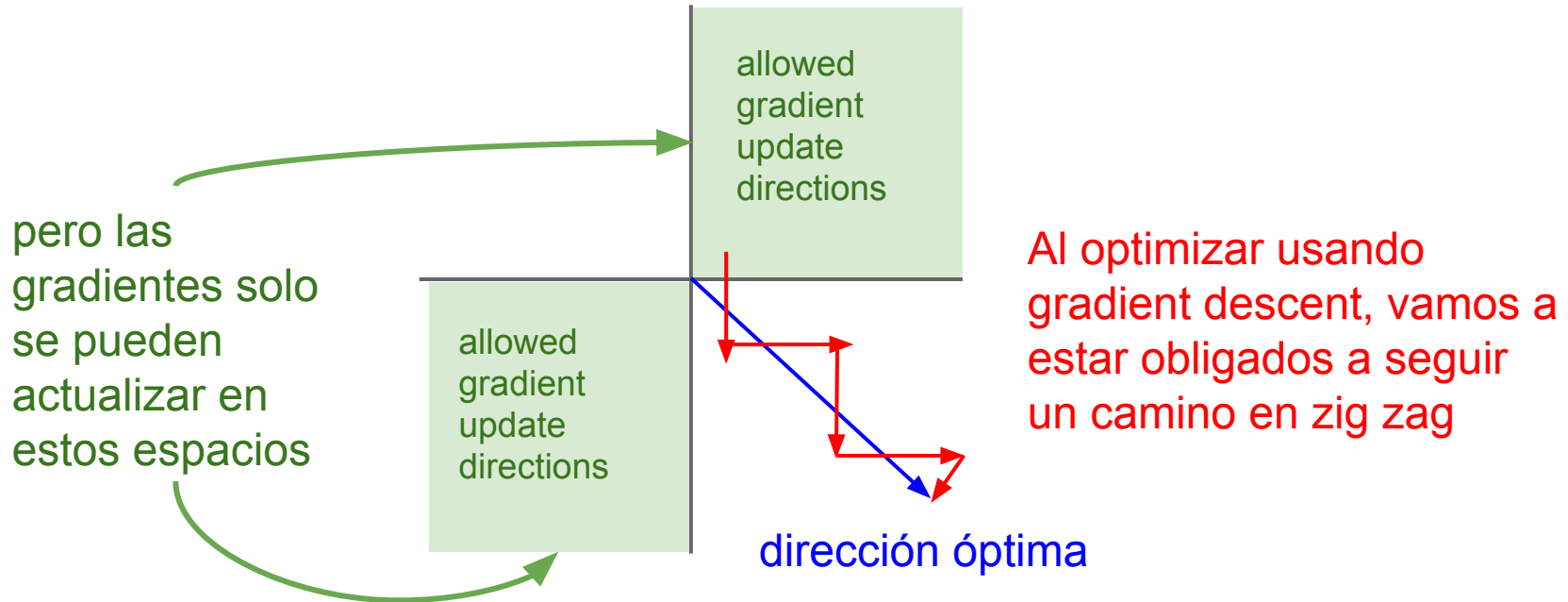
Cuando el input (x) de una neurona siempre es positivo, las gradientes van a ser todas positivas o todas negativas.

Espacio de optimización de la gradiente:



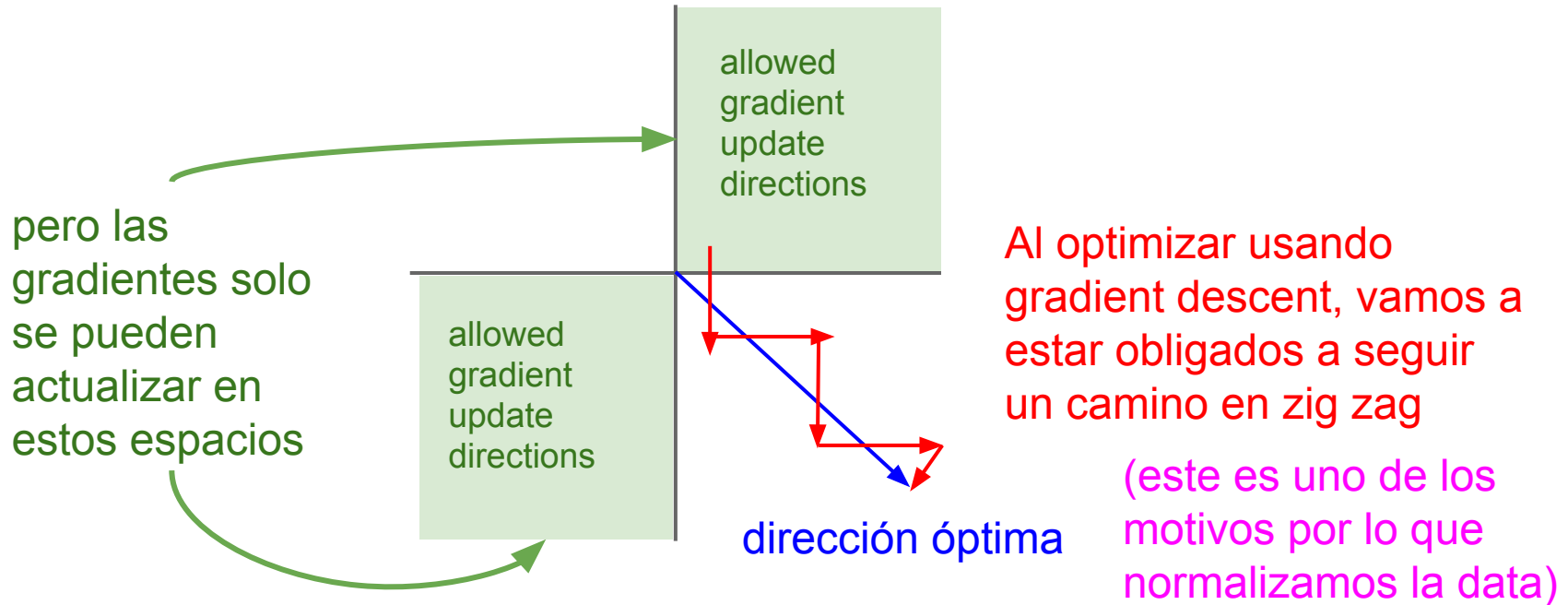
Cuando el input (x) de una neurona siempre es positivo, las gradientes van a ser todas positivas o todas negativas.

Espacio de optimización de la gradiente:



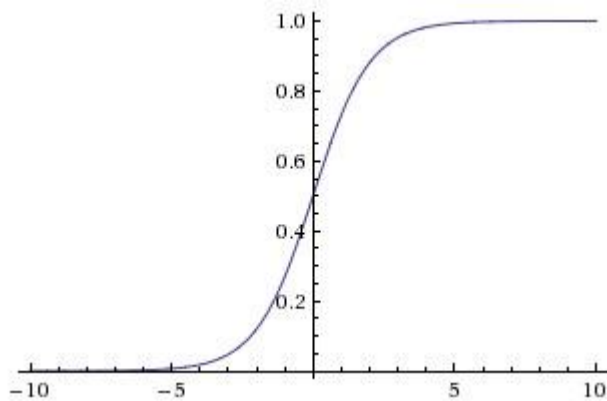
Cuando el input (x) de una neurona siempre es positivo, las gradientes van a ser todas positivas o todas negativas.

Espacio de optimización de la gradiente:



Funciones de Activación

$$\sigma(x) = 1/(1 + e^{-x})$$



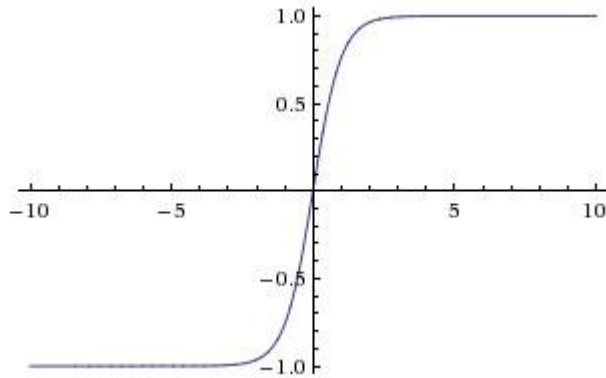
Sigmoid

- Reduce el rango a $[0,1]$
- Históricamente popular, por la interpretación de “activar” una neurona.

3 problemas:

1. Las neuronas saturadas “matan” las gradientes.
2. El outputs no está centrado en cero.
3. La función exponencial es pesada de calcular.

Funciones de Activación

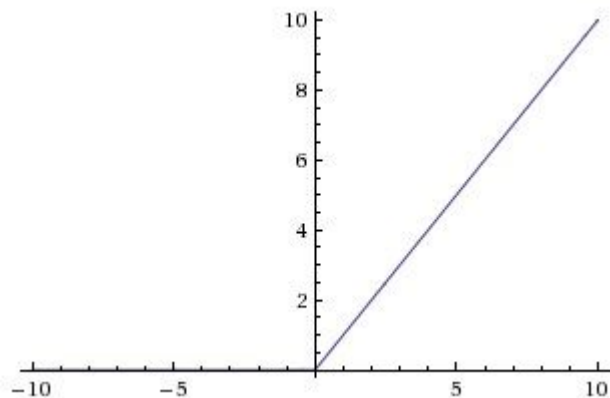


- Reduce el rango a $[-1,1]$
- Centrada en cero :D
- Neuronas saturadas :(

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

Funciones de Activación

$$\text{relu}(x) = \max(0, x)$$



- No se satura (en los positivos).
- Computacionalmente eficiente.
- No está centrada en cero.
- Las gradientes cuando $x < 0$ es ?

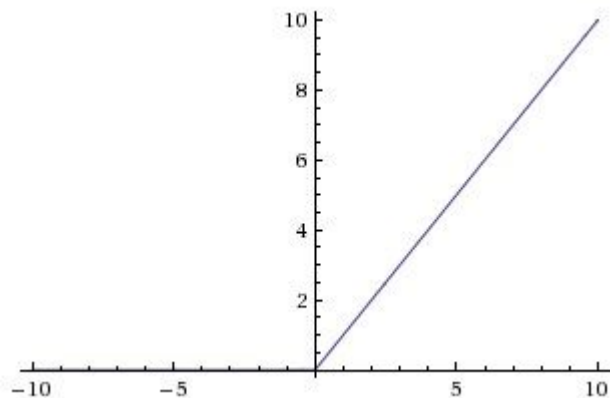
ReLU

(Rectified Linear Unit)

[Krizhevsky et al., 2012]

Funciones de Activación

$$\text{relu}(x) = \max(0, x)$$



- No se satura (en los positivos).
- Computacionalmente eficiente.
- No está centrada en cero.
- Las gradientes cuando $x < 0$ es 0.

ReLU

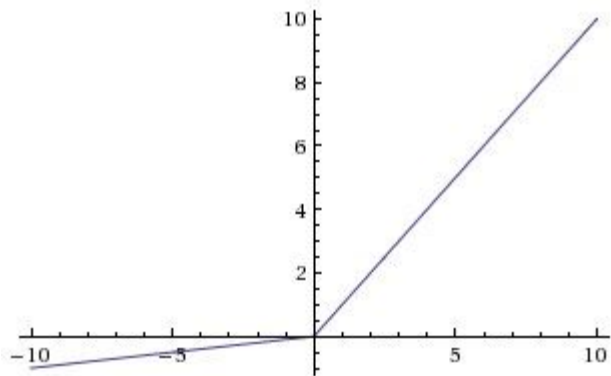
(Rectified Linear Unit)

[Krizhevsky et al., 2012]

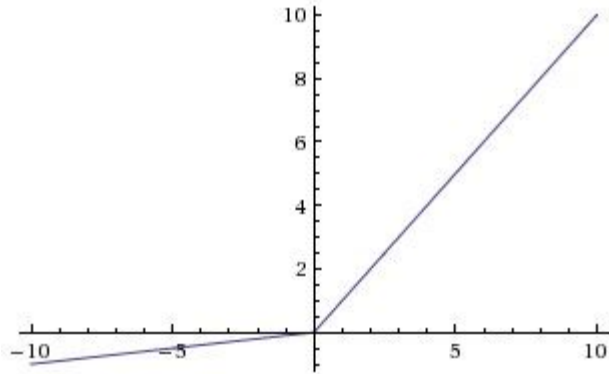
Funciones de Activación

Leaky ReLU

$$f(x) = \max(0.01x, x)$$



Funciones de Activación



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

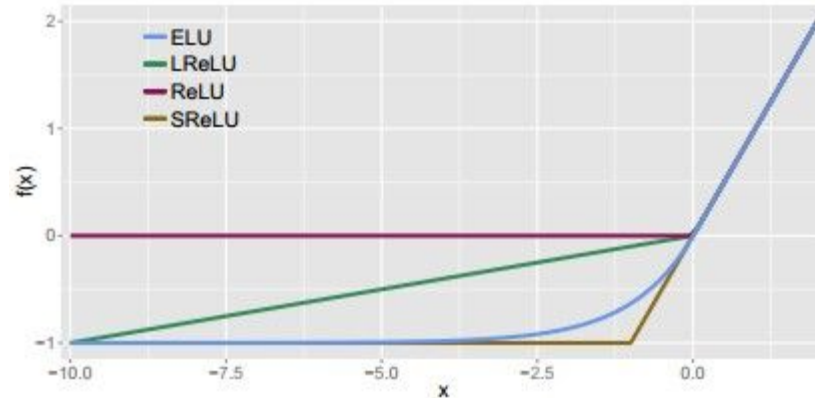
alpha es un parámetro a optimizar

[Mass et al., 2013]

[He et al., 2015]

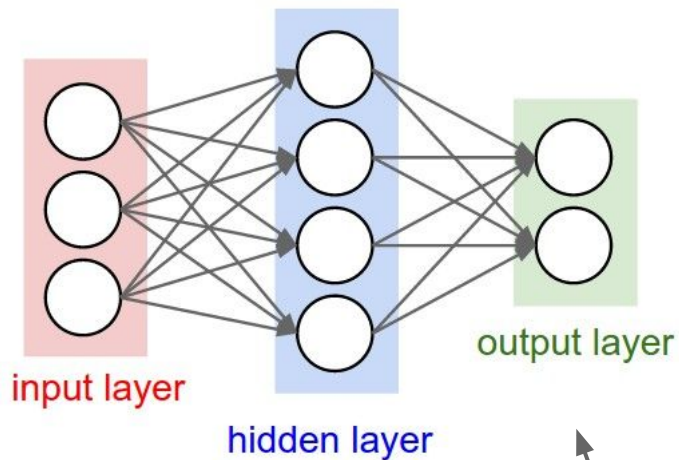
Funciones de Activación

Exponential Linear Units (ELU)

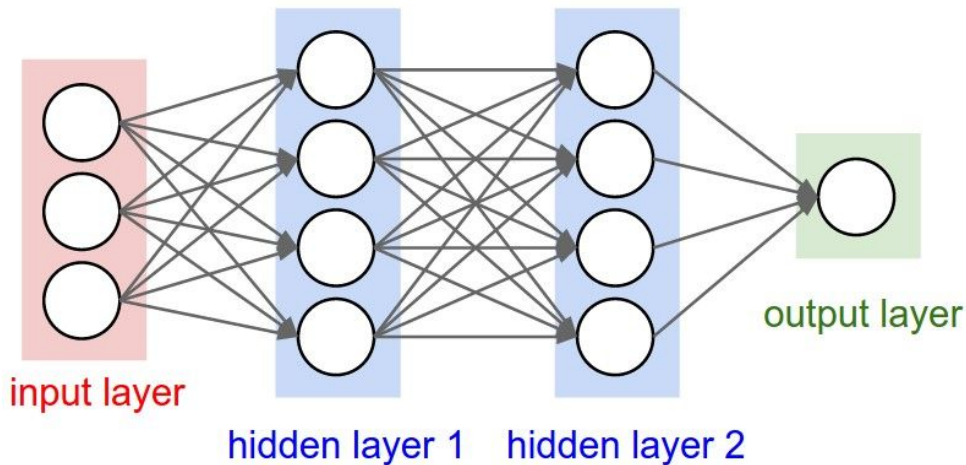


$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

Redes Neuronales



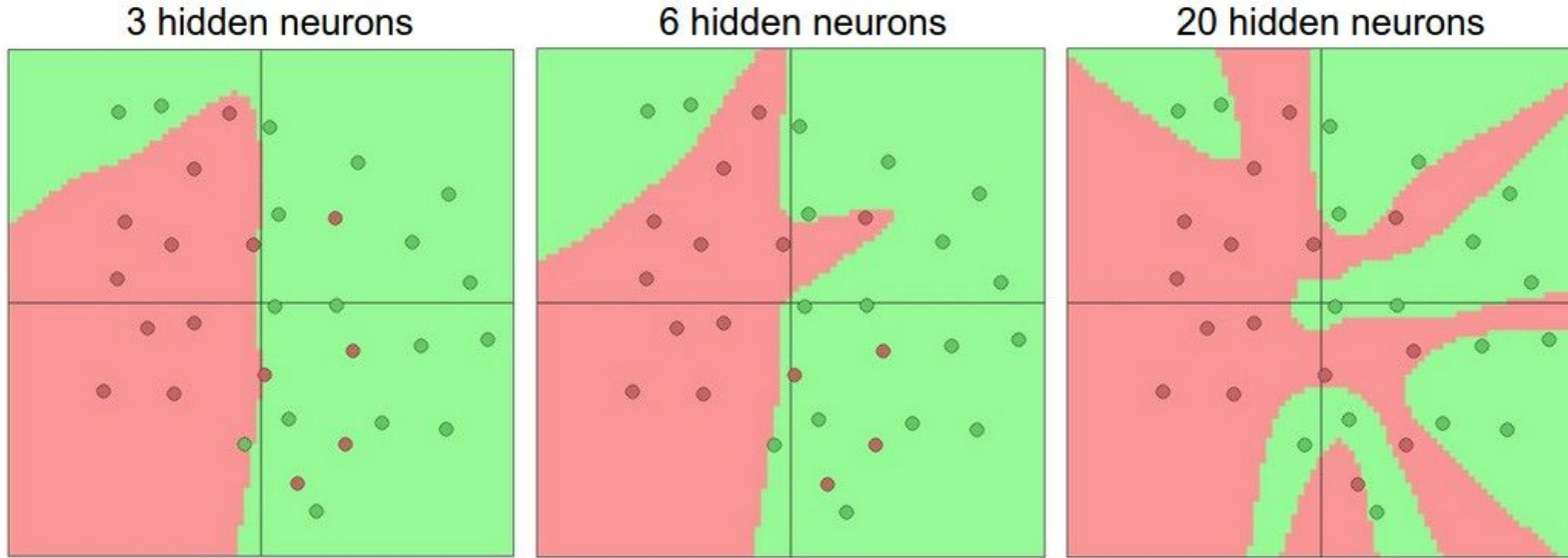
“2-layer Neural Net”, or
“1-hidden-layer Neural Net”



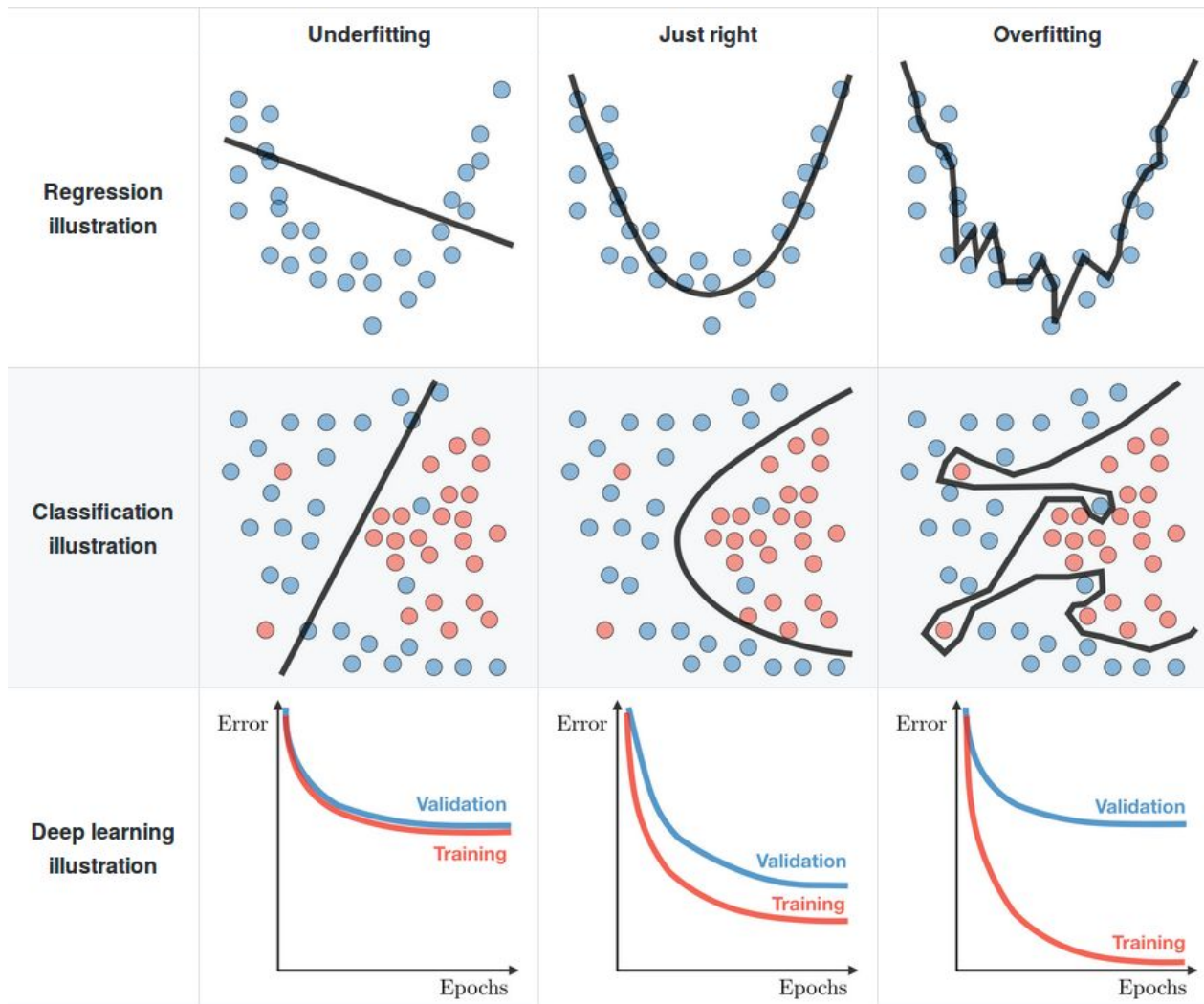
“3-layer Neural Net”, or
“2-hidden-layer Neural Net”

“Fully-connected” layers

¿Número de capas, de neuronas?



↑
de neuronas?.... *hiperparametros*



	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none"> • High training error • Training error close to test error • High bias 	<ul style="list-style-type: none"> • Training error slightly lower than test error 	<ul style="list-style-type: none"> • Very low training error • Training error much lower than test error • High variance
Possible remedies	<ul style="list-style-type: none"> • Complexify model • Add more features • Train longer 		<ul style="list-style-type: none"> • Perform regularization • Get more data

Redes Neuronales

