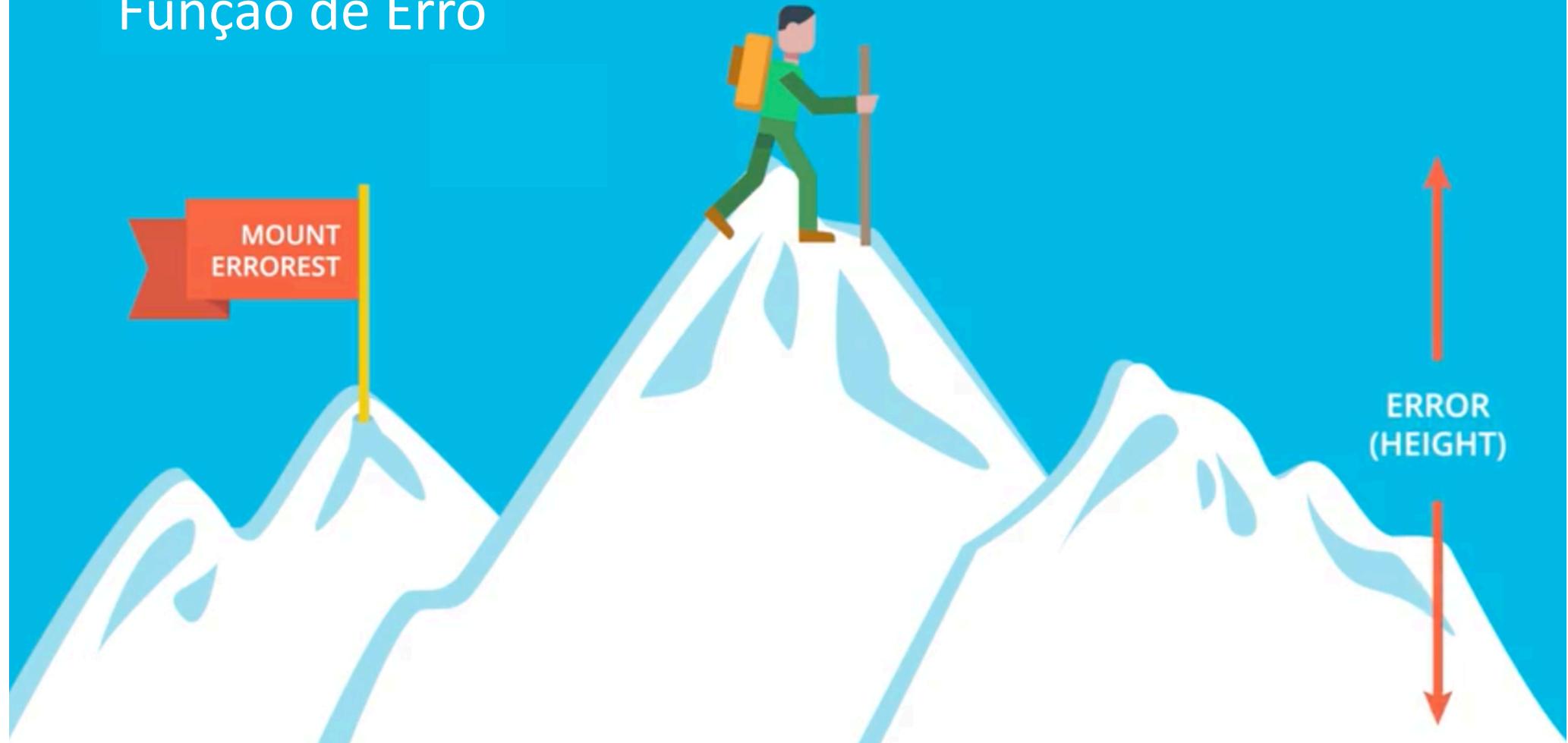


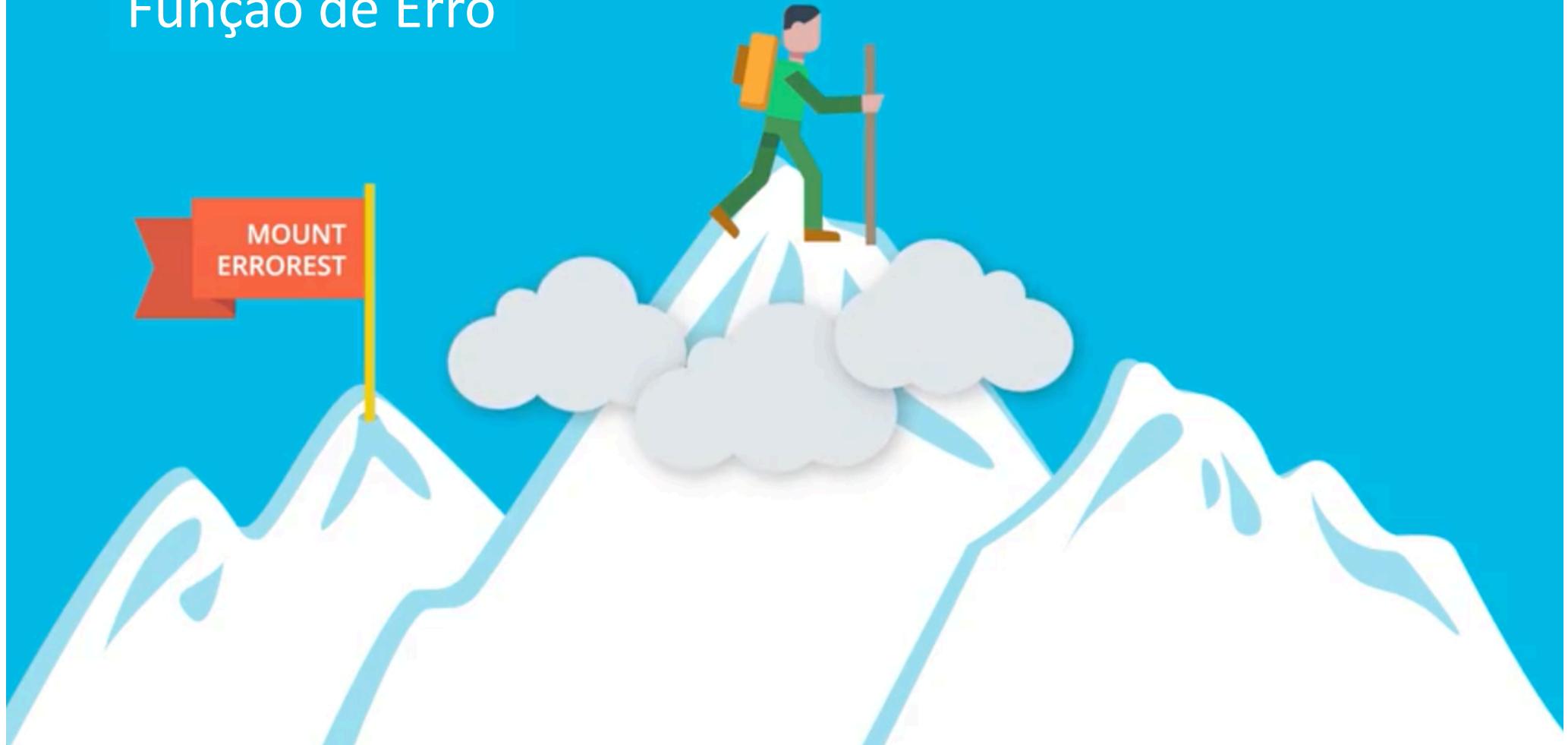
CURSO PRÁTICO DE DEEP LEARNING COM PYTHON E KERAS

Redes com Múltiplas Camadas =S

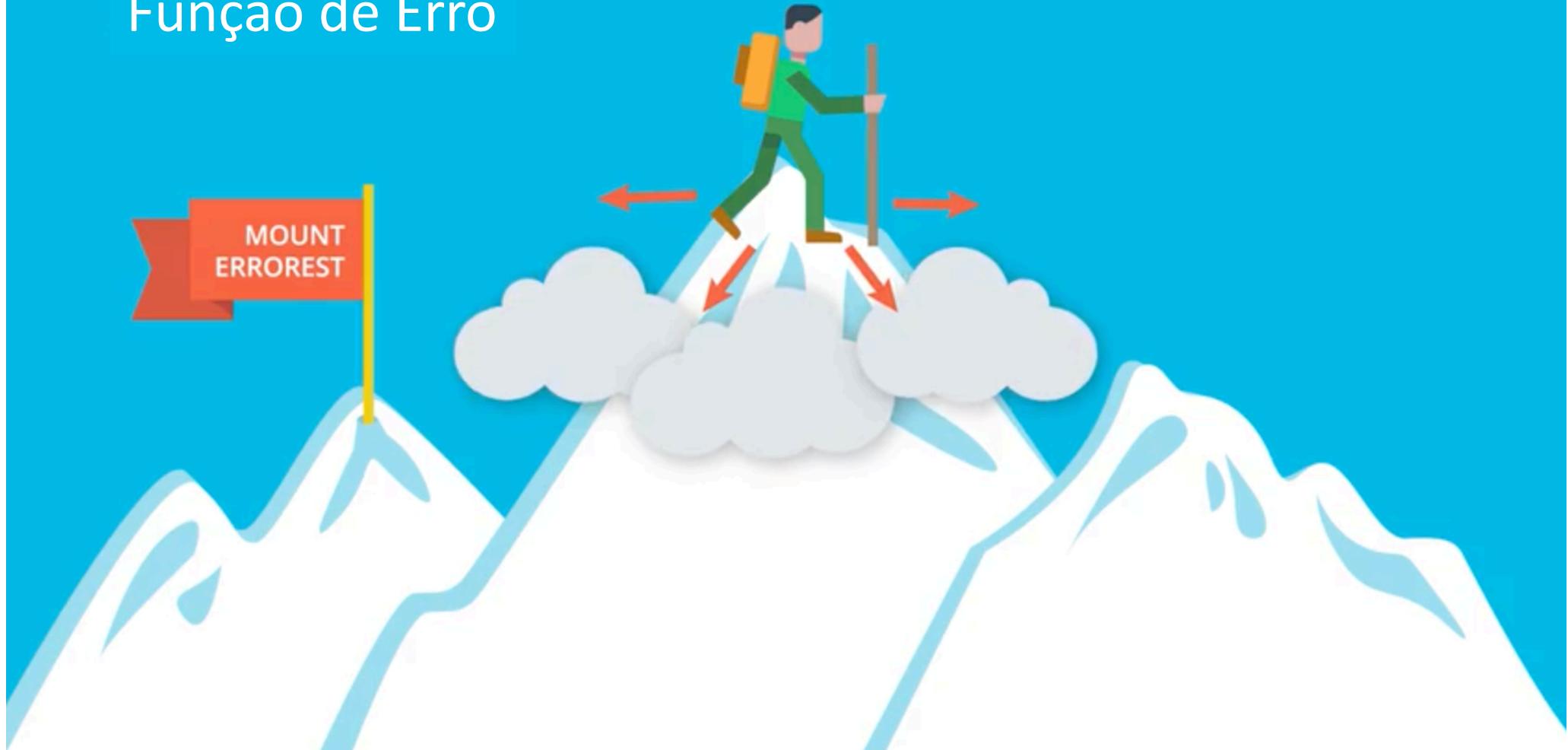
Função de Erro



Função de Erro



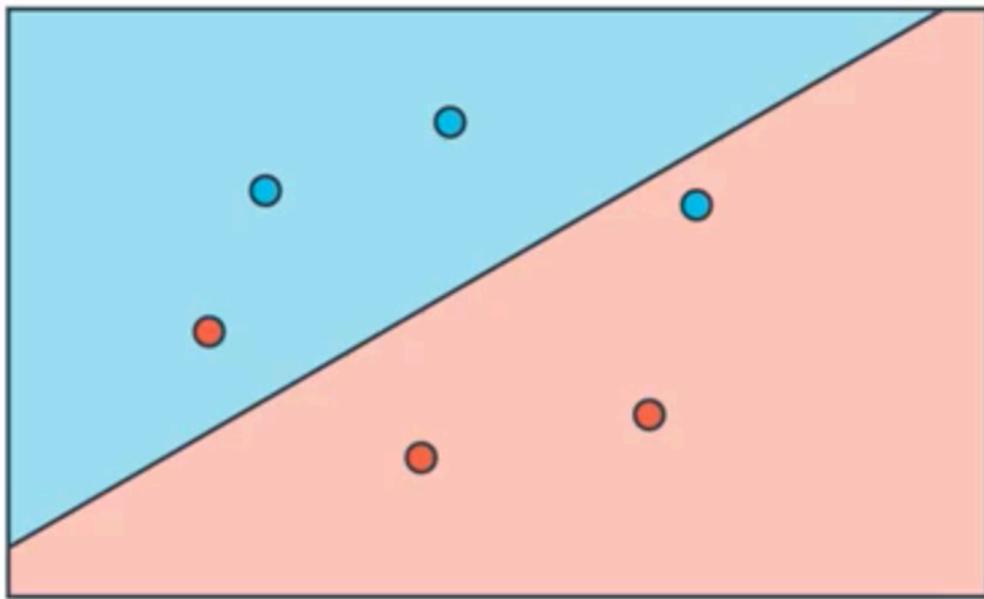
Função de Erro



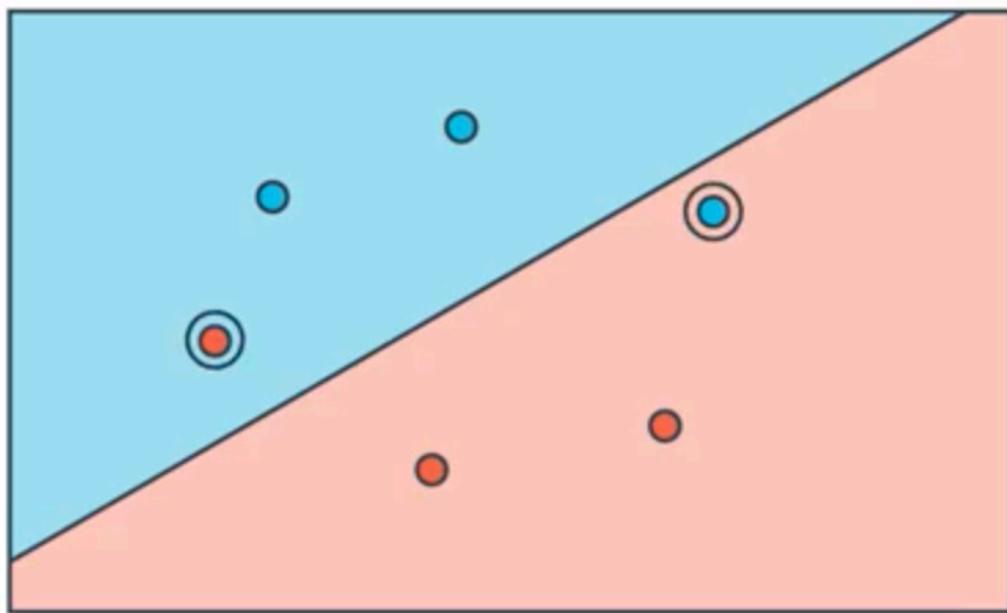
Função de Erro



▷ O quão errado é o modelo?

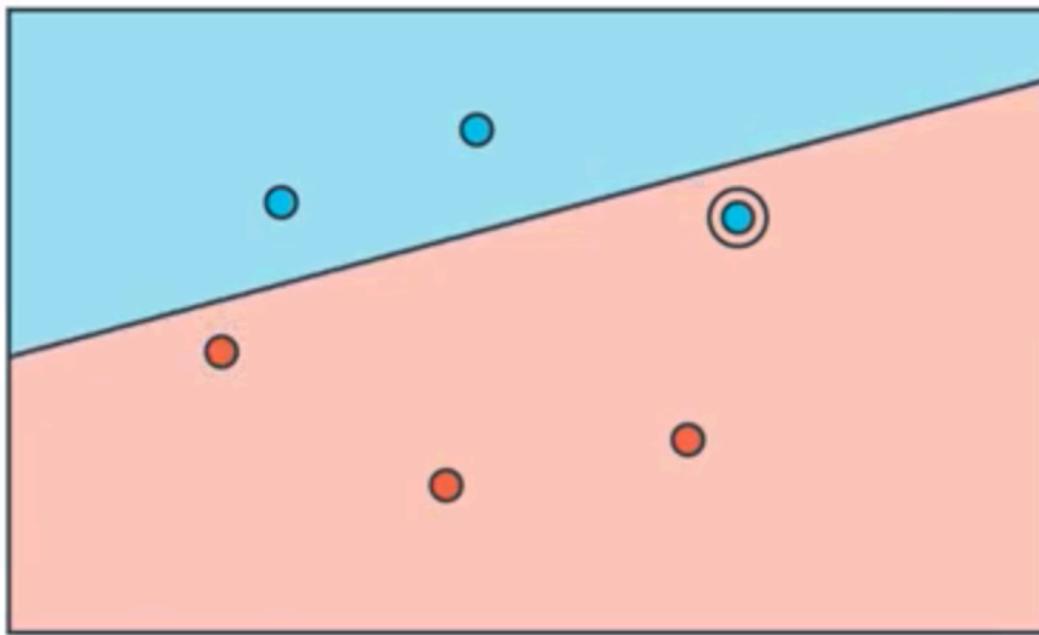


▷ O quão errado é o modelo?



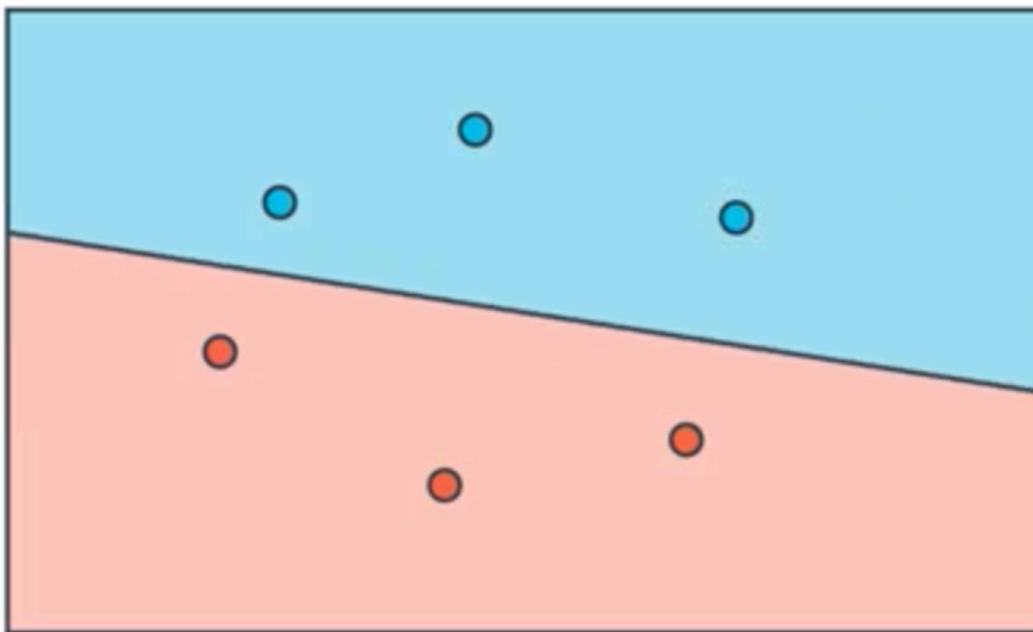
2 erros

▷ O quão errado é o modelo?



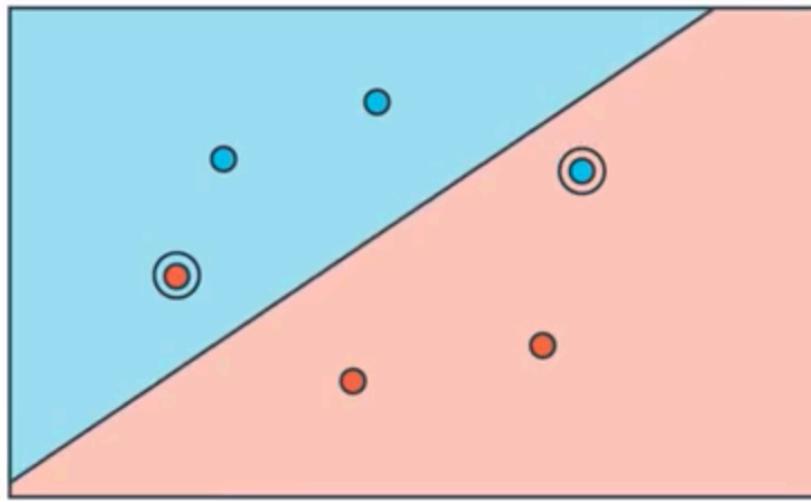
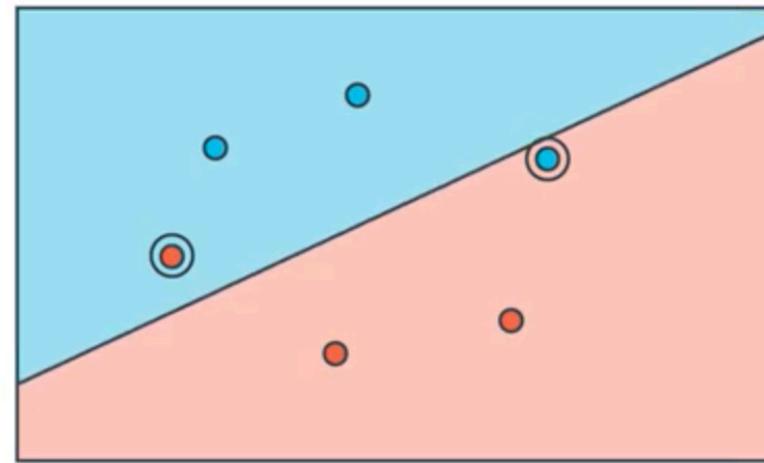
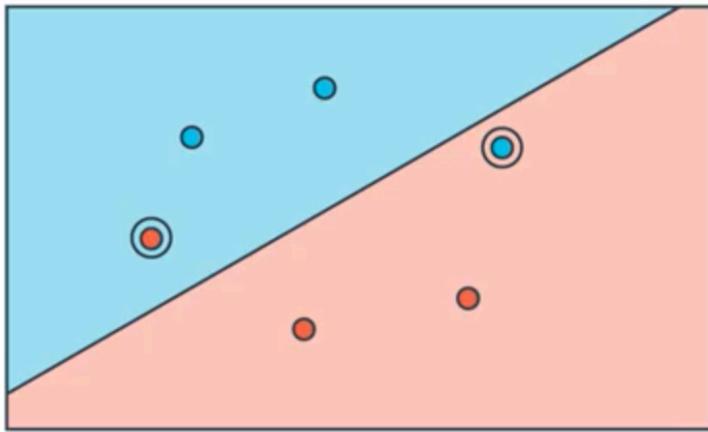
1 erro

▷ O quão errado é o modelo?



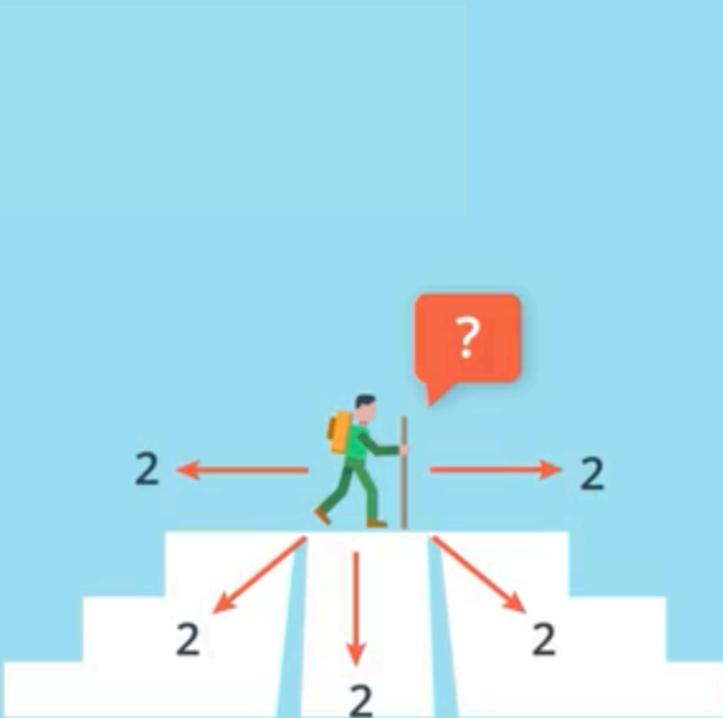
0 erros

▷ Problema!

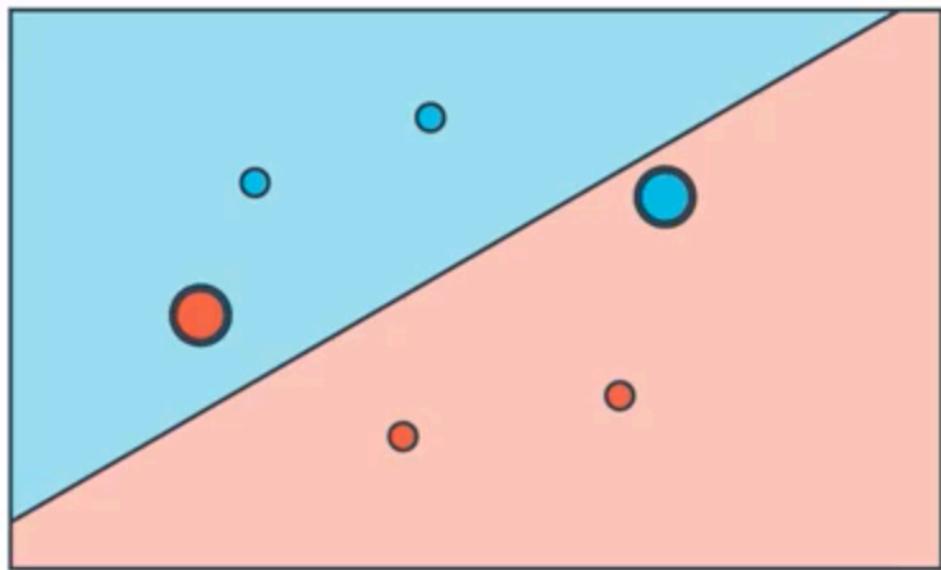


2 erros

Discreto vs Contínuo

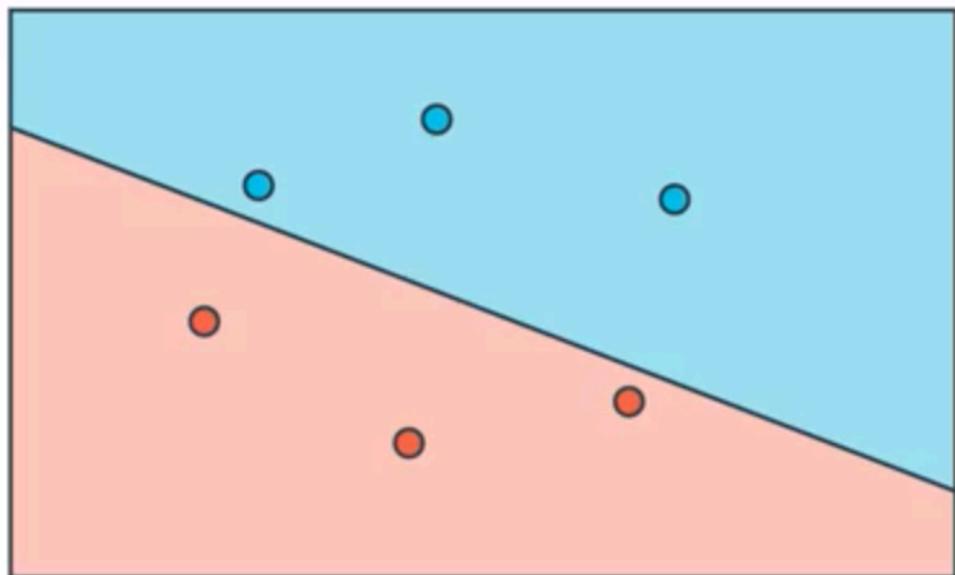


Função de Erro



ERROR = ● + ● + ● + ● + ● + ●

Função de Erro



$$\text{ERRO} = \textcolor{blue}{\bullet} + \textcolor{blue}{\bullet} + \textcolor{blue}{\bullet} + \textcolor{red}{\bullet} + \textcolor{red}{\bullet} + \textcolor{red}{\bullet}$$

$$\text{ERRO} = \textcolor{blue}{\bullet} + \textcolor{blue}{\bullet} + \textcolor{blue}{\bullet} + \textcolor{red}{\bullet} + \textcolor{red}{\bullet} + \textcolor{red}{\bullet}$$

MINIMIZAR ERRO



Gradient Descent



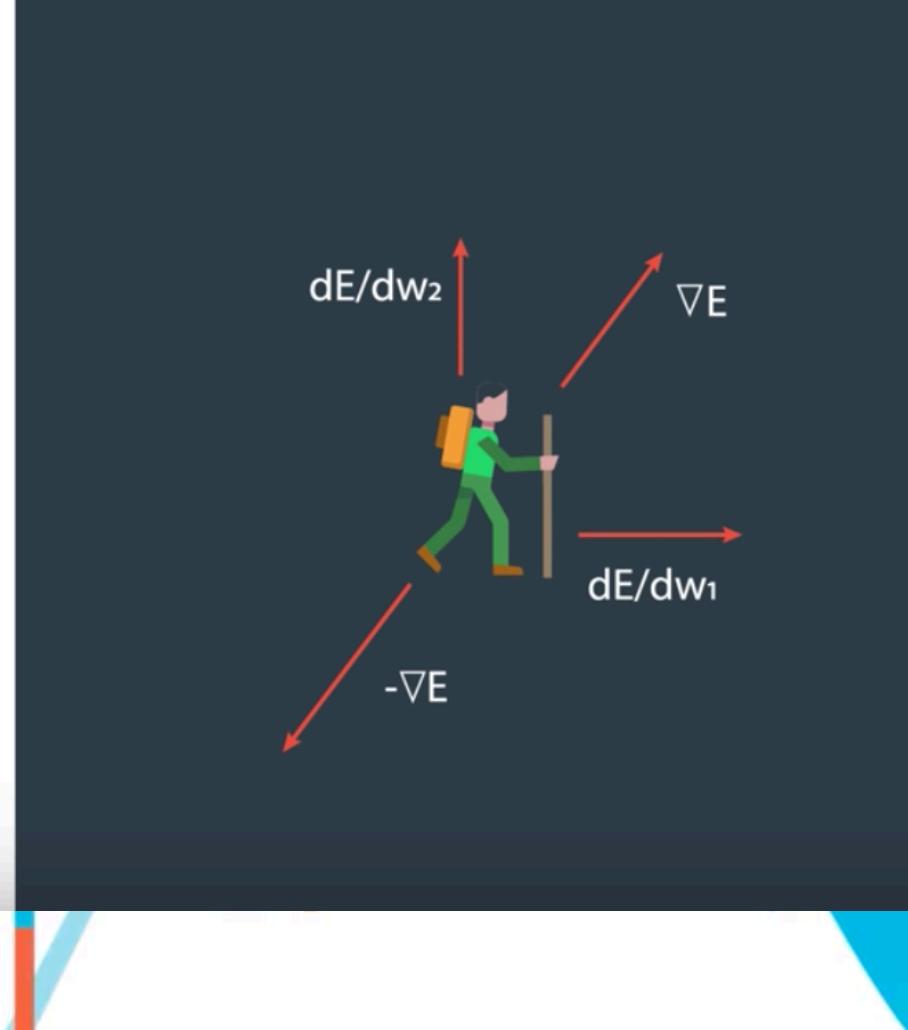
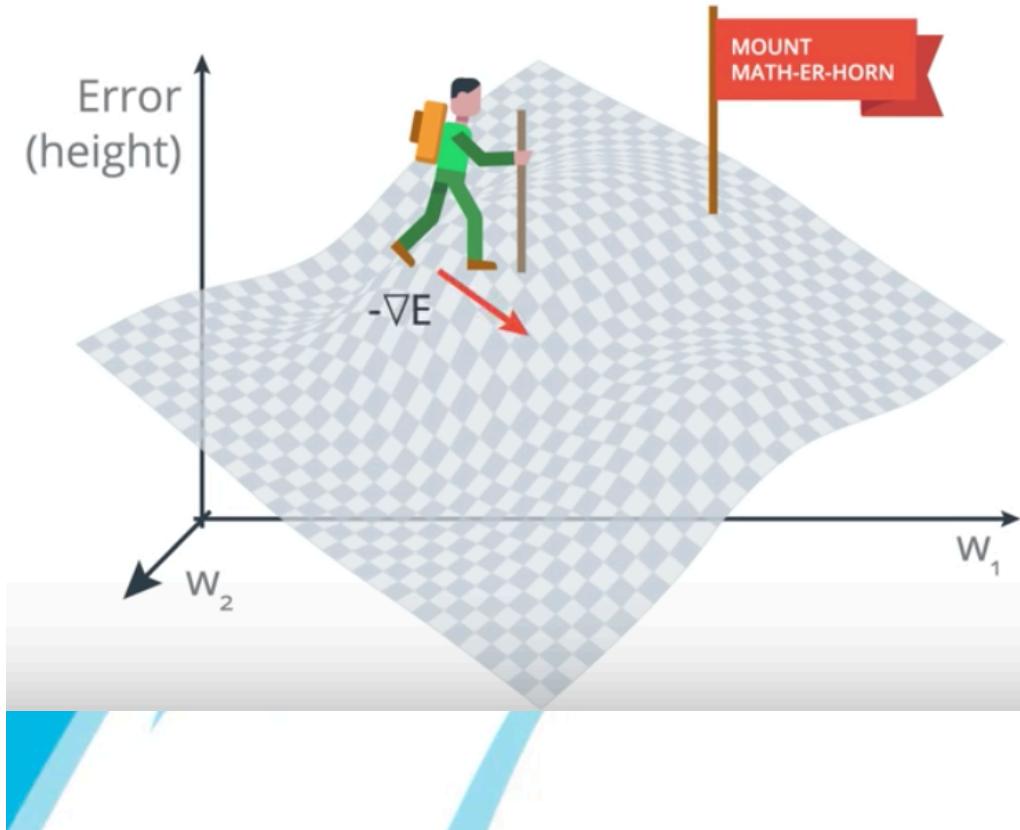
Gradient Descent



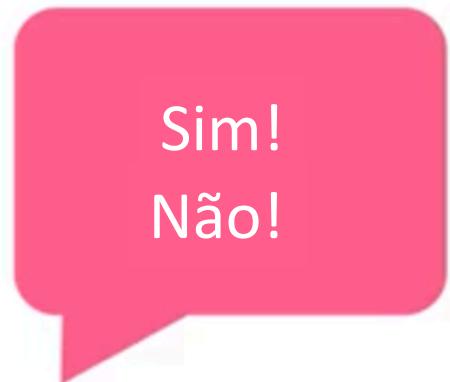
Gradient Descent



Gradient Descent

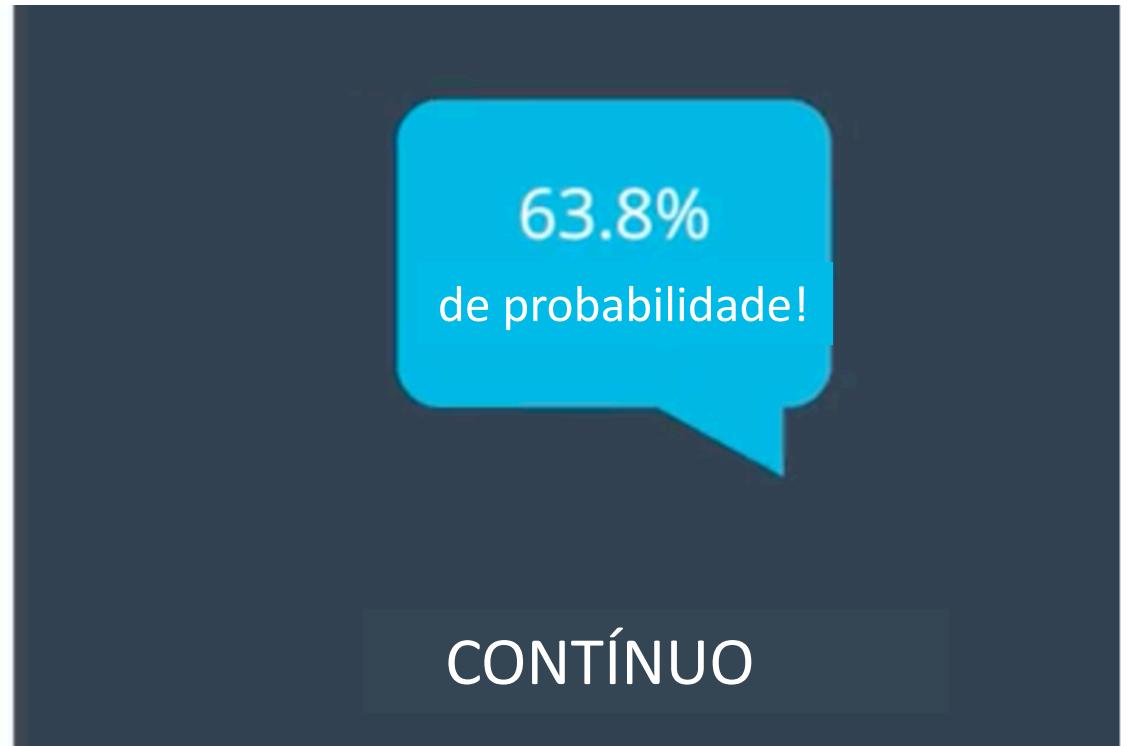


Predições – Discreta vs Contínua



Sim!
Não!

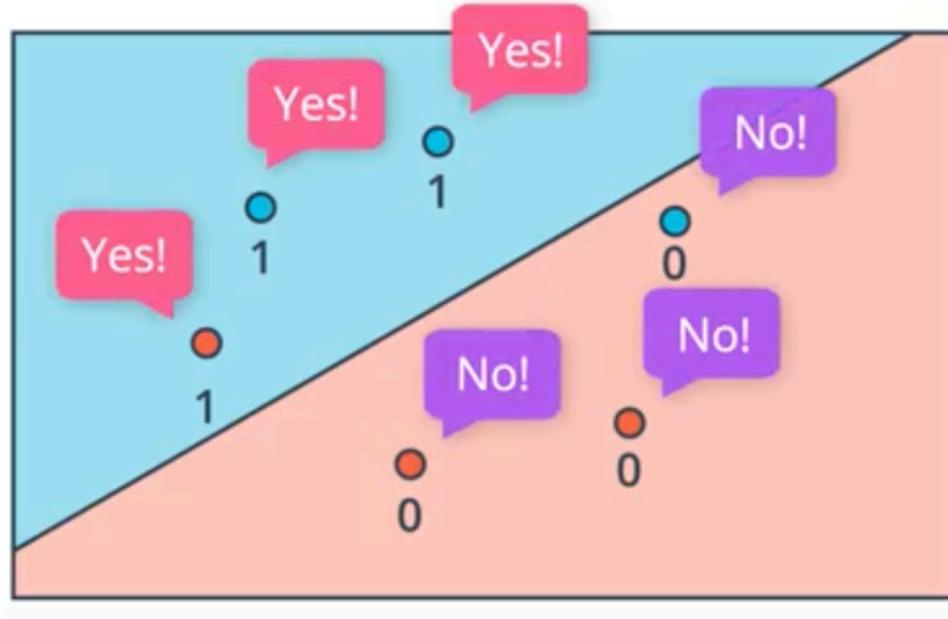
DISCRETO



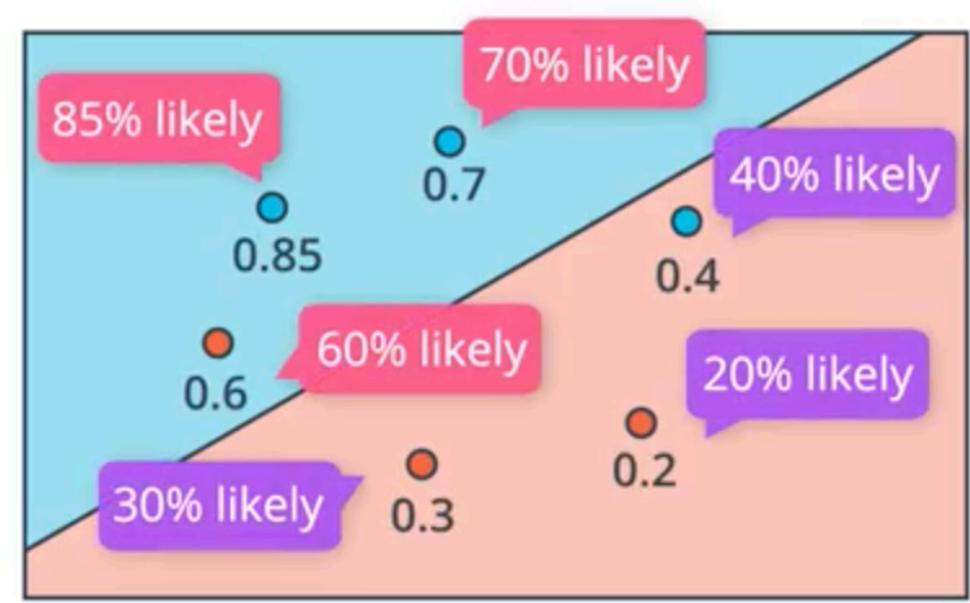
63.8%
de probabilidade!

CONTÍNUO

Predições

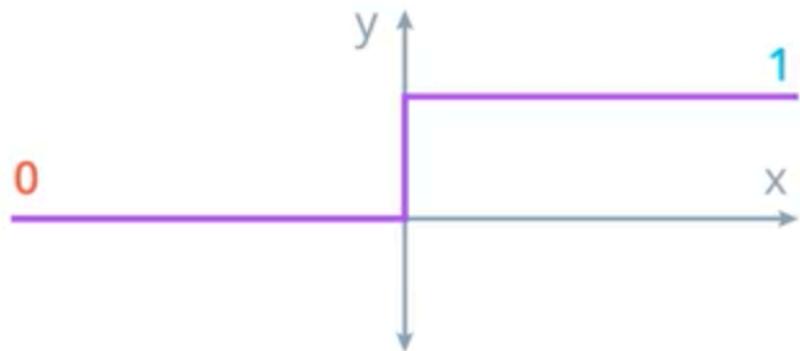


DISCRETO



CONTÍNUO

Funções de Ativação



DISCRETO:
Função Degrau

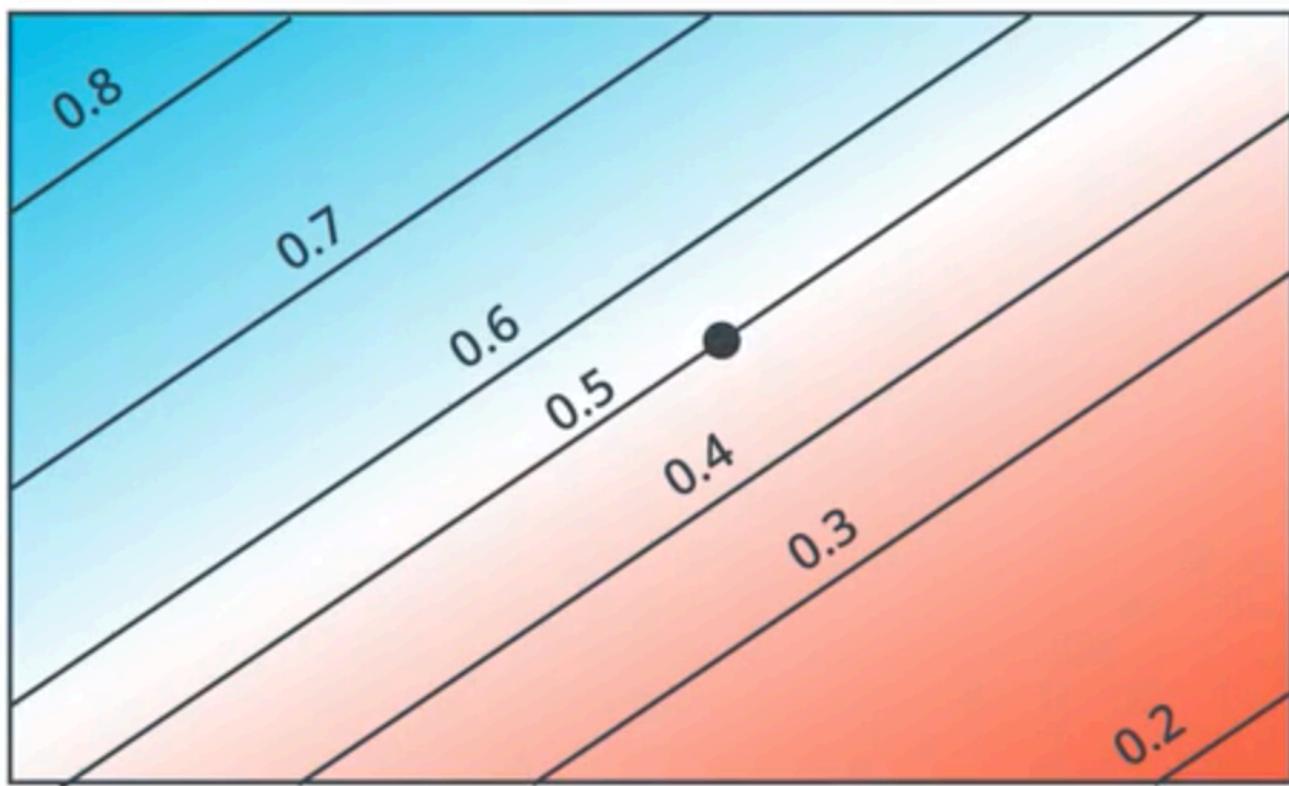
$$y = \begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$



CONTÍNUO
Função Sigmoide

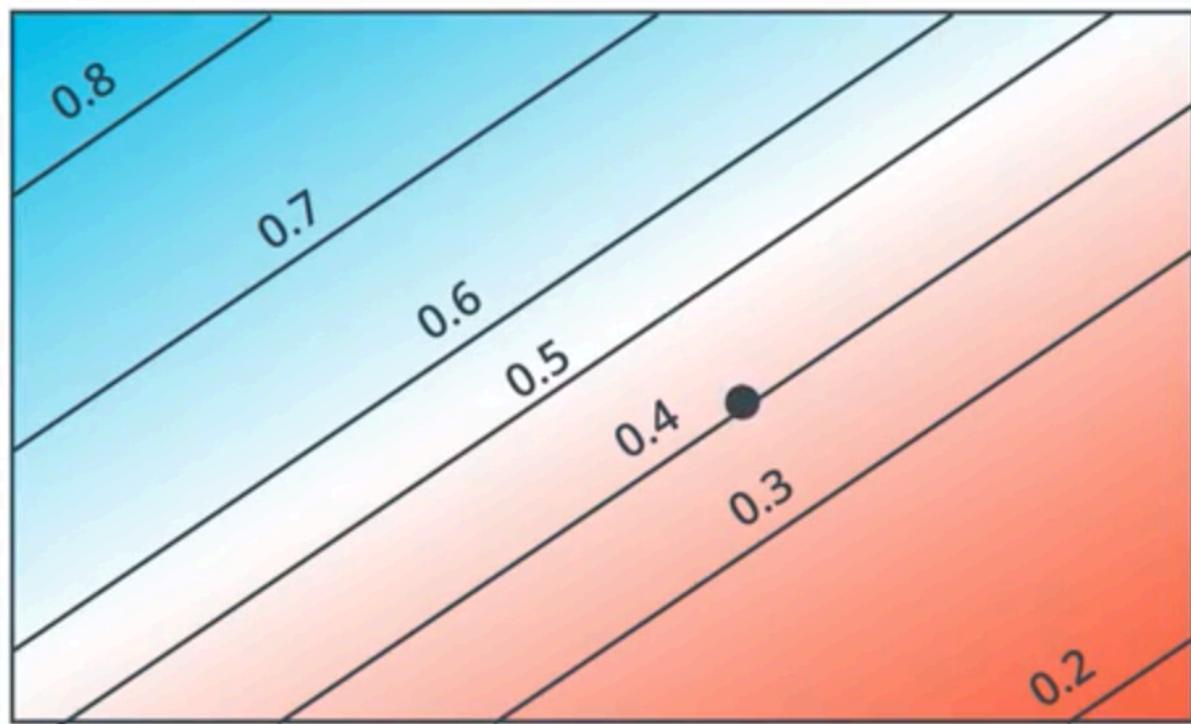
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Predições



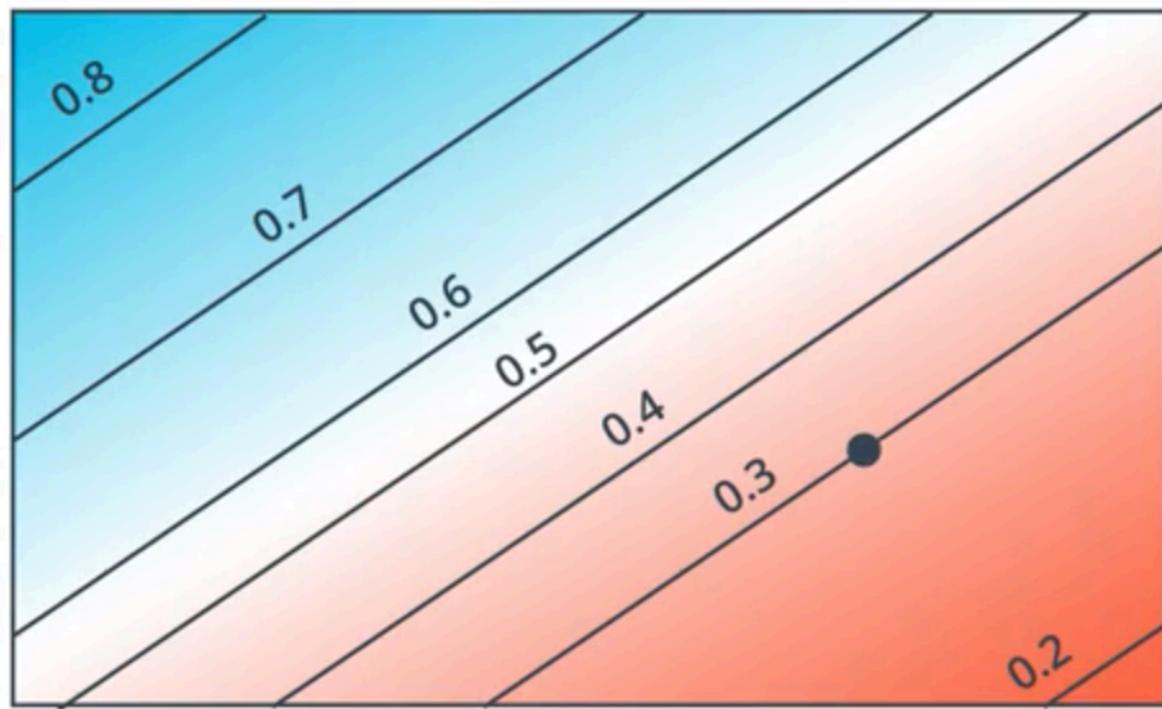
$P(\text{BLUE}) = 0.5$
 $P(\text{RED}) = 0.5$

Predições



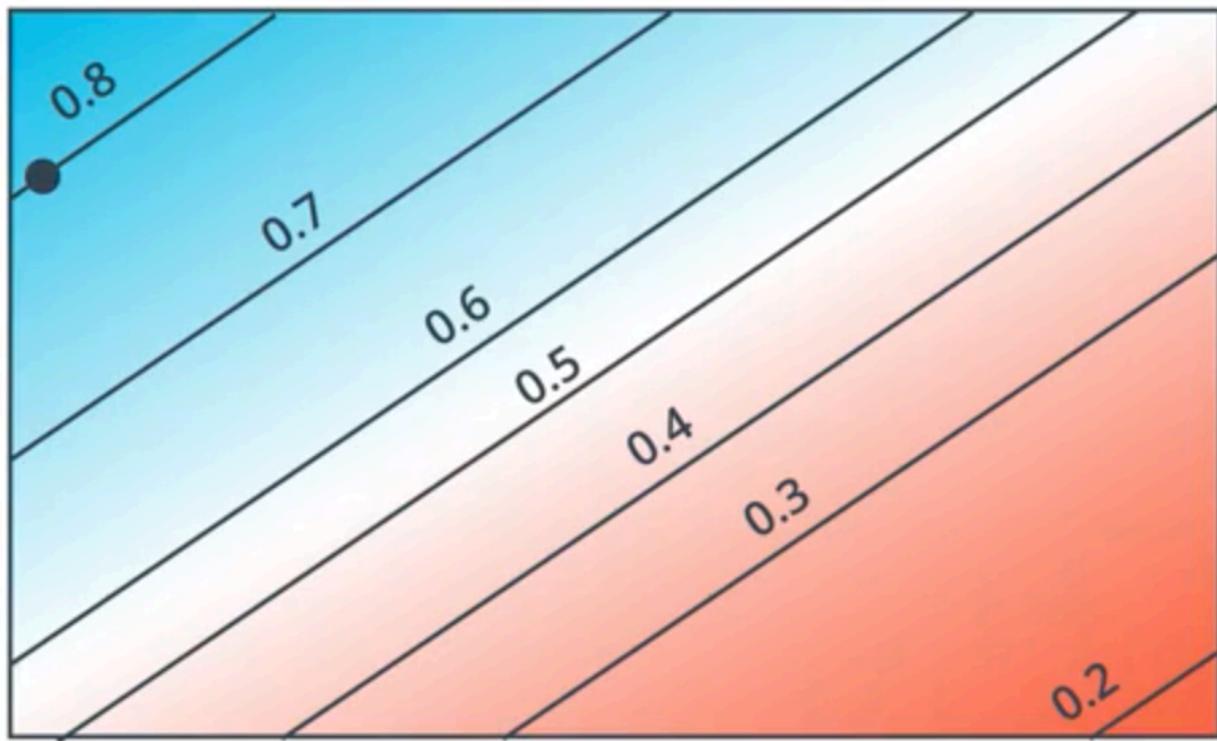
$P(\text{BLUE}) = 0.4$
 $P(\text{RED}) = 0.6$

Predicões



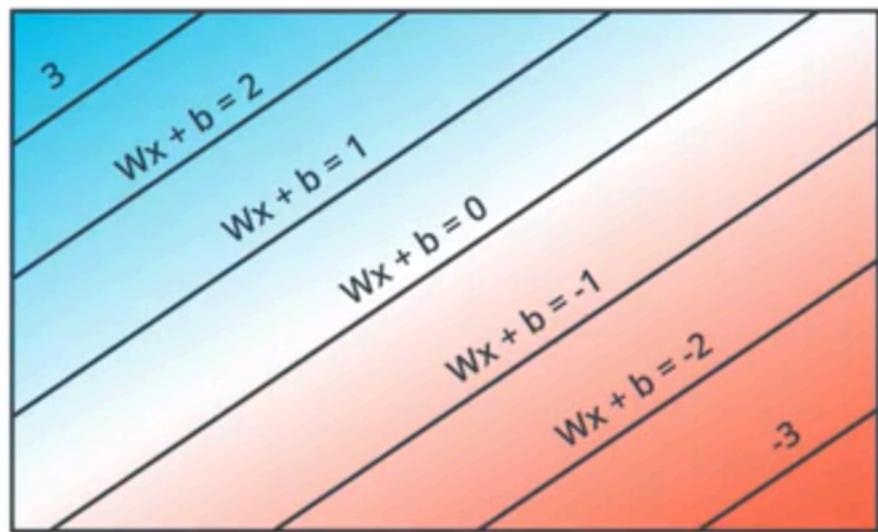
$P(\text{BLUE}) = 0.3$
 $P(\text{RED}) = 0.7$

Predições

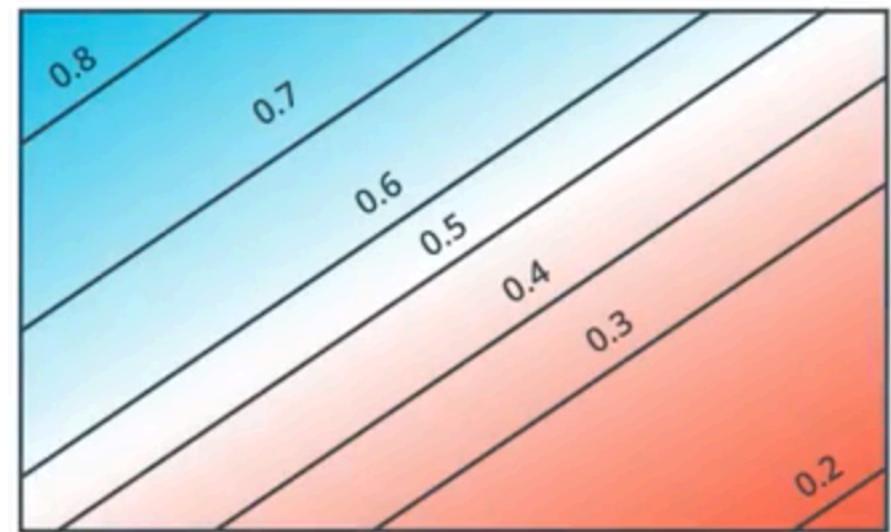


$P(\text{BLUE}) = 0.8$
 $P(\text{RED}) = 0.2$

Predições

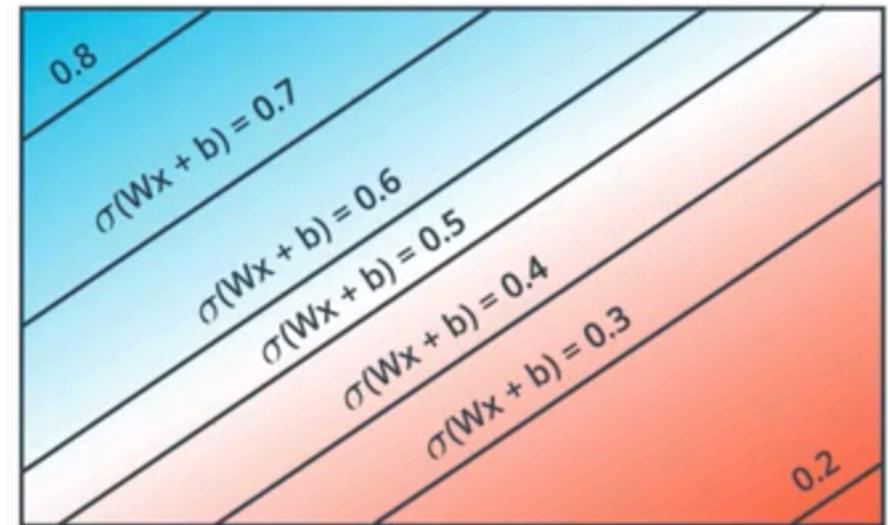
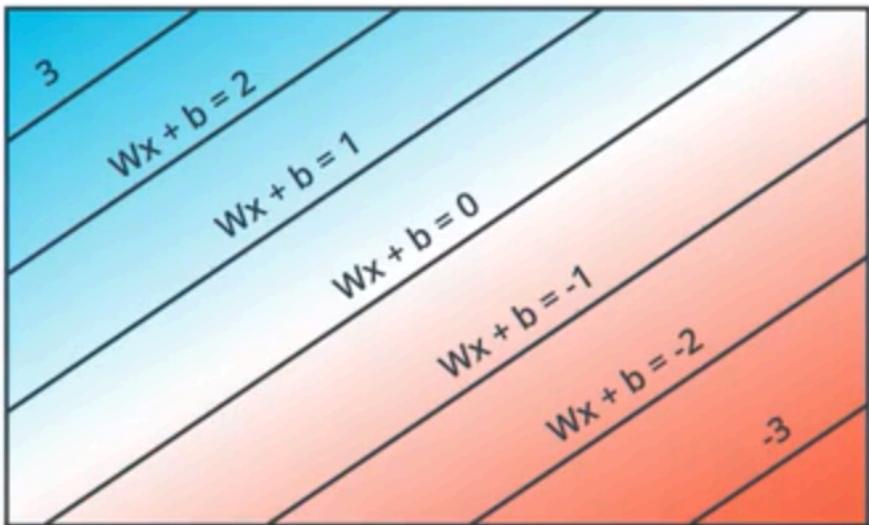


σ



$$Wx + b$$

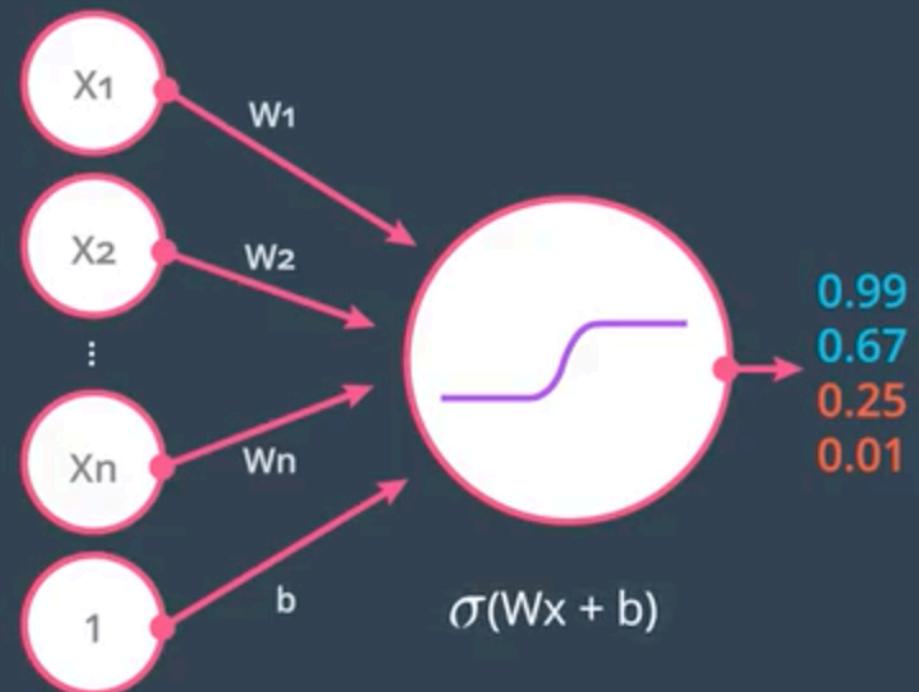
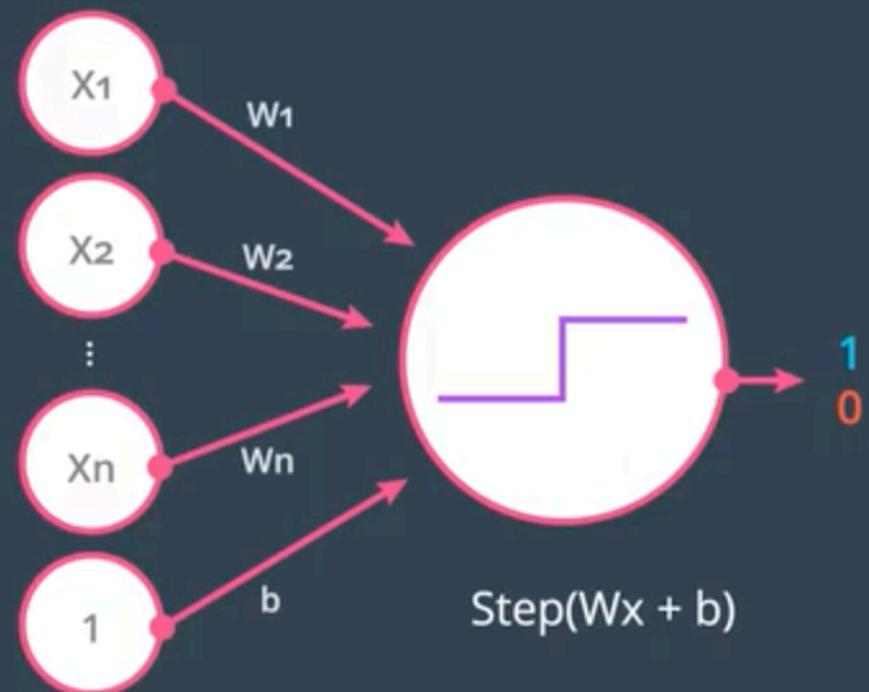
Predições



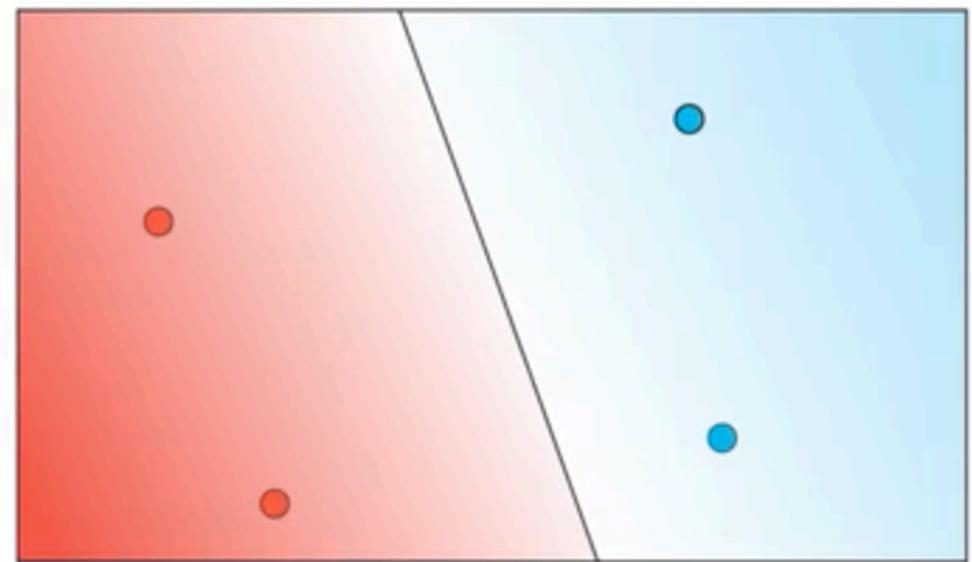
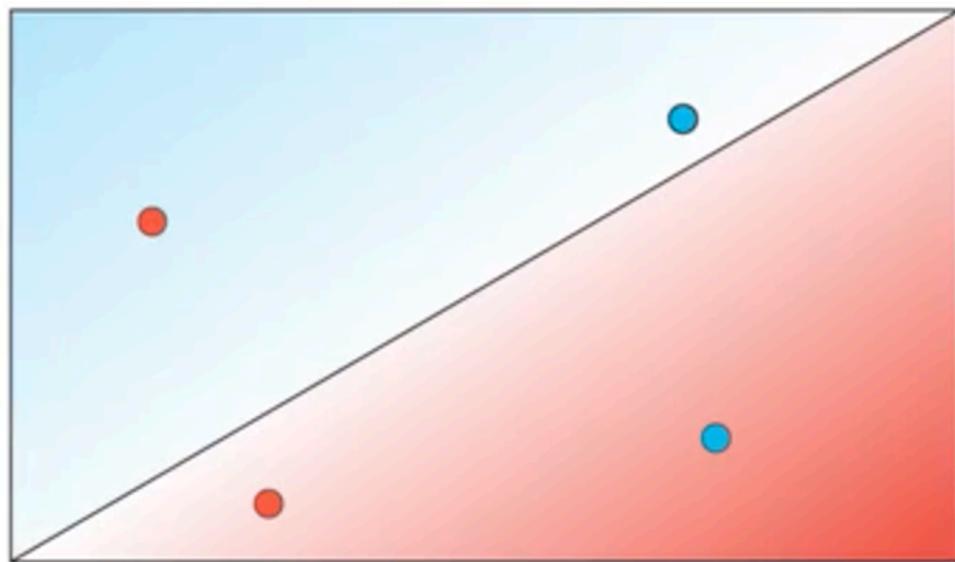
$$Wx + b$$

$$\hat{y} = \sigma(Wx + b)$$

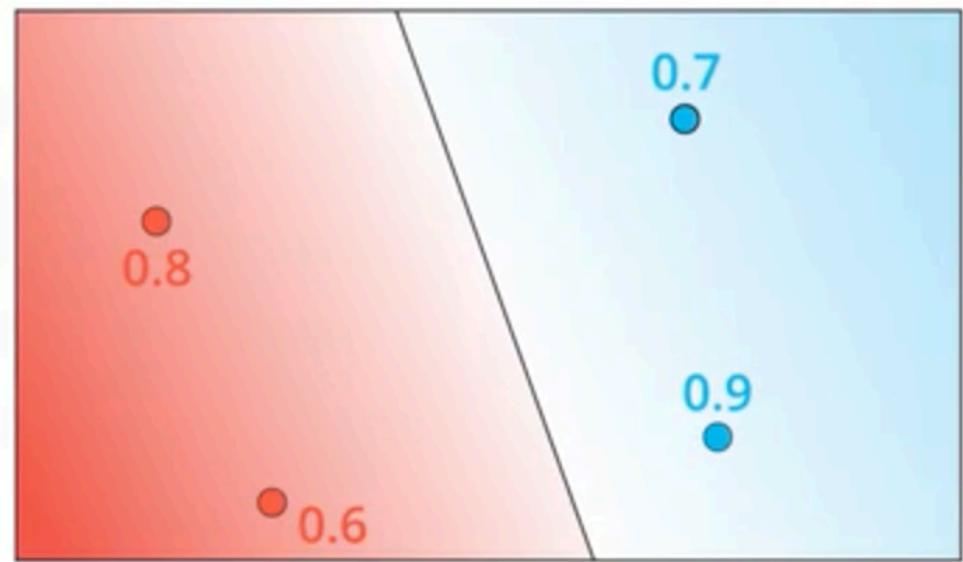
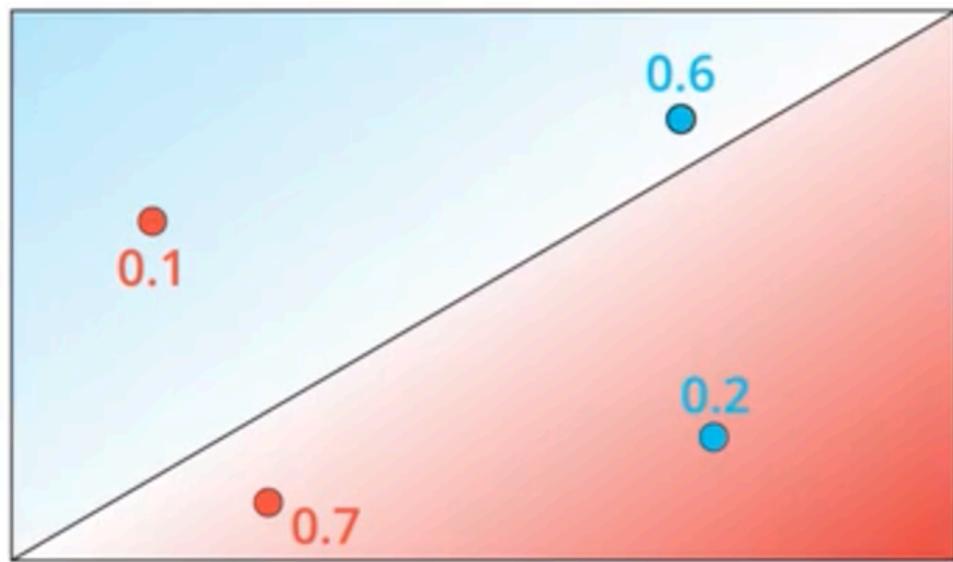
Perceptron



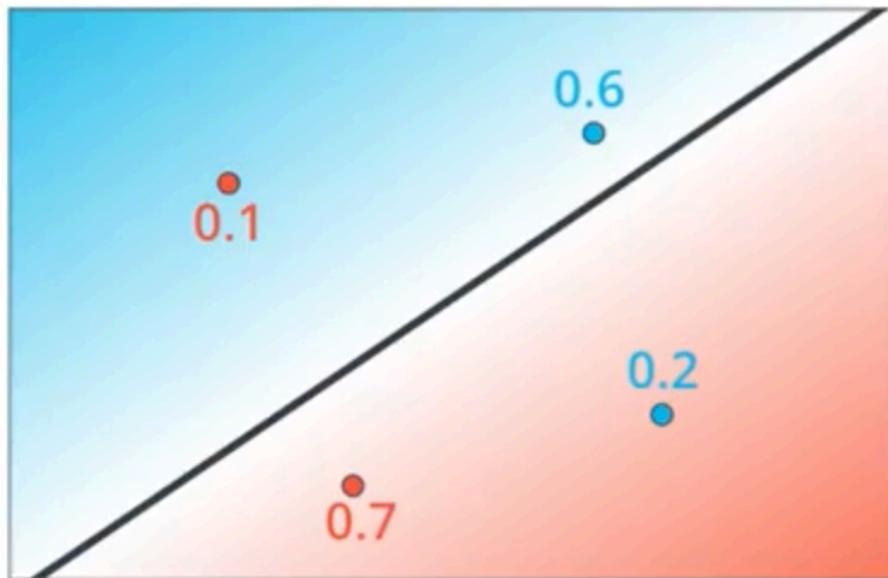
Qual desses modelos é melhor?



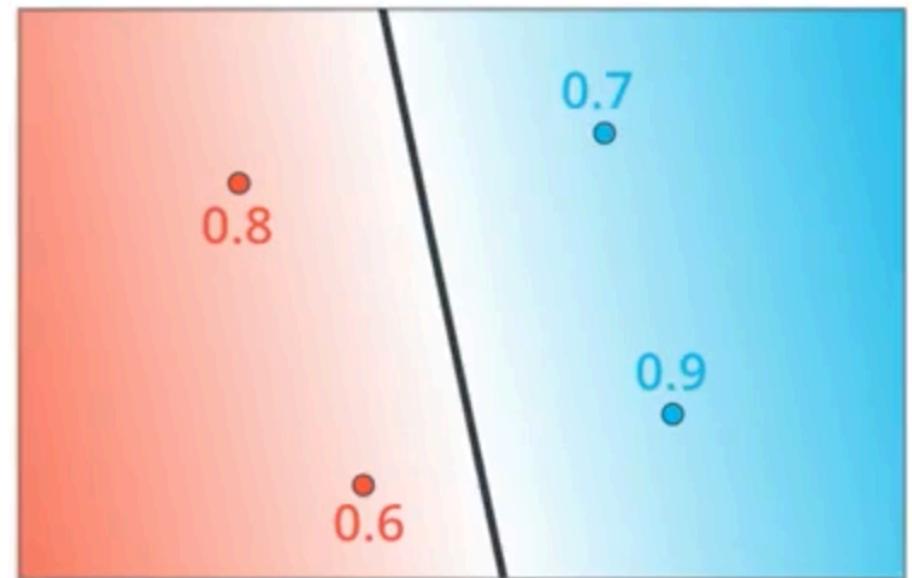
Qual desses modelos é melhor?



Qual desses modelos é melhor?

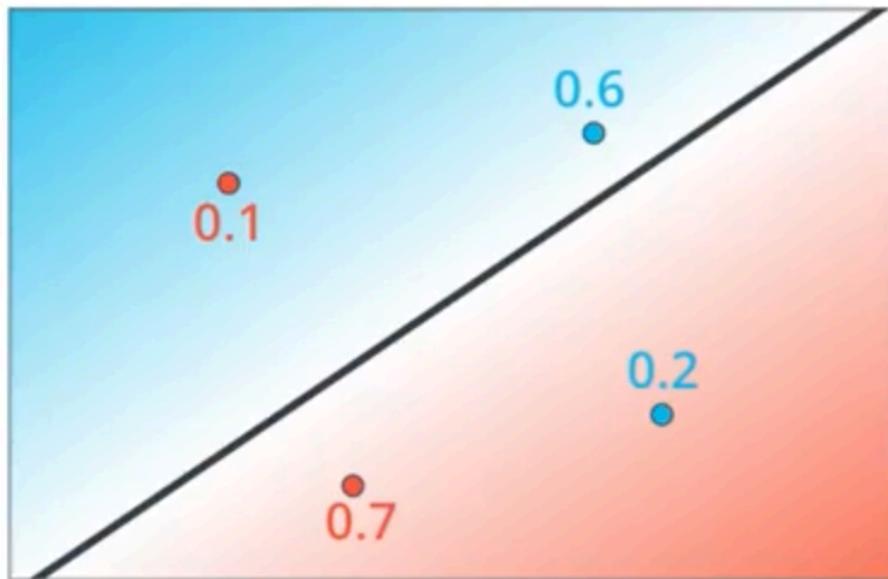


$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

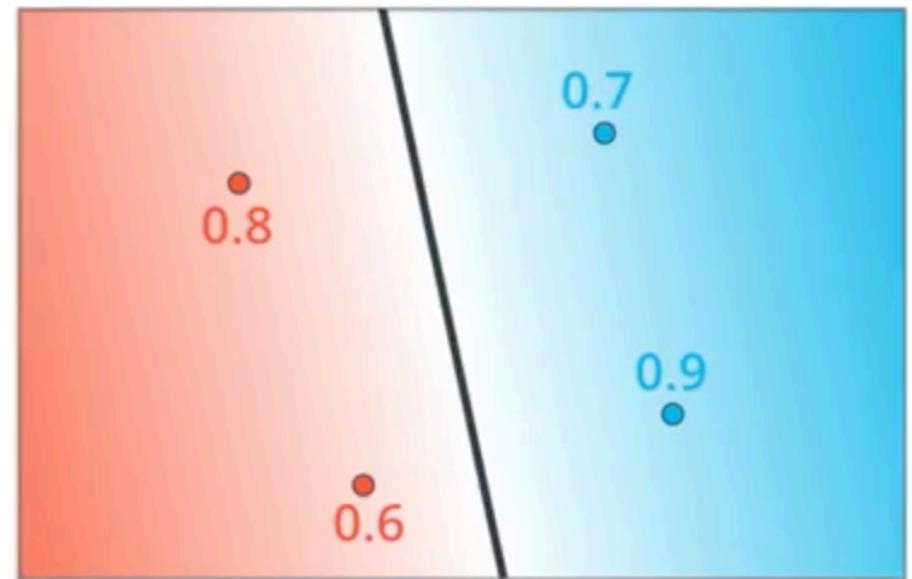


$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

Qual desses modelos é melhor?



$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$



$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$



Máxima Verossimilhança (Maximum Likelihood)

Problema

$$0.6*0.2*0.1*0.7 = 0.0084$$



Produto: Ruim!
Soma: Bom!

$$0.7*0.9*0.8*0.6 = 0.3024$$

Qual Função usar?

- sin
- cos
- log
- exp

Problema

$$0.6*0.2*0.1*0.7 = 0.0084$$



Produto: Ruim!
Soma: Bom!

$$0.7*0.9*0.8*0.6 = 0.3024$$

Qual Função usar?

- sin
- cos
- log
- exp

$$\log(ab) = \log(a) + \log(b)$$

Produtos

$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

Produtos

$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

$$\ln(0.6) + \ln(0.2) + \ln(0.1) + \ln(0.7)$$

-0.51 -1.61 -2.3 -0.36

$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

$$\ln(0.7) + \ln(0.9) + \ln(0.8) + \ln(0.6)$$

-0.36 -0.1 -.22 -0.51

Produtos

$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

$$\ln(0.6) + \ln(0.2) + \ln(0.1) + \ln(0.7)$$

-0.51 -1.61 -2.3 -0.36

$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

$$\ln(0.7) + \ln(0.9) + \ln(0.8) + \ln(0.6)$$

-0.36 -0.1 -.22 -0.51

$$-\ln(0.6) - \ln(0.2) - \ln(0.1) - \ln(0.7) = 4.8$$

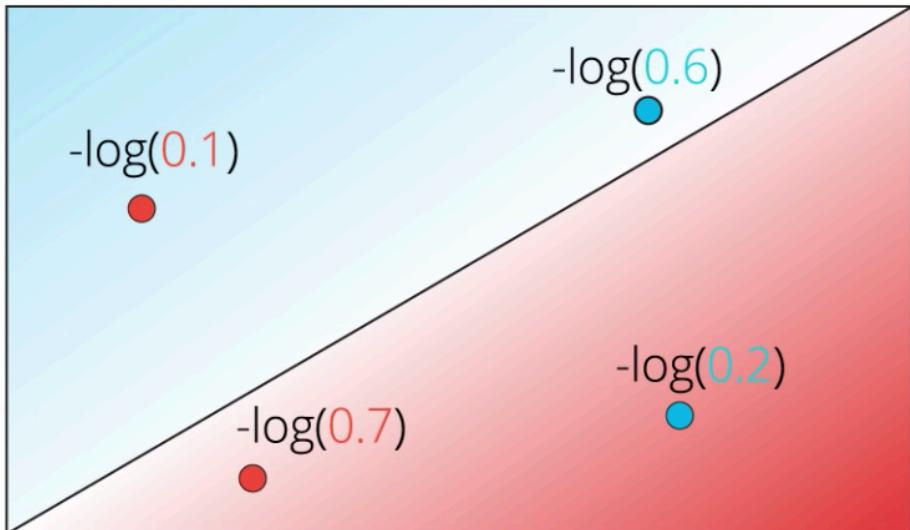
0.51 1.61 2.3 0.36

$$-\ln(0.7) - \ln(0.9) - \ln(0.8) - \ln(0.6) = 1.2$$

0.36 0.1 .22 0.51

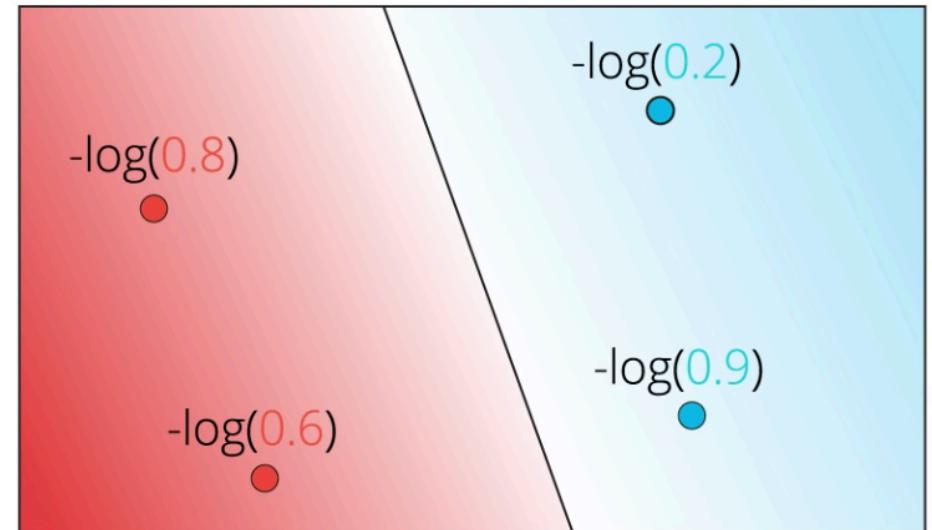
Cross Entropy

Cross Entropia



$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

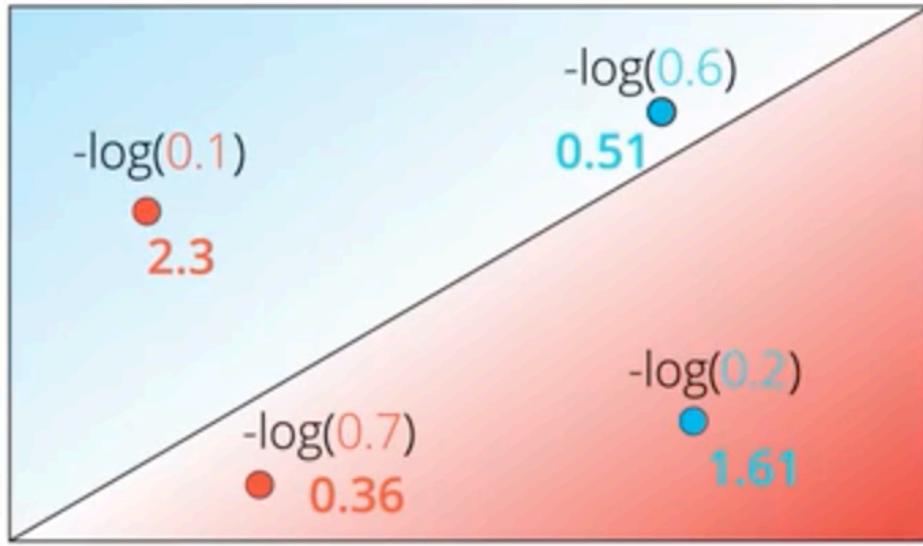
$$-\log(0.6) - \log(0.2) - \log(0.1) - \log(0.7) = 4.8$$



$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

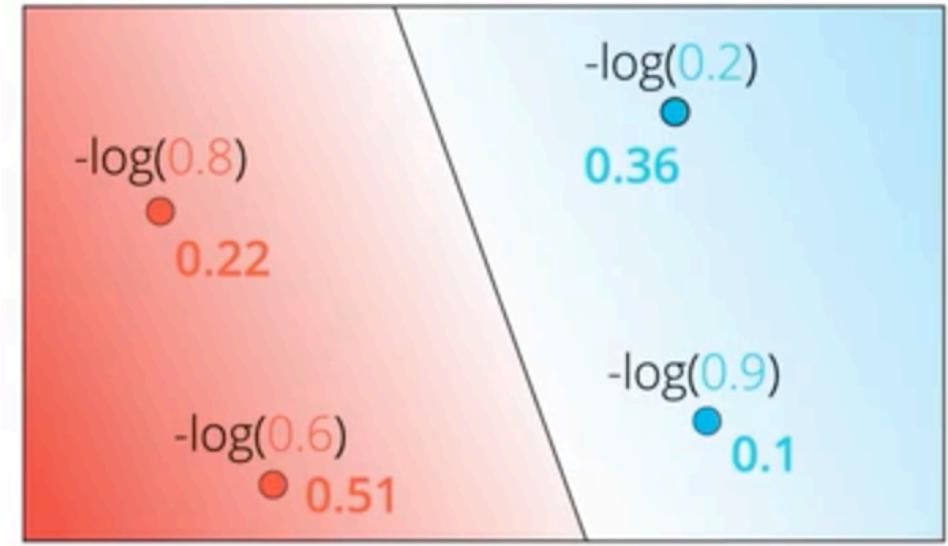
$$-\log(0.7) - \log(0.9) - \log(0.8) - \log(0.6) = 1.2$$

Cross Entropia



$$0.6 * 0.2 * 0.1 * 0.7 = 0.0084$$

$$-\log(0.6) - \log(0.2) - \log(0.1) - \log(0.7) = 4.8$$

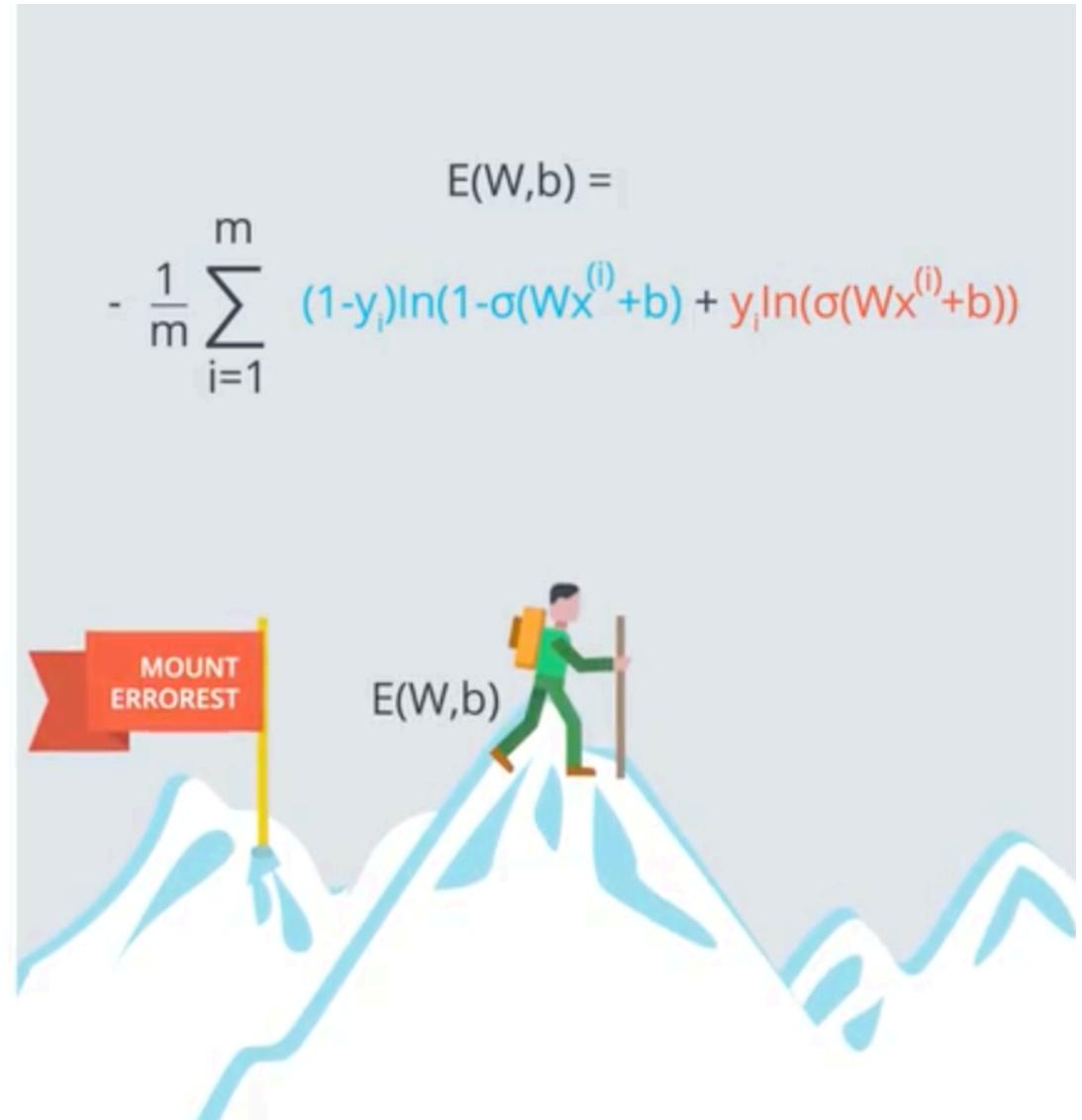
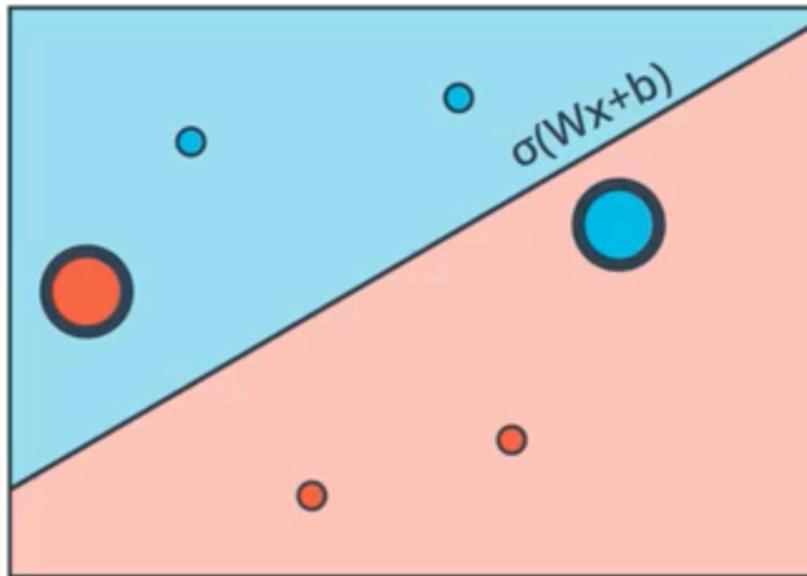


$$0.7 * 0.9 * 0.8 * 0.6 = 0.3024$$

$$-\log(0.7) - \log(0.9) - \log(0.8) - \log(0.6) = 1.2$$

Objetivo: Minimizar a Cross Entropia

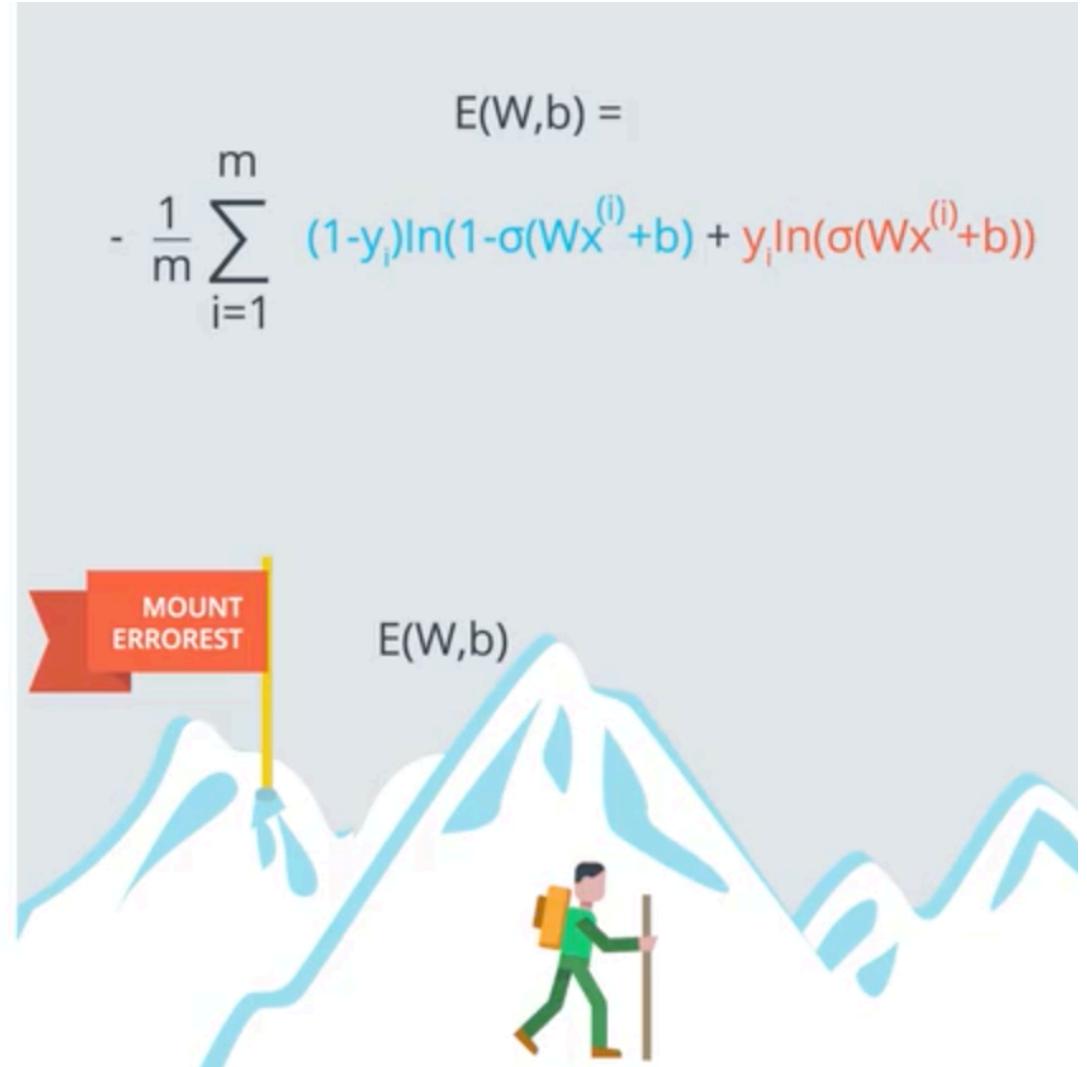
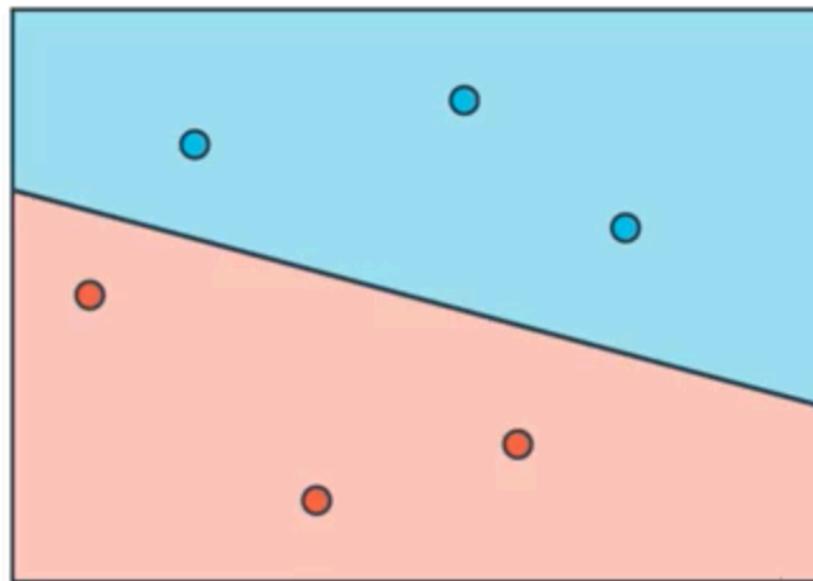
Objetivo: Minimizar Função de Erro



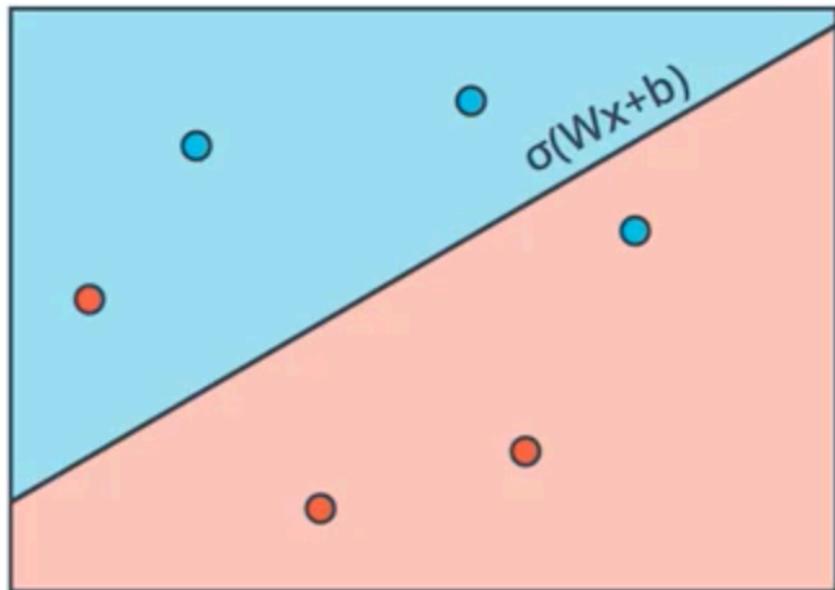
$$E(W, b) =$$

$$-\frac{1}{m} \sum_{i=1}^m (1-y_i)\ln(1-\sigma(Wx^{(i)}+b)) + y_i\ln(\sigma(Wx^{(i)}+b))$$

Objetivo: Minimizar Função de Erro



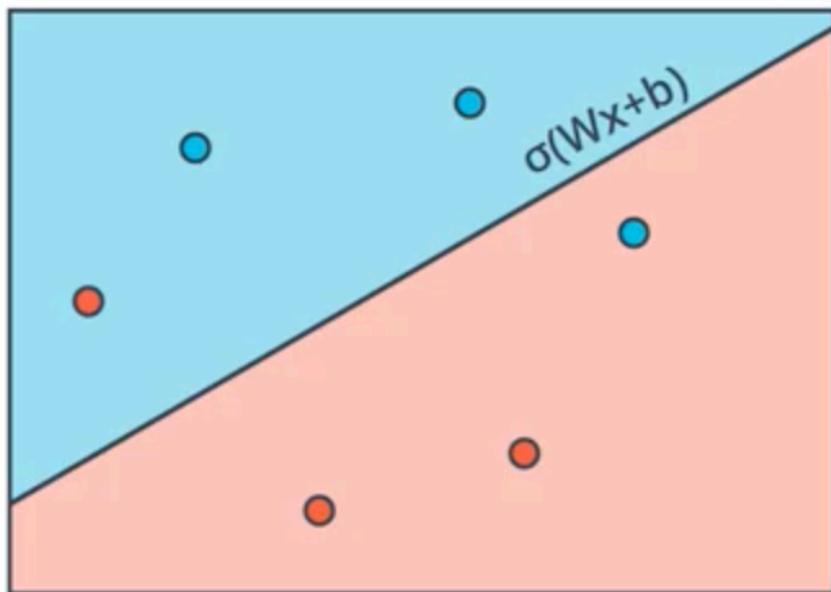
Algoritmo Gradiente Descendente



1. Comece com pesos aleatórios:

$$w_1, \dots, w_n, b$$

Algoritmo Gradiente Descendente



1. Comece com pesos aleatórios:

$$w_1, \dots, w_n, b$$

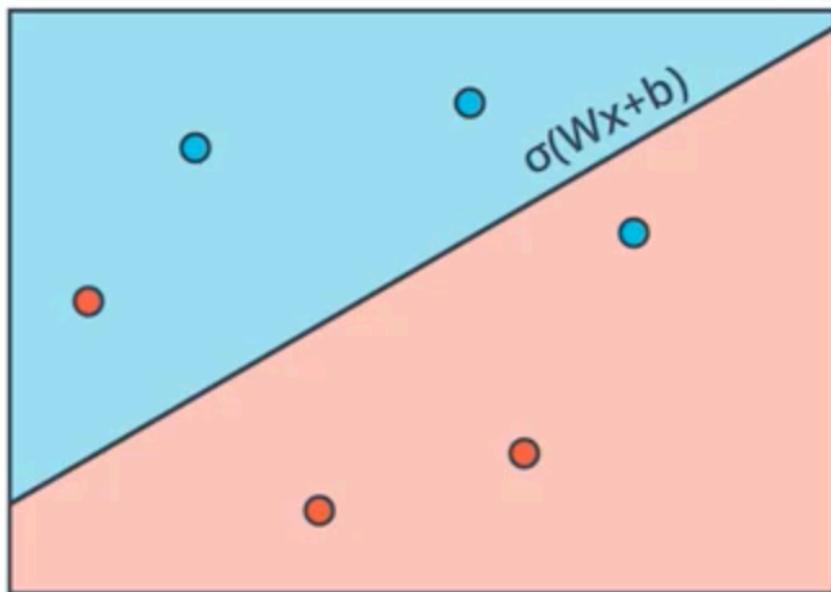
2. Para cada ponto (x_1, \dots, x_n) :

2.1 Para $i=1 \dots n$

2.1.1 Atualiza $w'_i \leftarrow w_i - \alpha \partial E / \partial w_i$

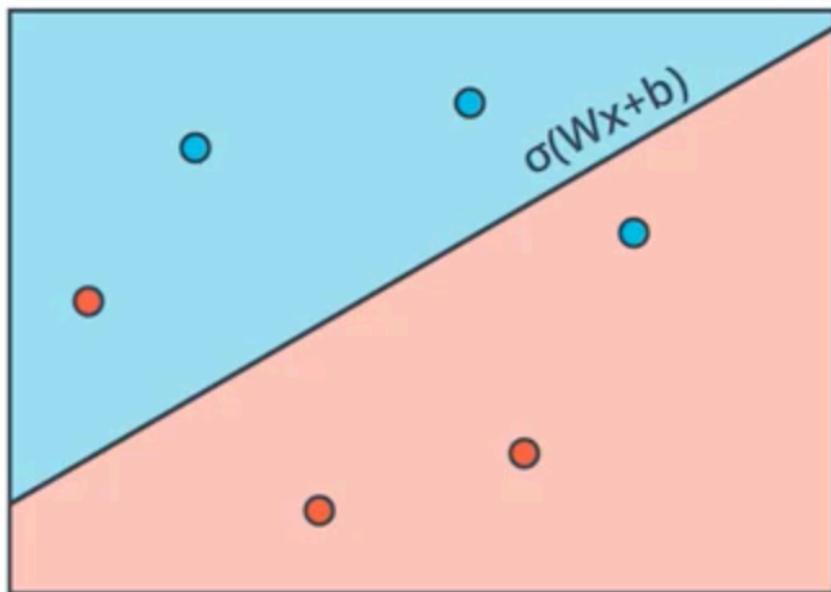
2.1.2 Atualiza $b' \leftarrow b - \alpha \partial E / \partial b$

Algoritmo Gradiente Descendente



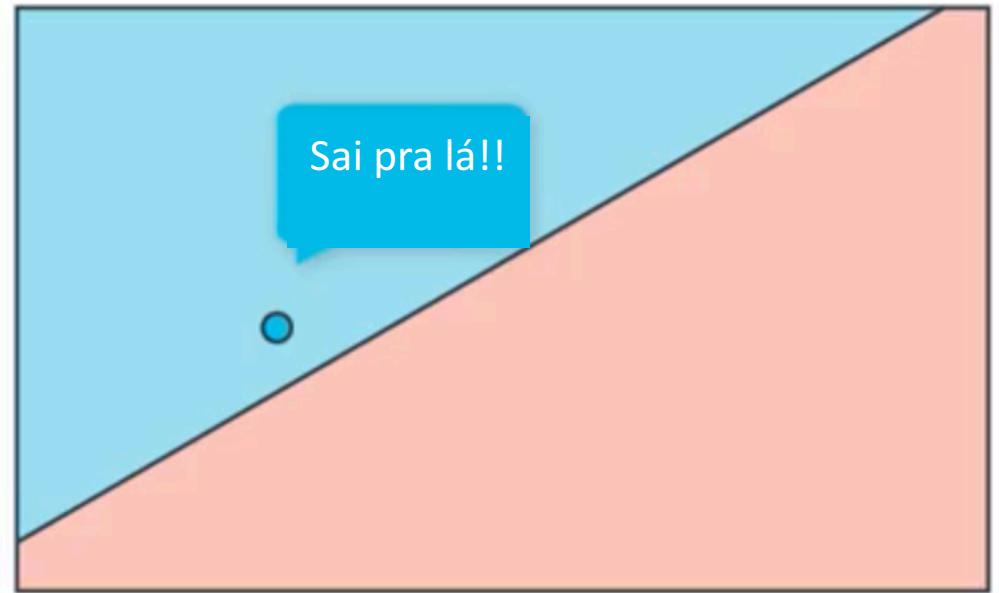
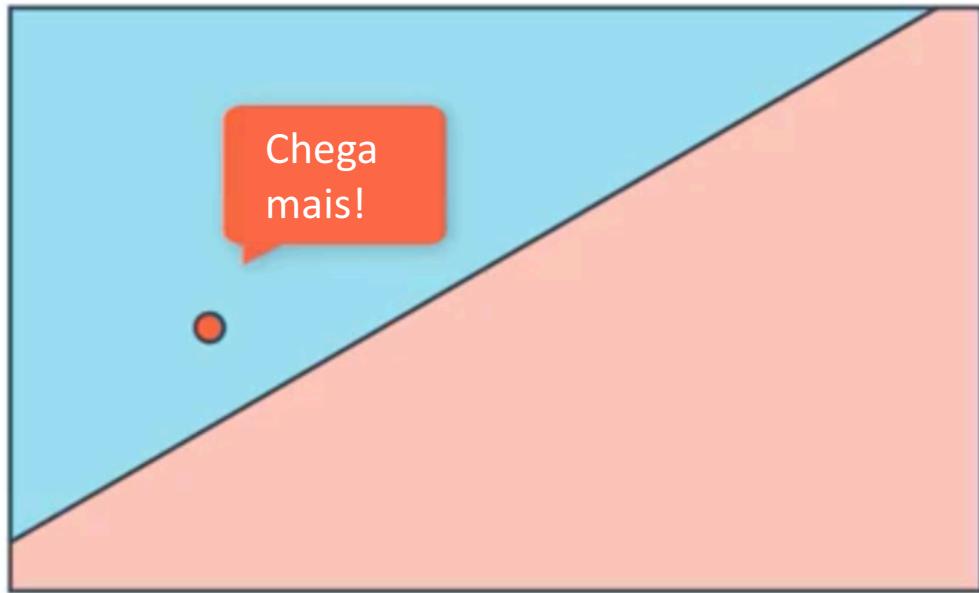
1. Comece com pesos aleatórios:
 w_1, \dots, w_n, b
2. Para cada ponto (x_1, \dots, x_n) :
 - 2.1 Para $i=1 \dots n$
 - 2.1.1 Atualiza $w'_i \leftarrow w_i - \alpha (\hat{y} - y)x_i$
 - 2.1.2 Atualiza $b' \leftarrow b - \alpha(\hat{y} - y)$

Algoritmo Gradiente Descendente



1. Comece com pesos aleatórios:
 w_1, \dots, w_n, b
2. Para cada ponto (x_1, \dots, x_n) :
 - 2.1 Para $i=1 \dots n$
 - 2.1.1 Atualiza $w'_i \leftarrow w_i - \alpha (\hat{y} - y)x_i$
 - 2.1.2 Atualiza $b' \leftarrow b - \alpha(\hat{y} - y)$
3. Repetir até erro ser pequeno

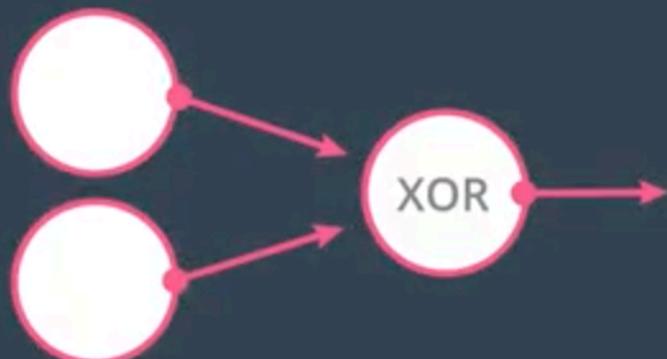
Perceptron vs Gradiente Descendente



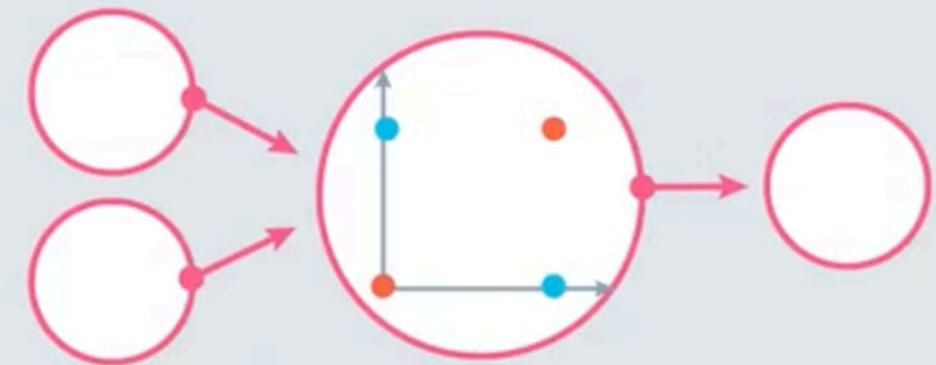
Exercício

- ▷ Implementar algoritmo de Gradiente Descendente.
 - ex_gradient_descent.ipynb

XOR Perceptron



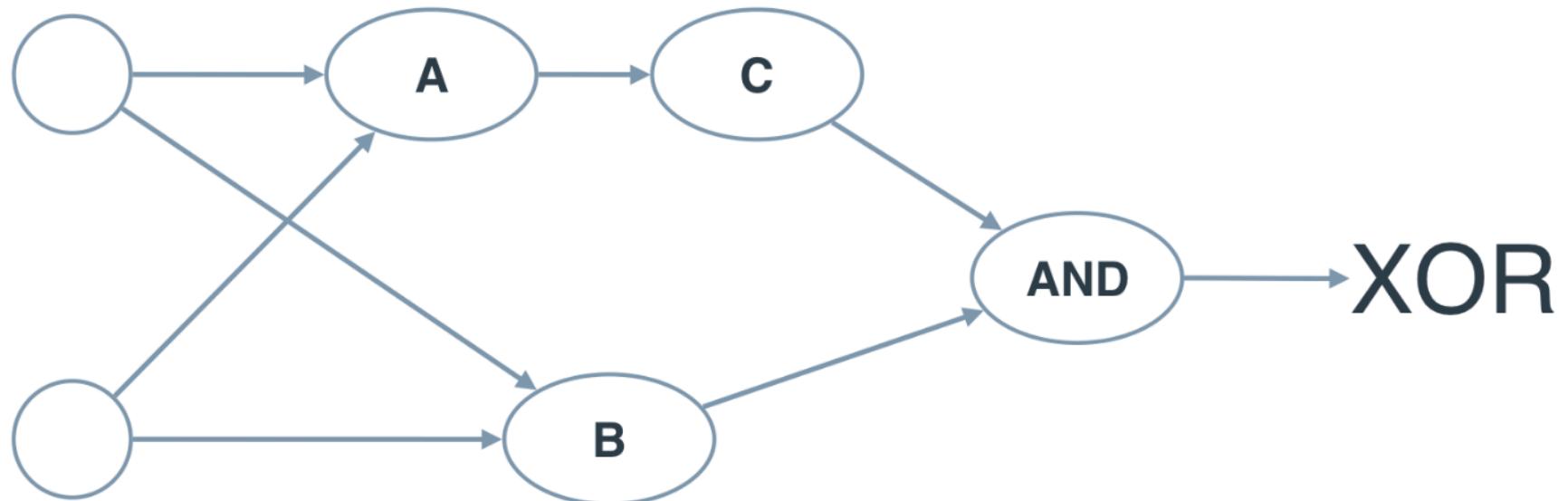
IN	IN	OUT
✓	✓	✗
✓	✗	✓
✗	✓	✓
✗	✗	✗



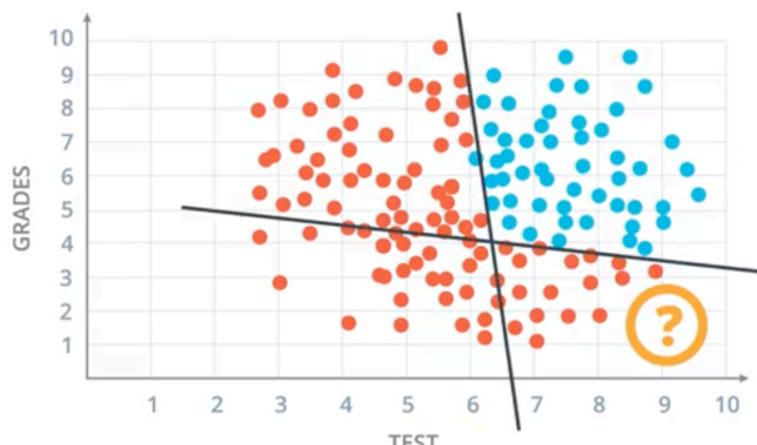
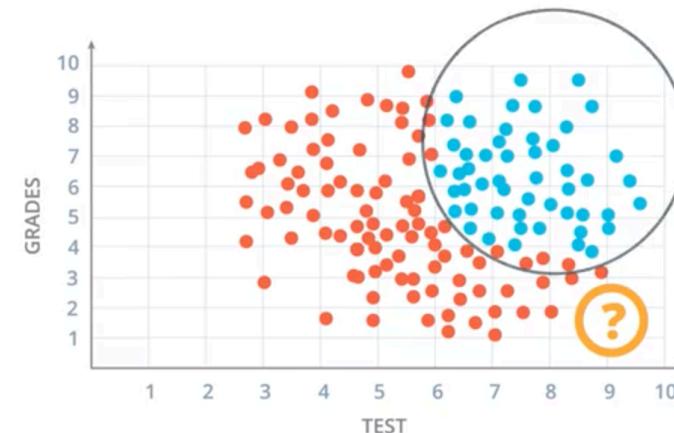
IN	IN	OUT
1	1	0
1	0	1
0	1	1
0	0	0

Quiz – Construir um XOR MultiLayer Perceptron

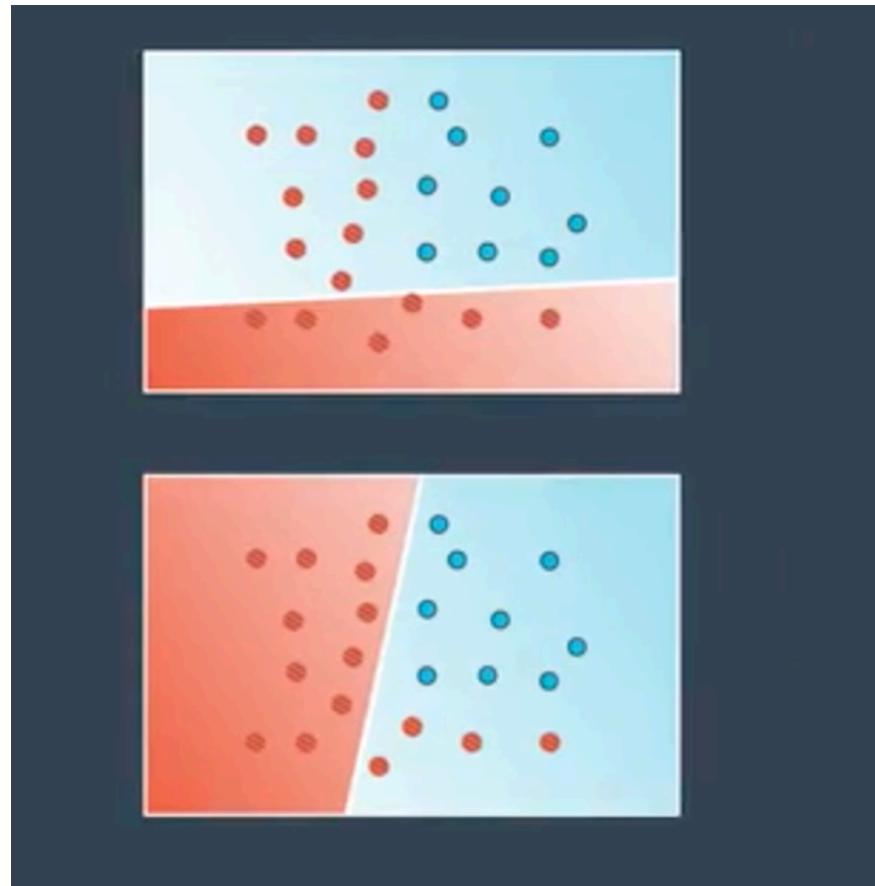
- ▷ A rede neural abaixo contém 3 perceptrons, A, B, e C. O último (AND) foi fornecido. A entrada da rede neural é através do primeiro nó. A saída é pelo último nó.
- ▷ O perceptron multi-camadas abaixo calcula a operação XOR. Cada perceptron é um operação lógica de AND, OR e NOT. Diga os operações de cada perceptron para calcular o XOR.



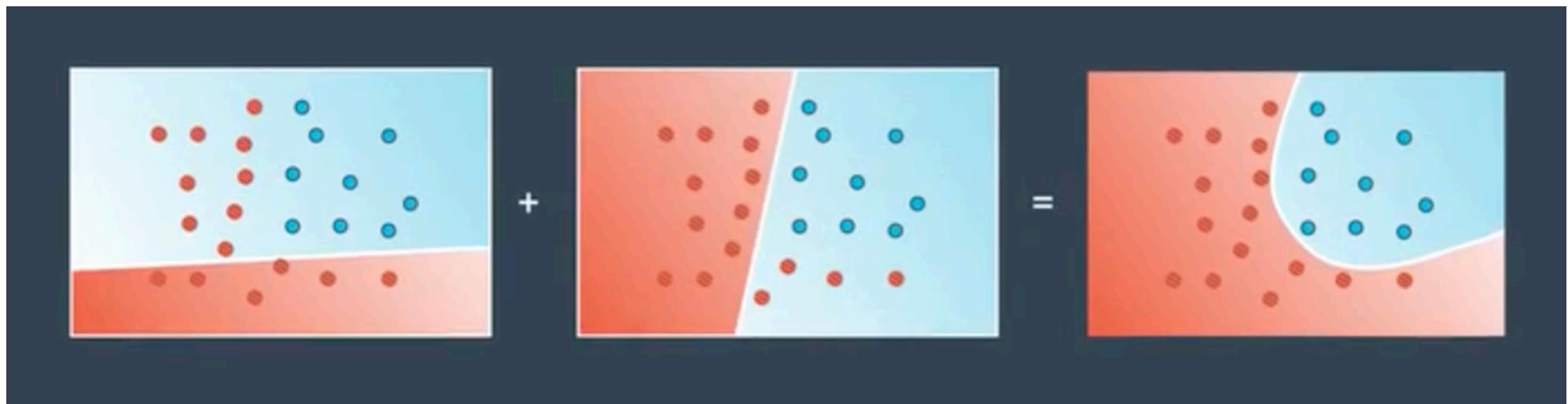
Regiões Não Lineares



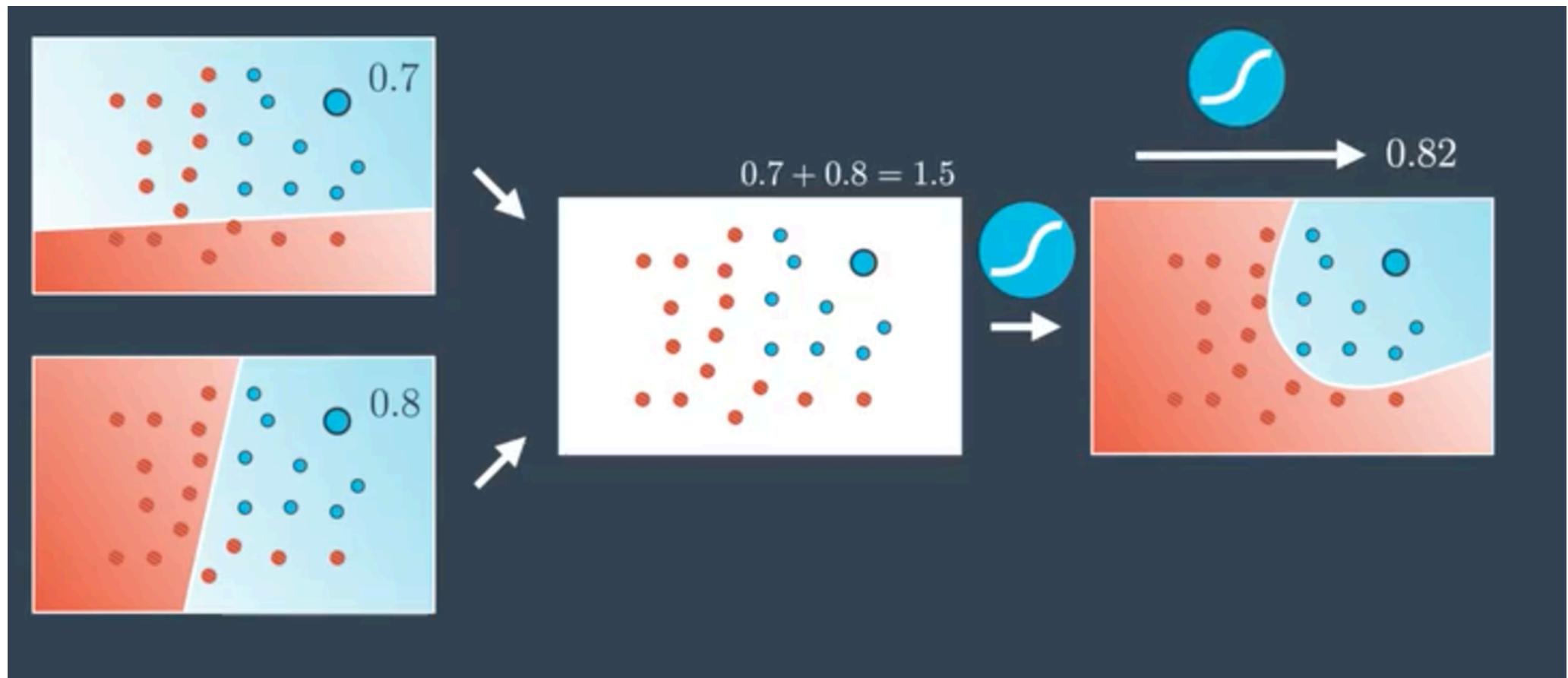
Combinando Regiões



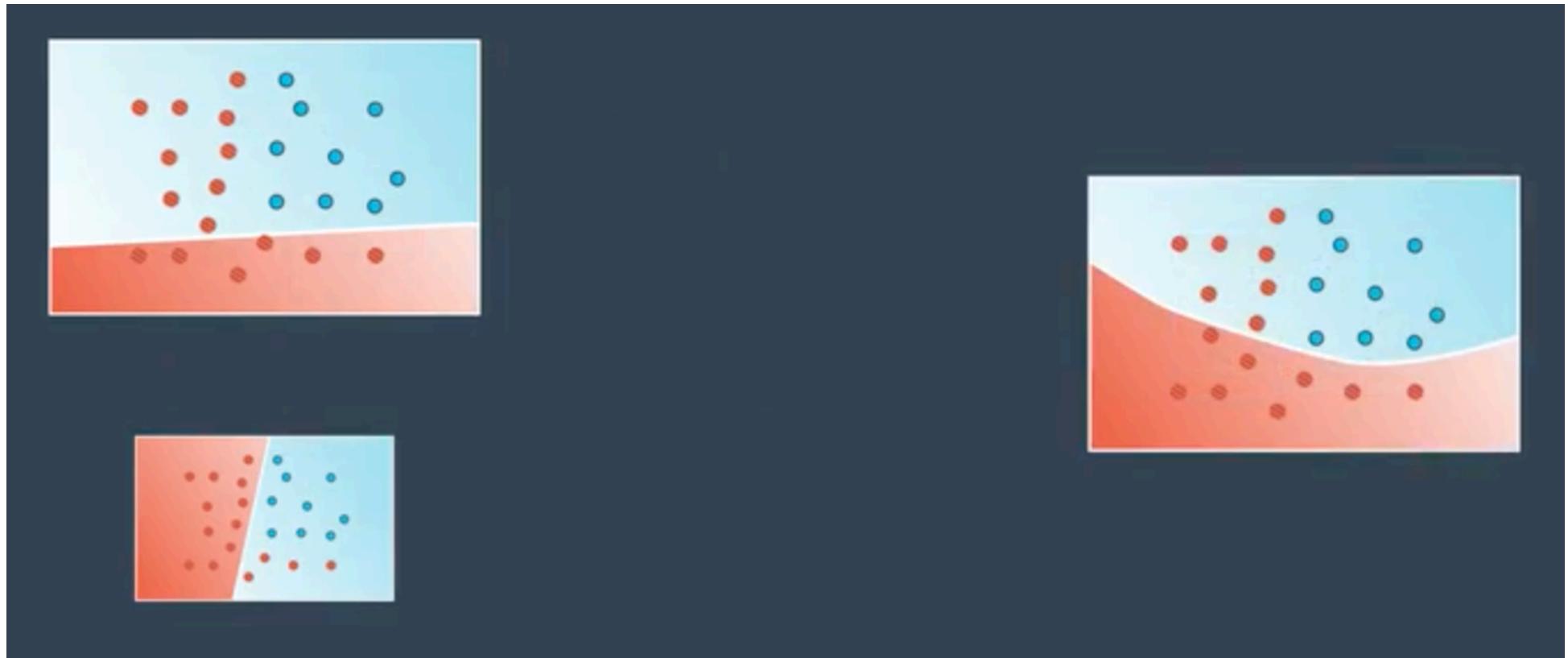
Combinando Regiões



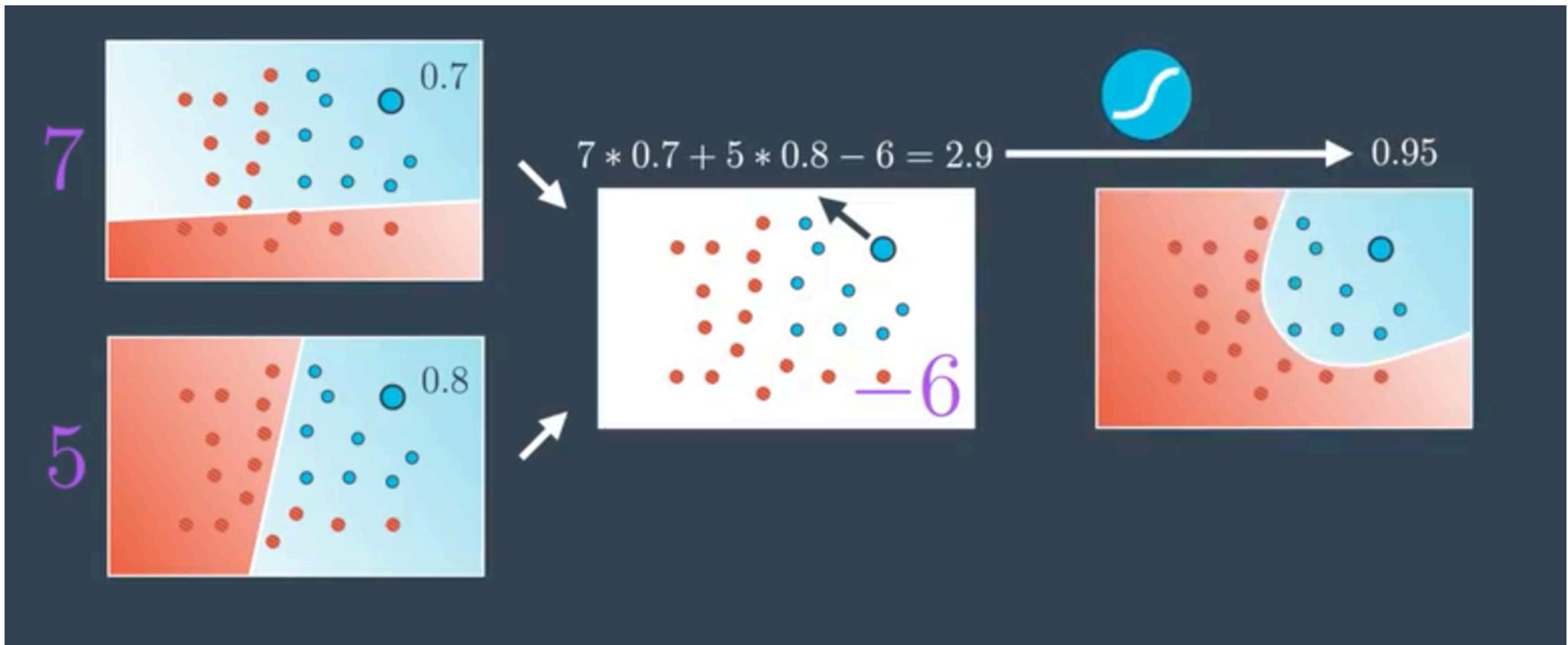
Combinando Regiões



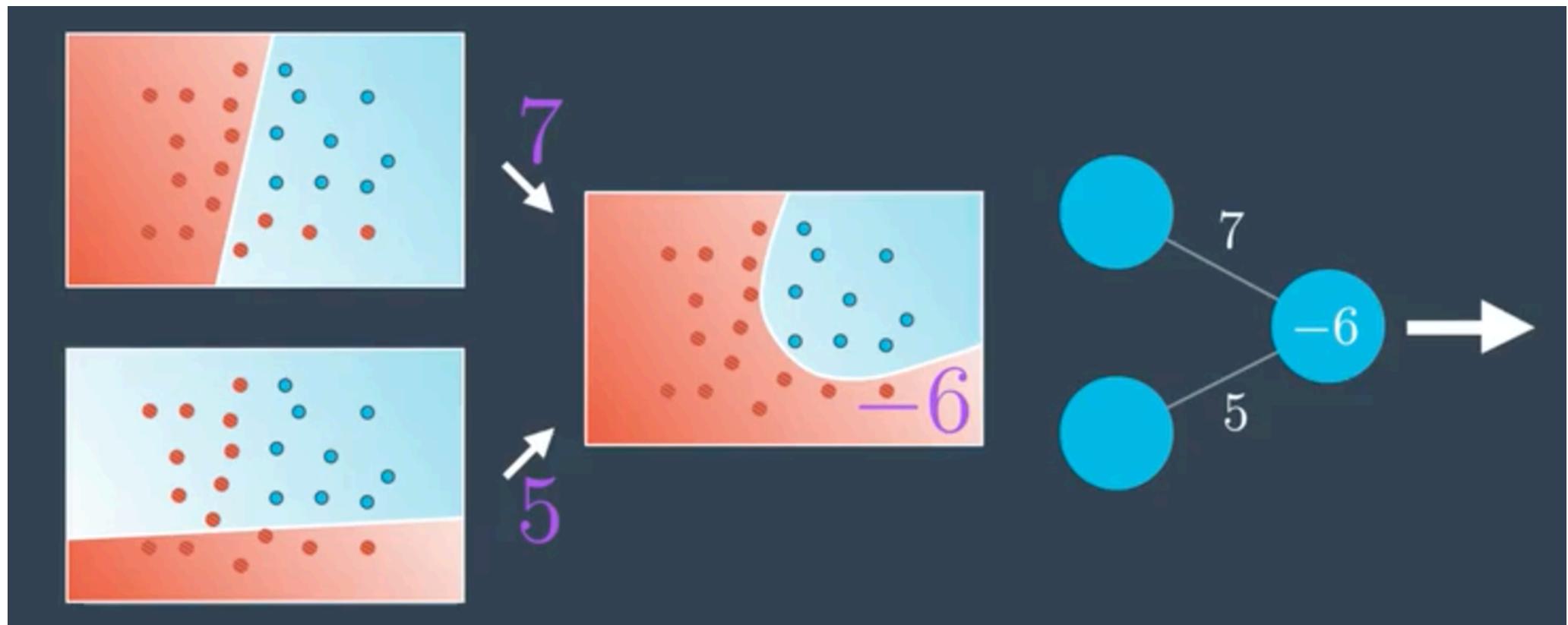
Combinando Regiões



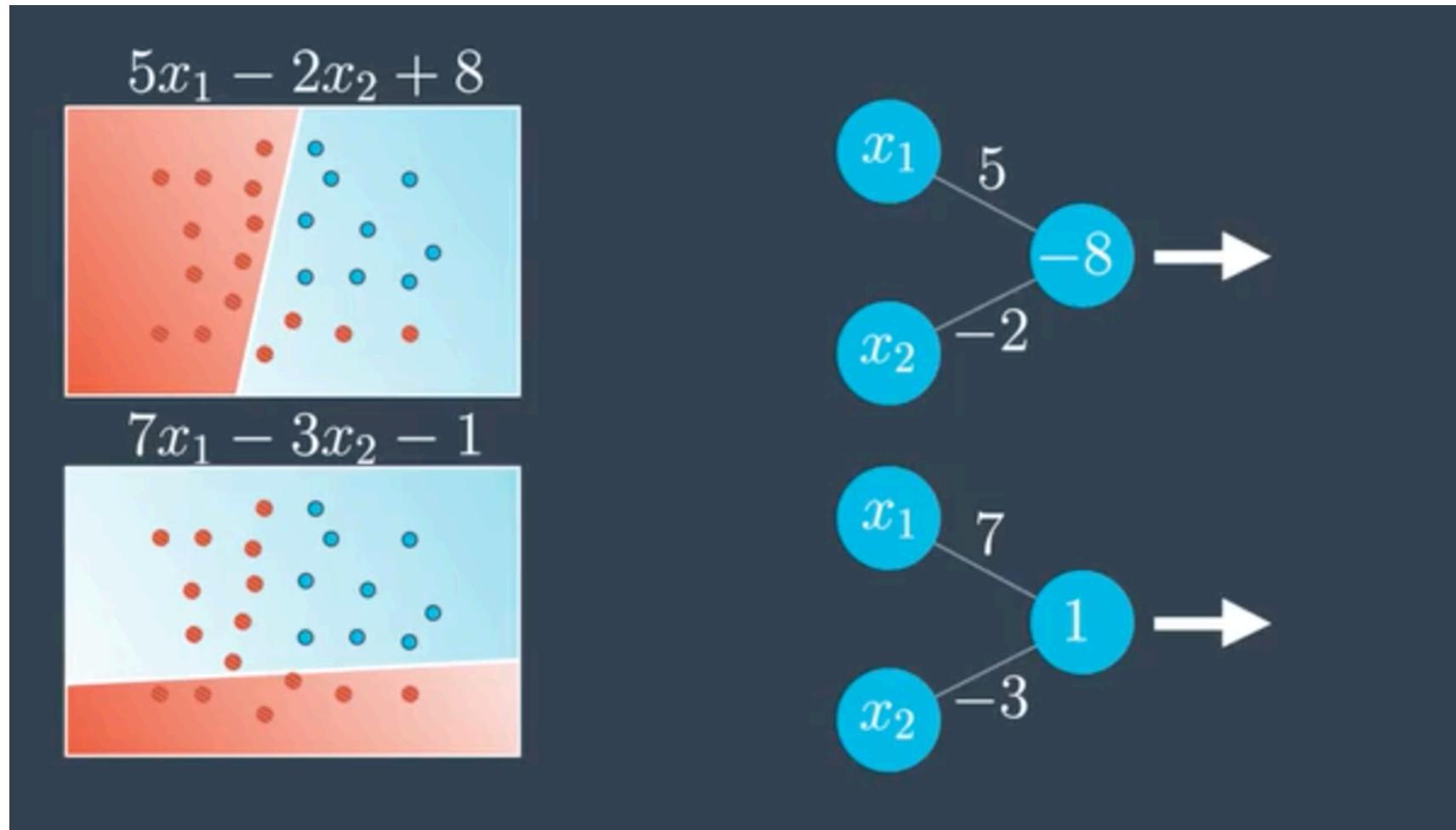
Combinando Regiões



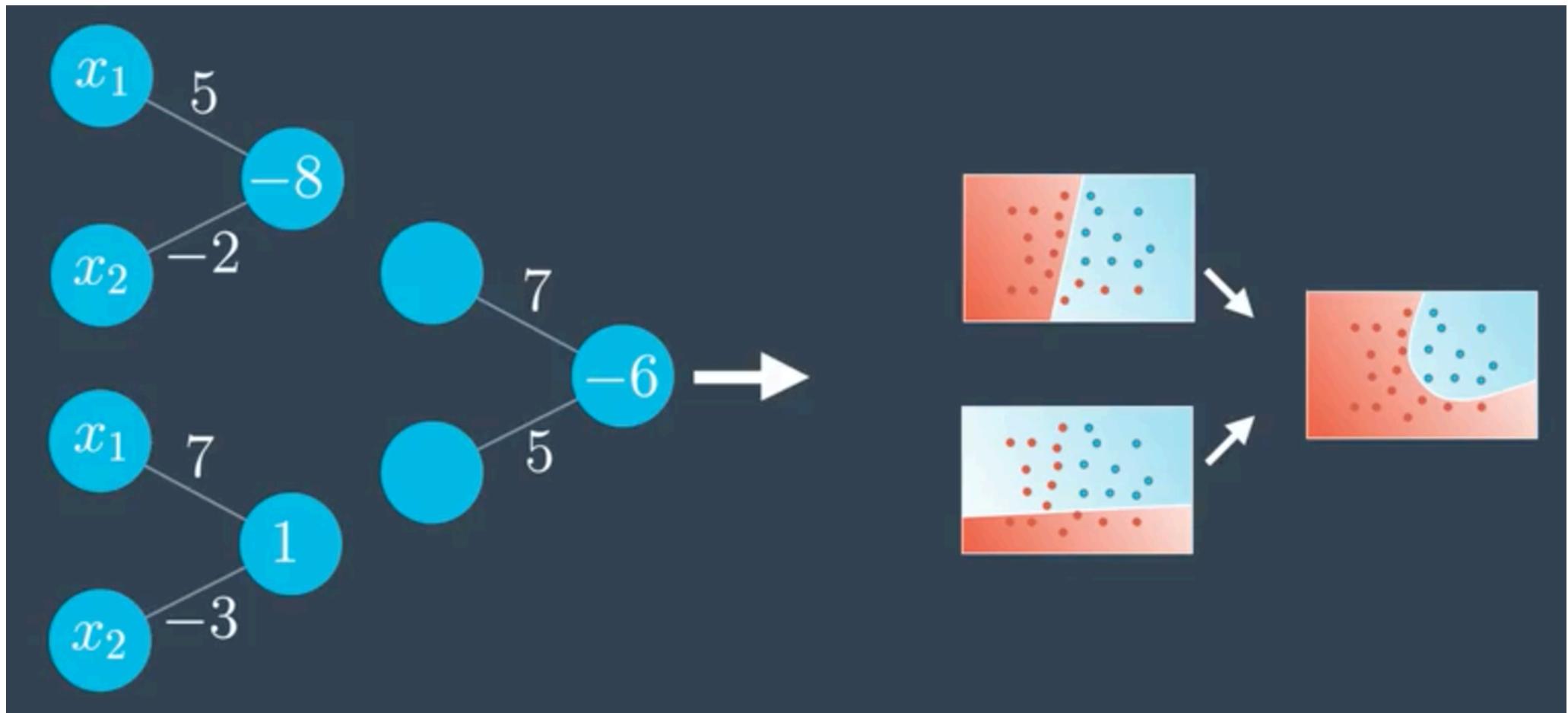
Combinando Regiões



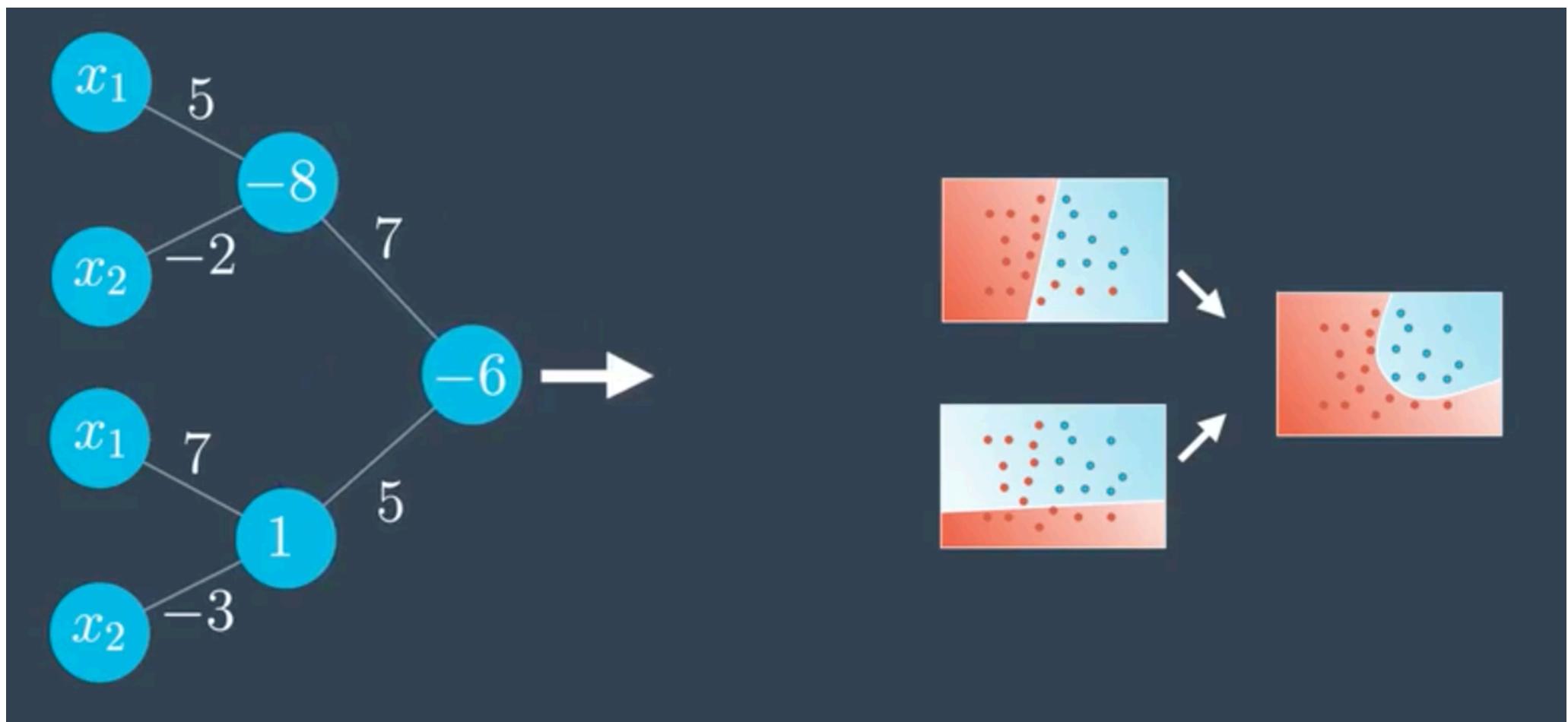
Combinando Regiões



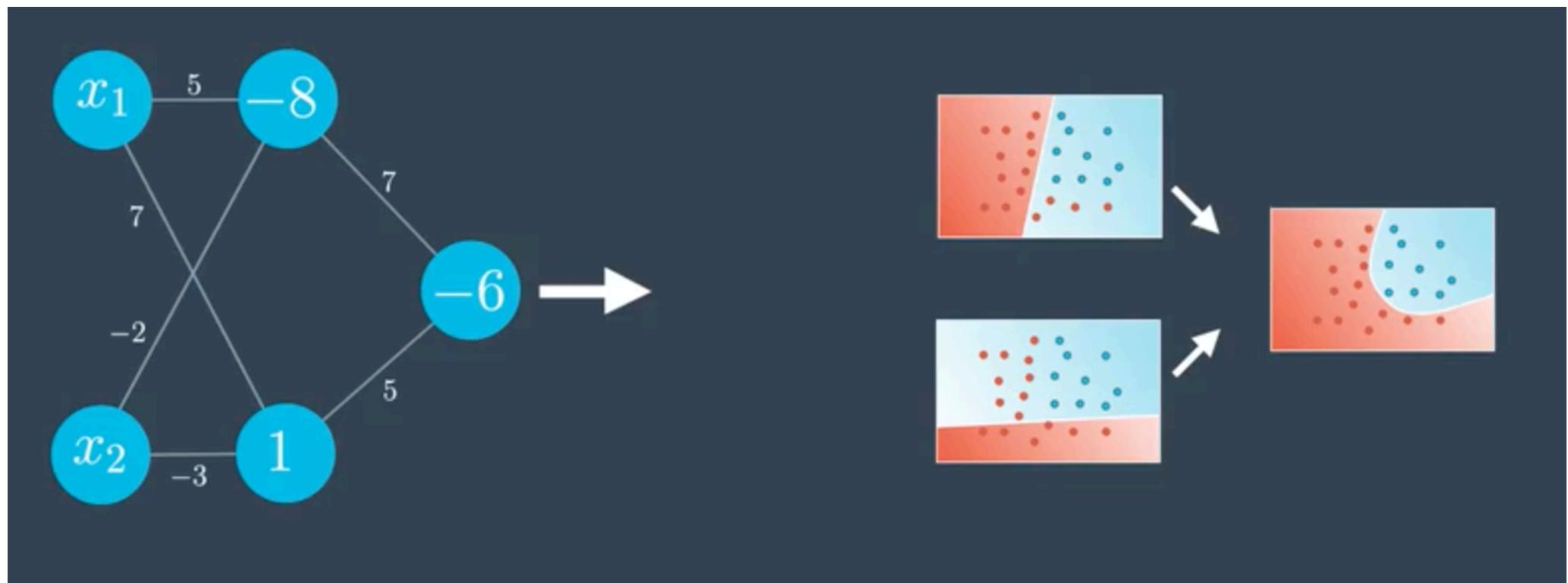
Combinando Regiões



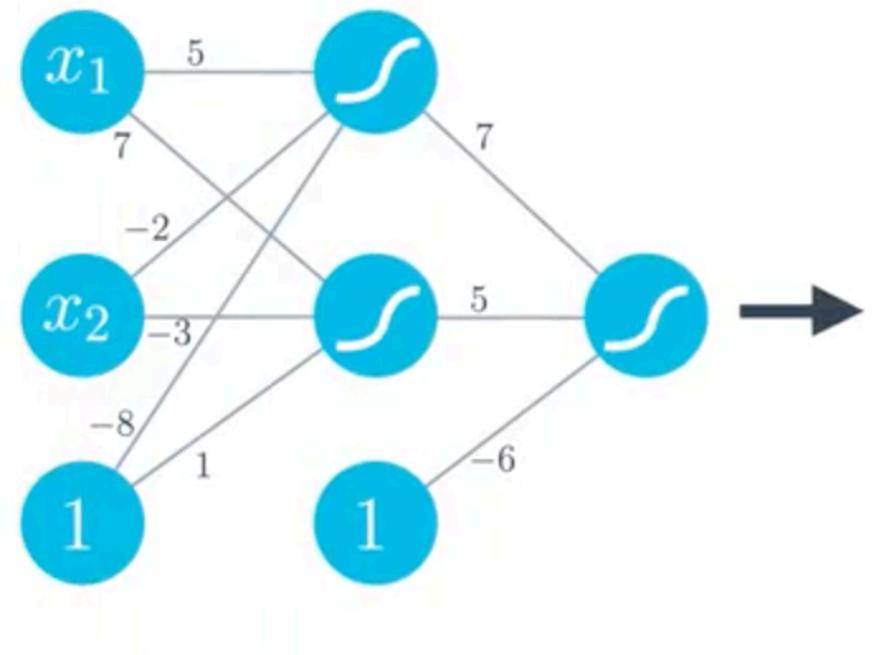
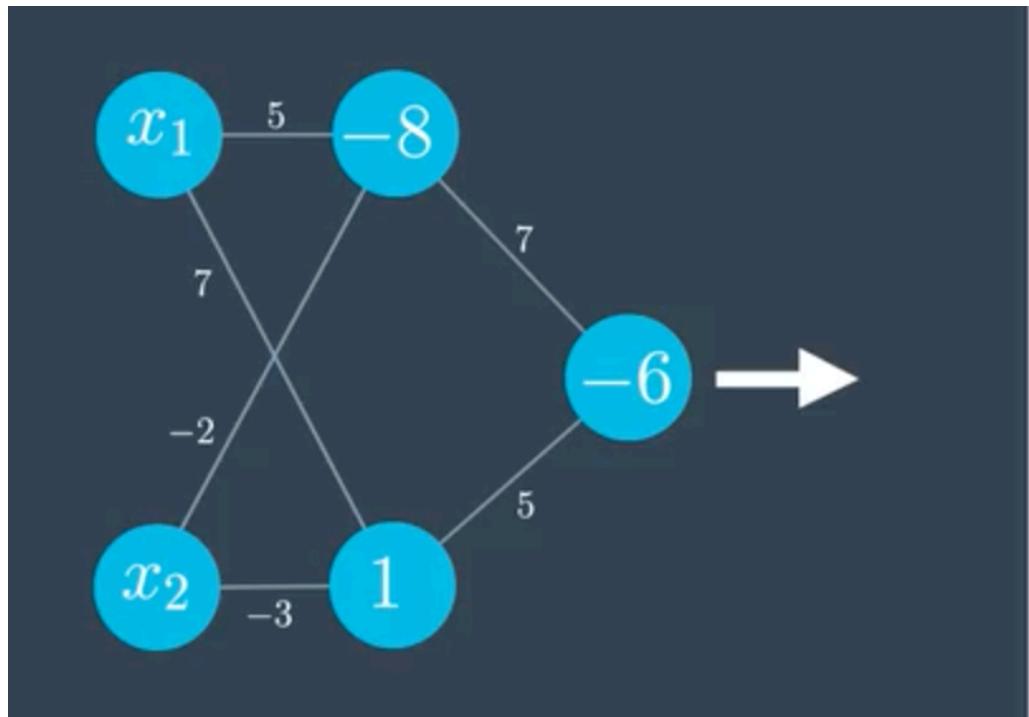
Combinando Regiões



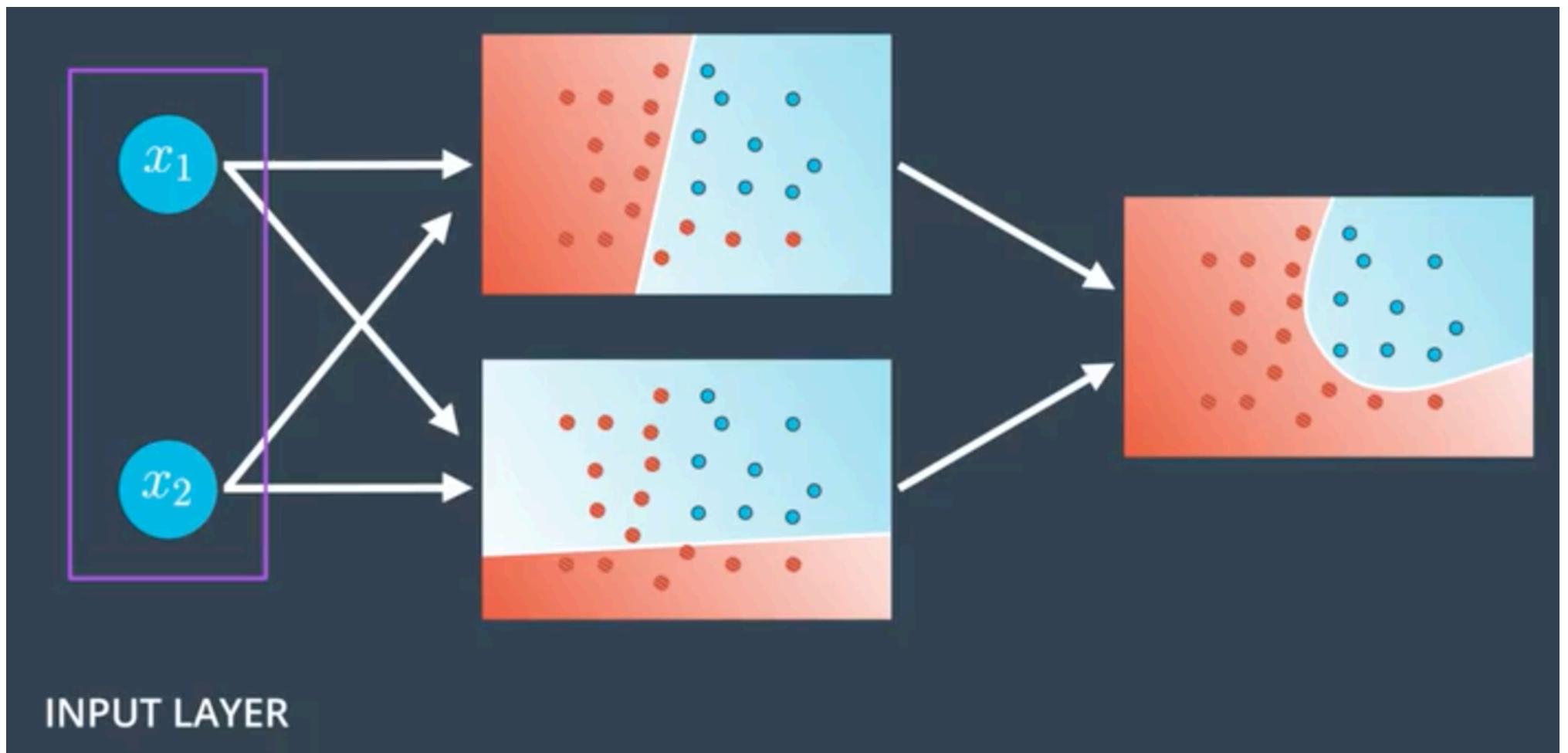
Combinando Regiões



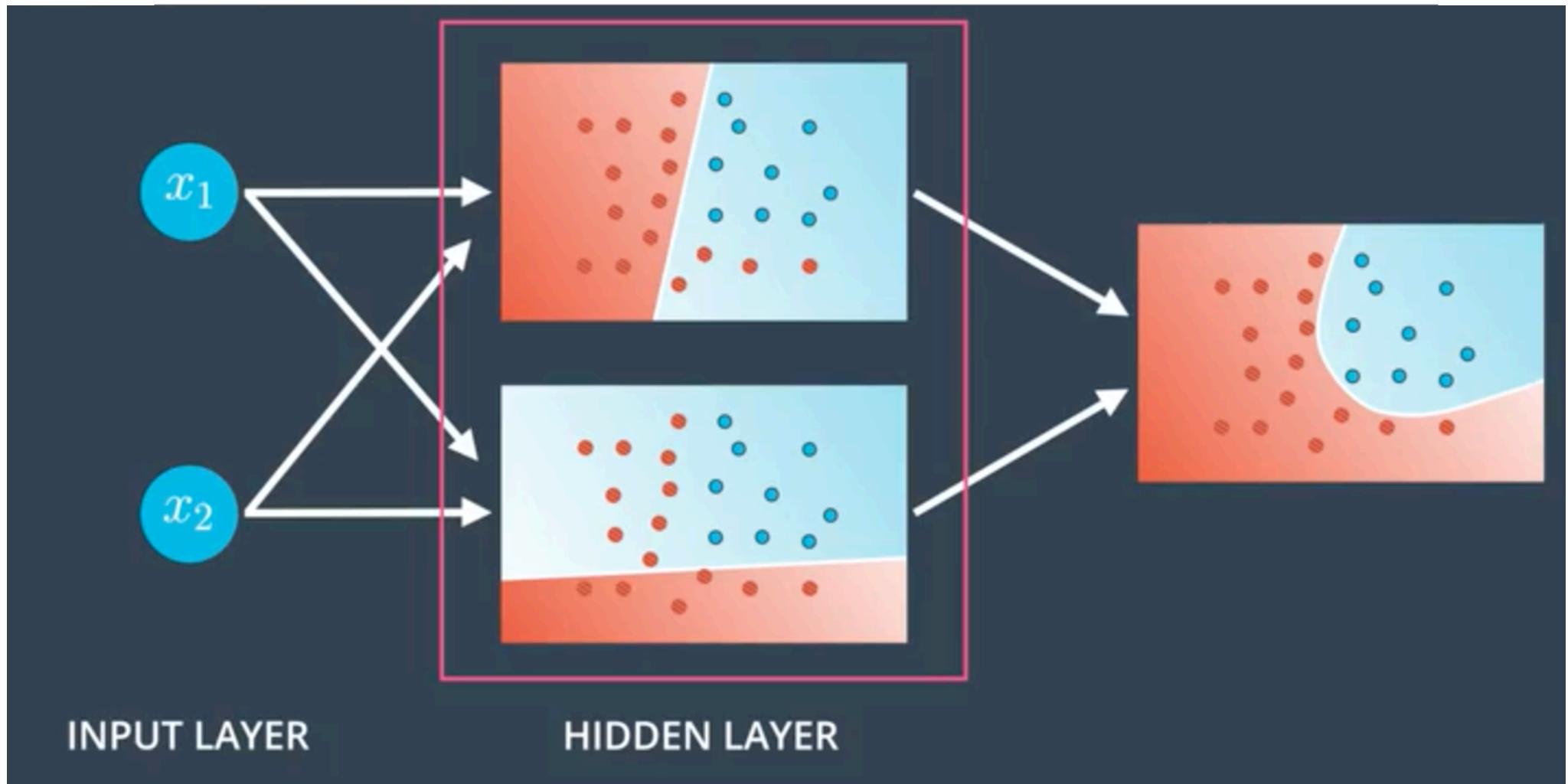
Combinando Regiões



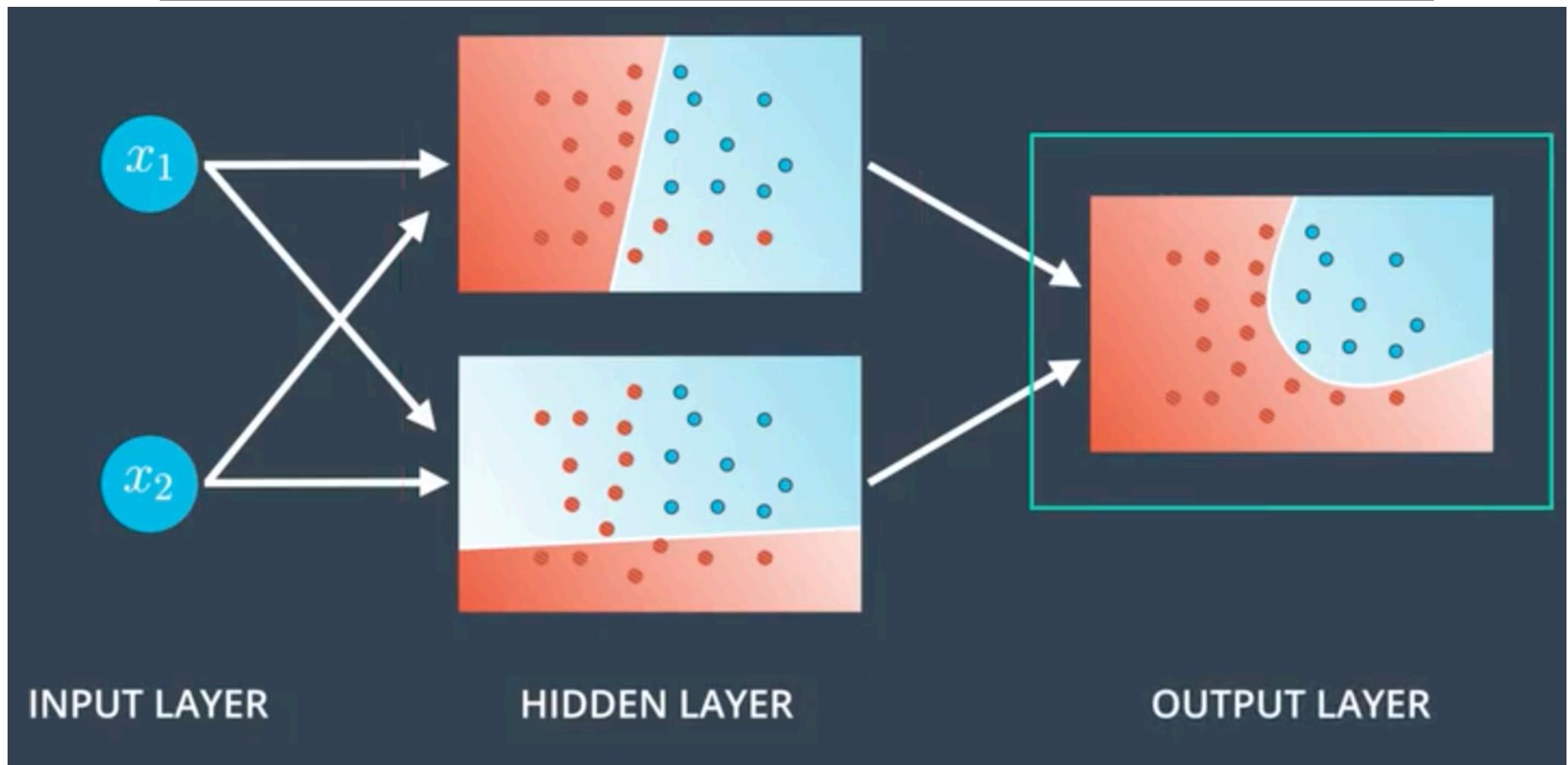
Multilayer Perceptron (MLP)

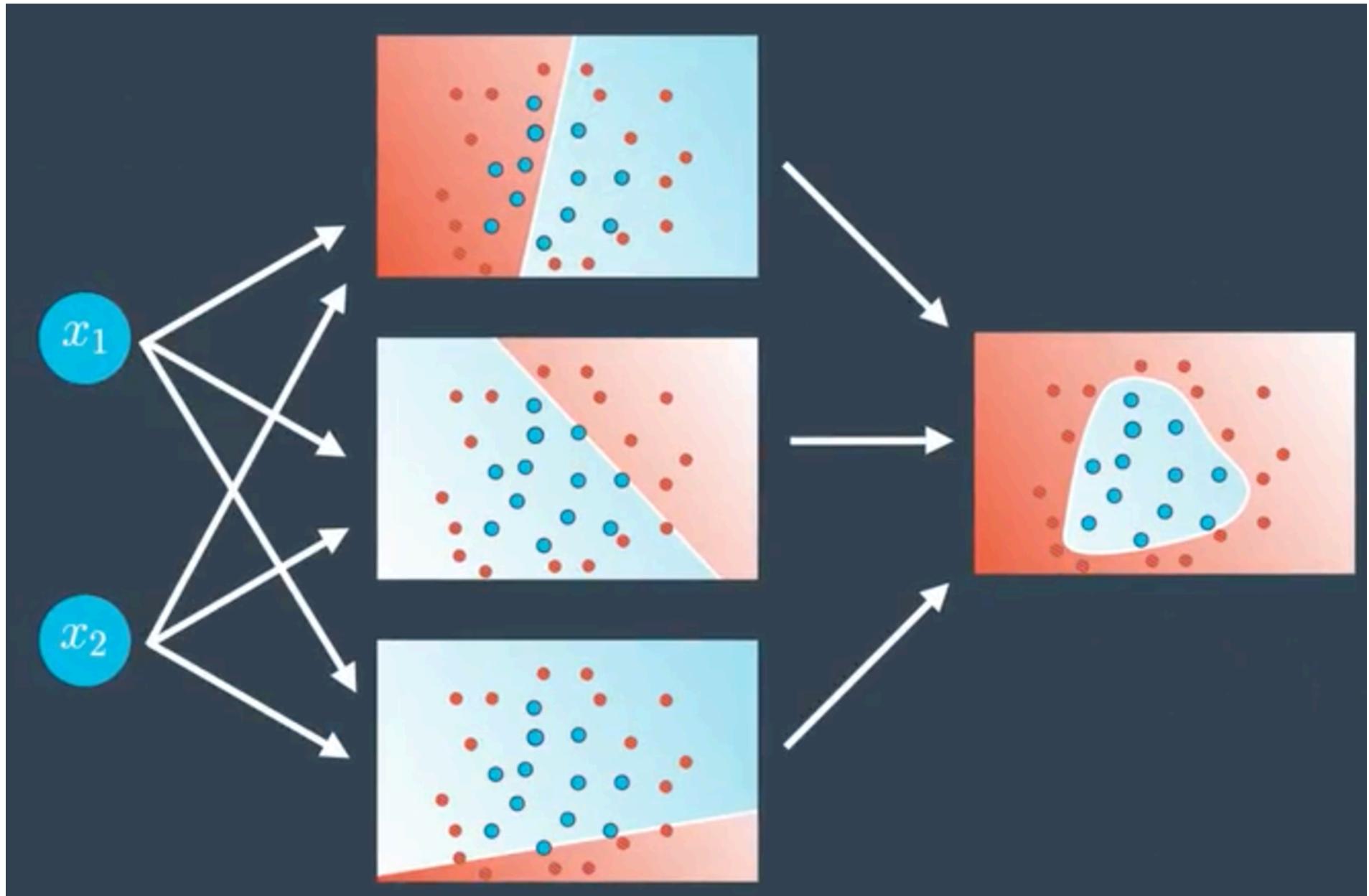


Multilayer Perceptron (MLP)

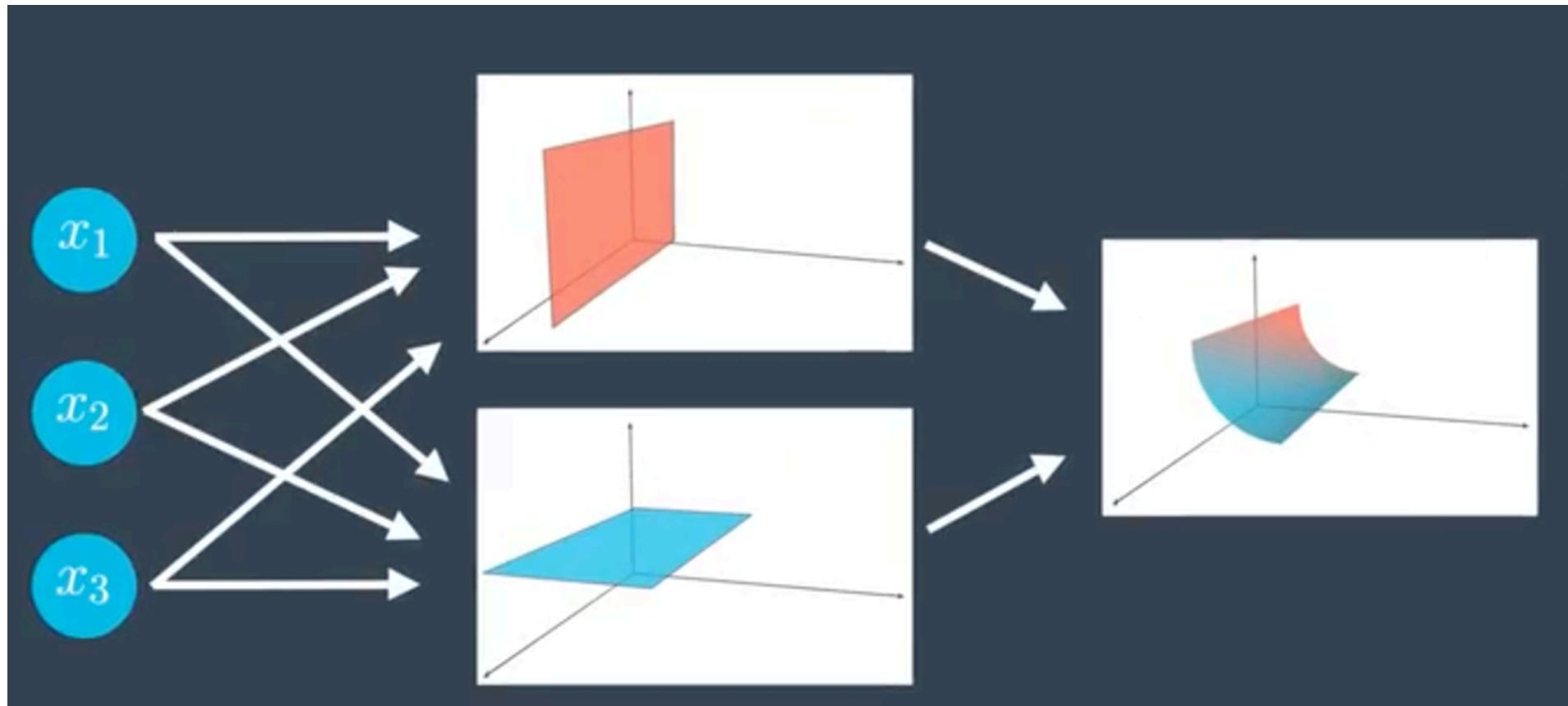


Multilayer Perceptron (MLP)

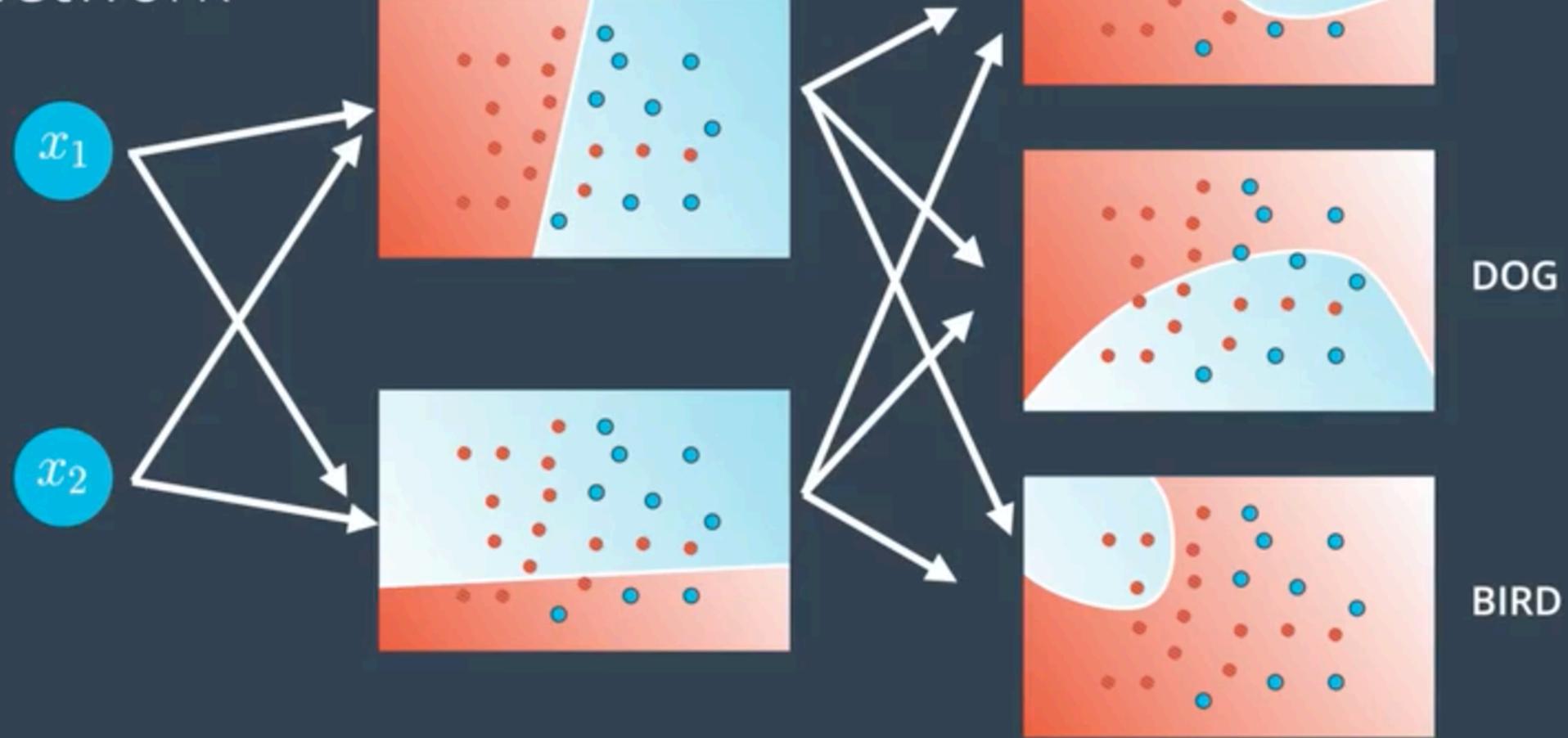




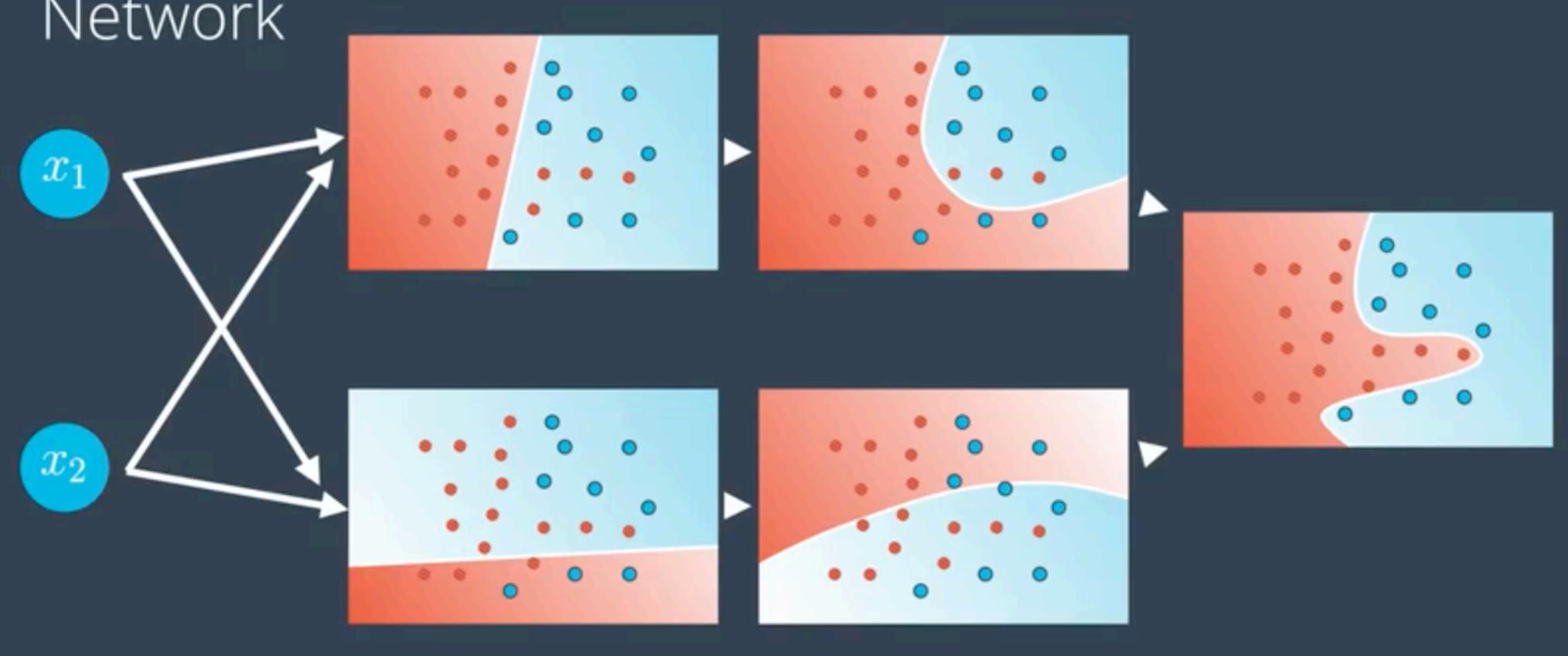
Multilayer Perceptron (MLP)



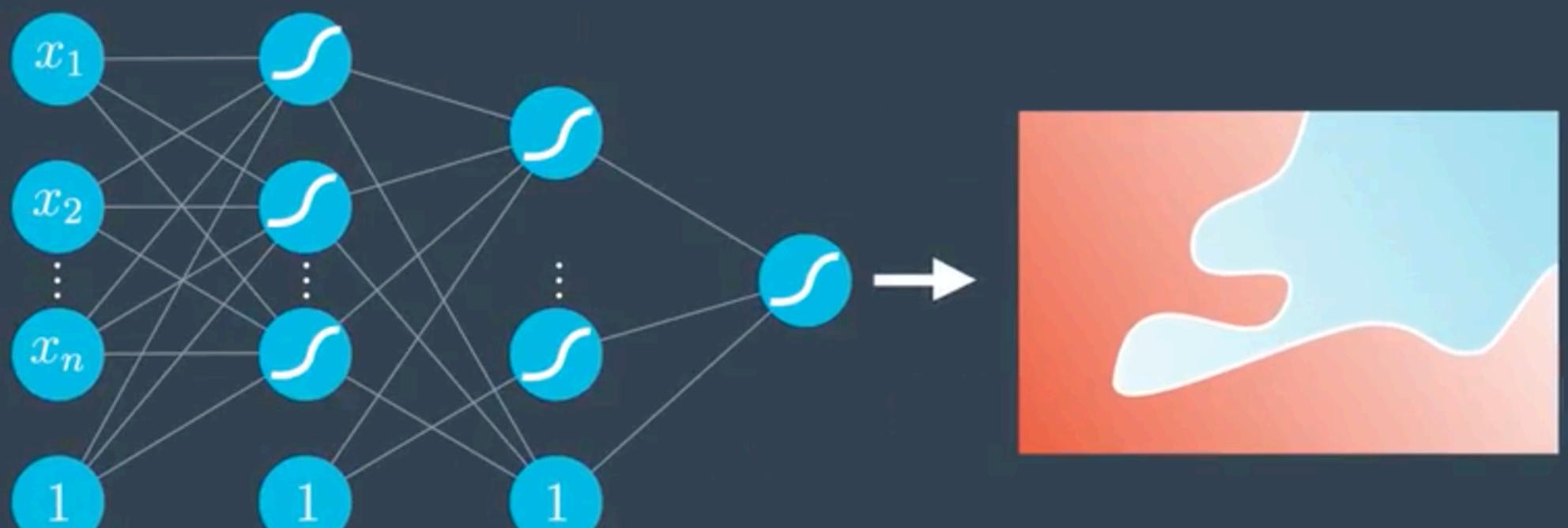
Deep Neural Network



Deep Neural Network

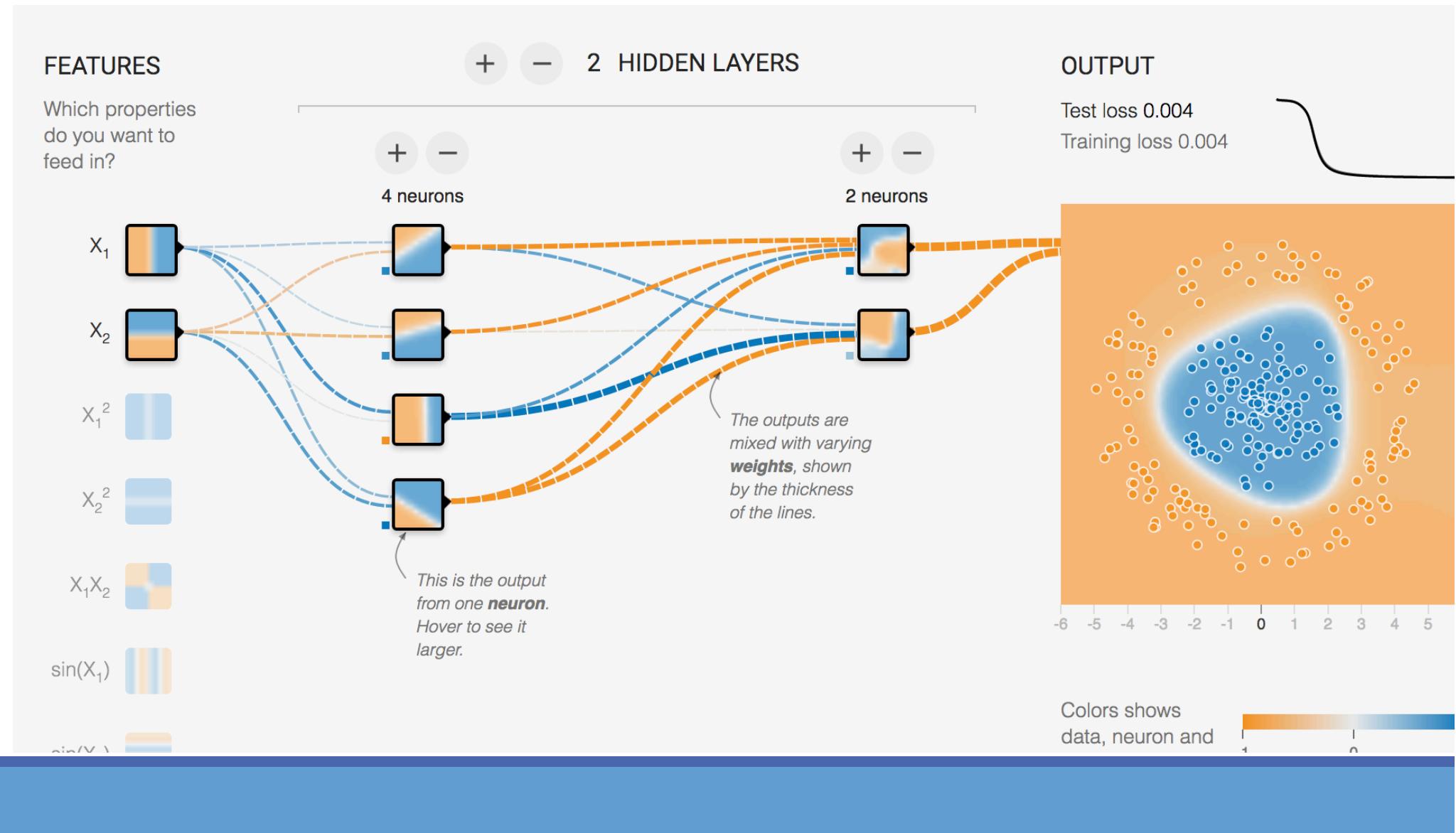


Neural Network

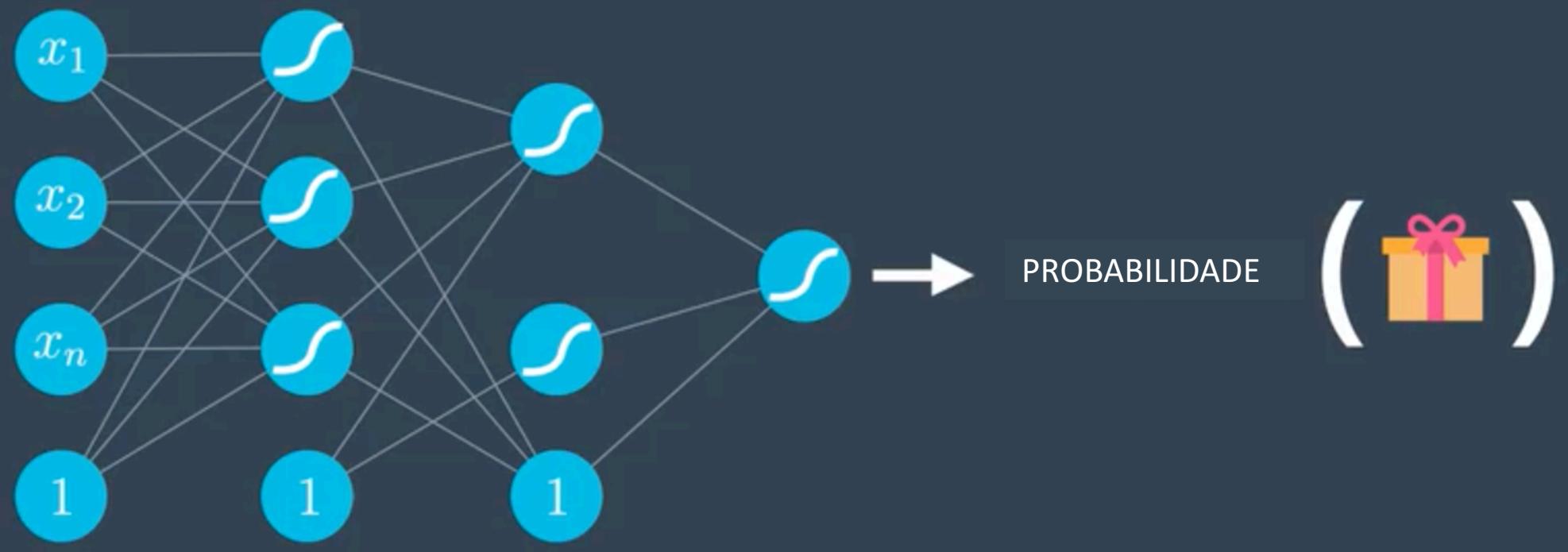


Exercicio

▷ <http://playground.tensorflow.org>



Classificação Binária



Classificação Multi-Classe



Pato



Castor



Morsa

Classificação Multi-Classe



PROBABILIDADE



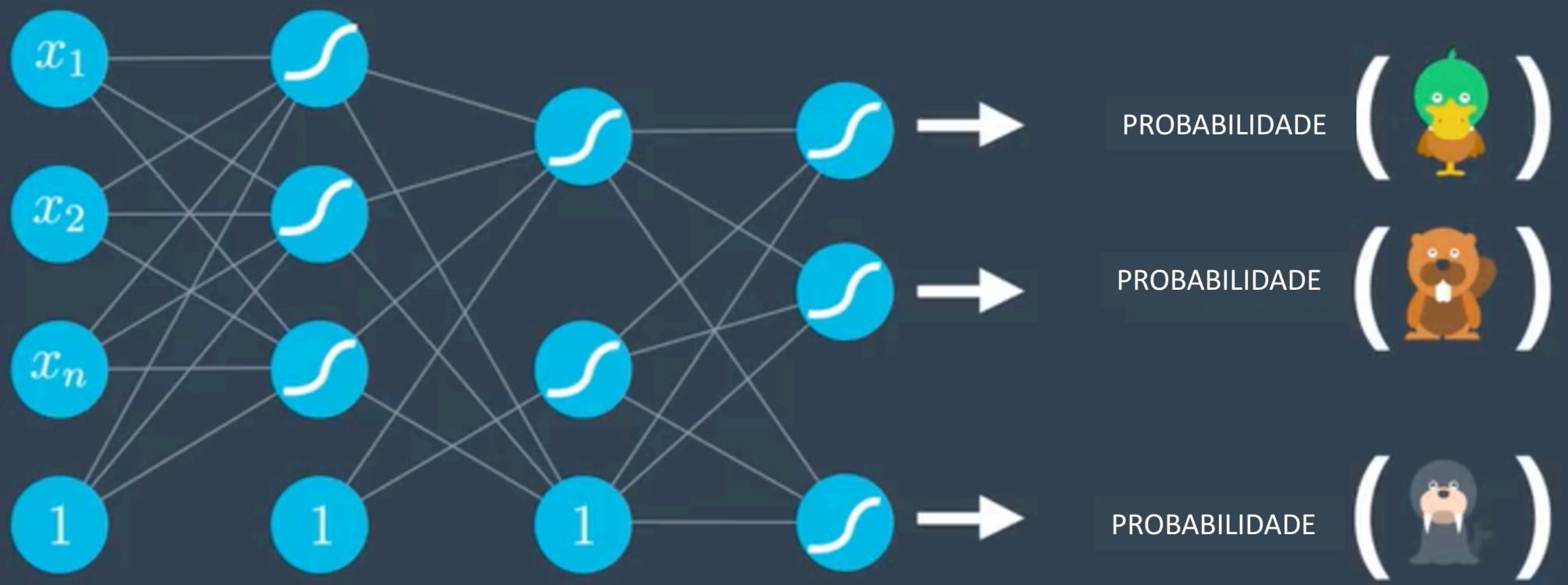
PROBABILIDADE



PROBABILIDADE



Neural Network



Classificação Multi-Classe



$P(\text{Pato}) =$
0.67



$P(\text{Castor}) =$
0.24



$P(\text{Morsa}) =$
0.09

Função Softmax

$$P(\text{classe } i) = \frac{e^{z_i}}{e^{z_1} + \dots + e^{z_n}}$$

z_i = resultado do somatório do neurônio, antes de aplicar a função de ativação



$$\frac{e^2}{e^2 + e^1 + e^0} = 0.67 \quad P(\text{Pato})$$



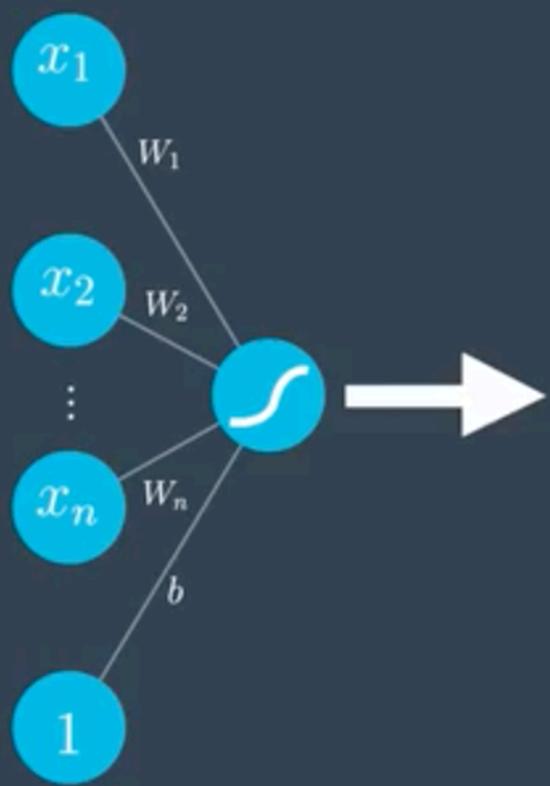
$$\frac{e^1}{e^2 + e^1 + e^0} = 0.24 \quad P(\text{Castor})$$



$$\frac{e^0}{e^2 + e^1 + e^0} = 0.09 \quad P(\text{Morsa})$$

Revisando

Perceptron

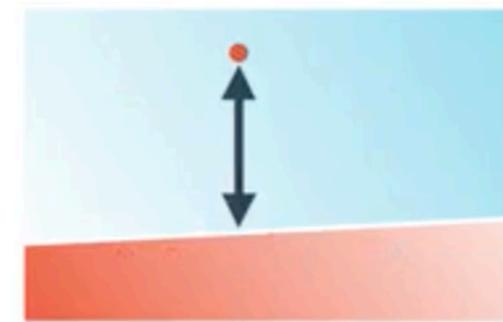


PREDIÇÃO

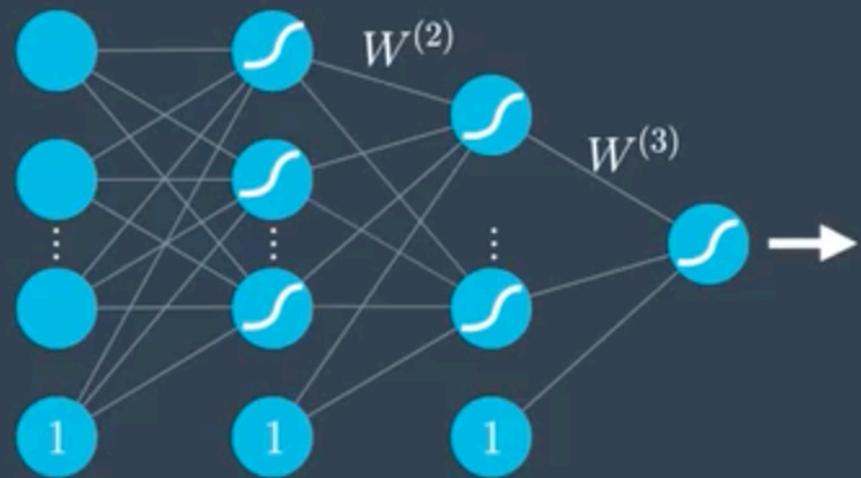
$$\hat{y} = \sigma(Wx + b)$$

FUNÇÃO DE ERRO

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$



Multi-layer Perceptron

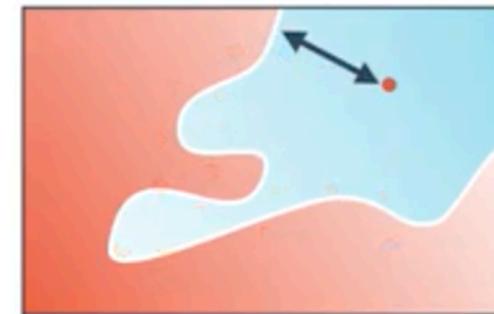


PREDIÇÃO

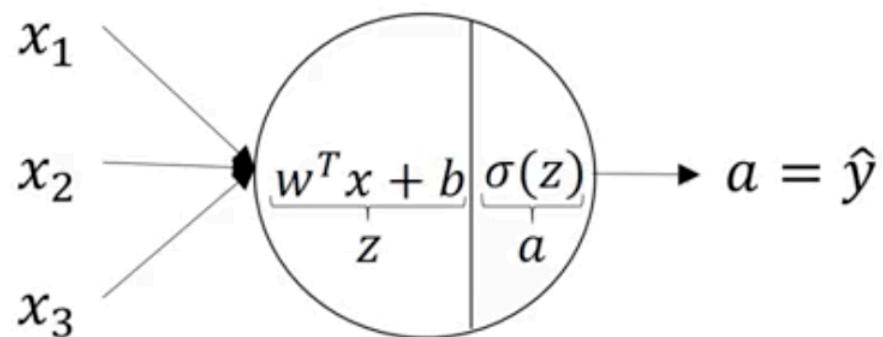
?

FUNÇÃO DE ERRO

$$E(W) = -\frac{1}{m} \sum_{i=1}^m y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$



Representação

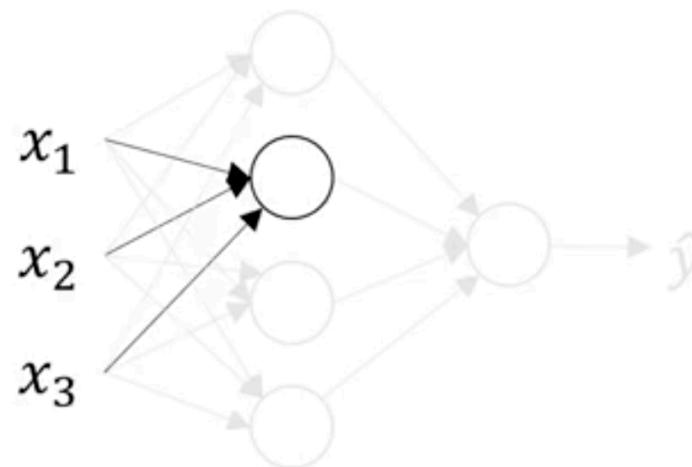
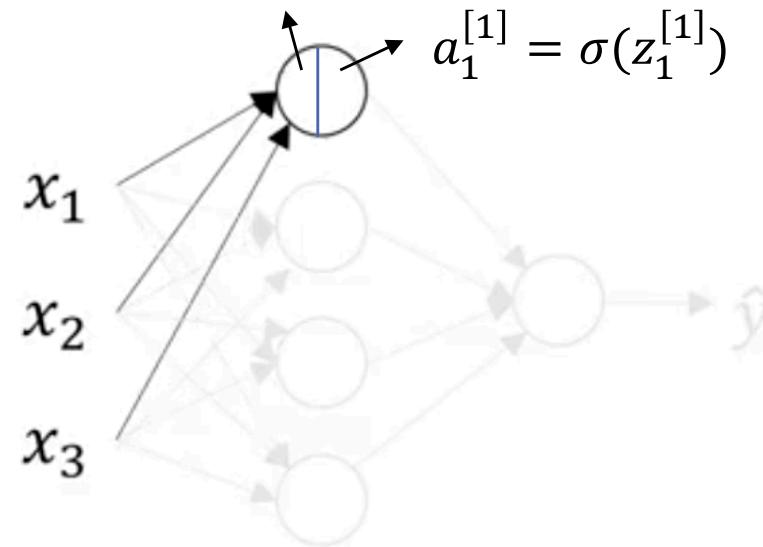


$$z = w^T x + b$$

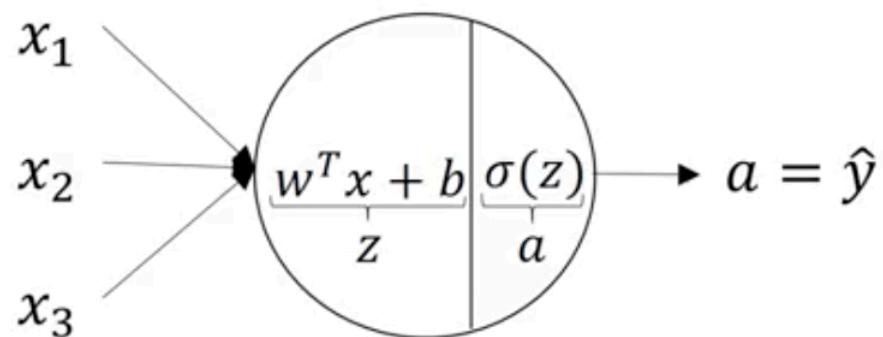
$$a = \sigma(z)$$

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$



Representação

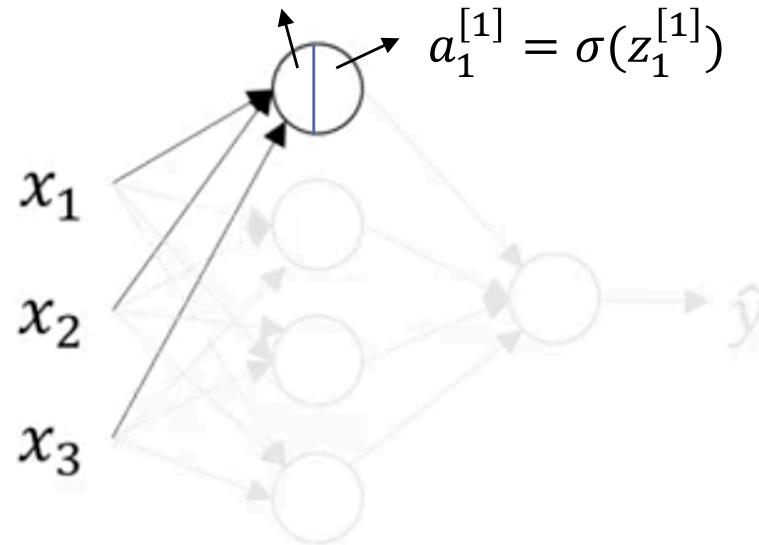


$$z = w^T x + b$$

$$a = \sigma(z)$$

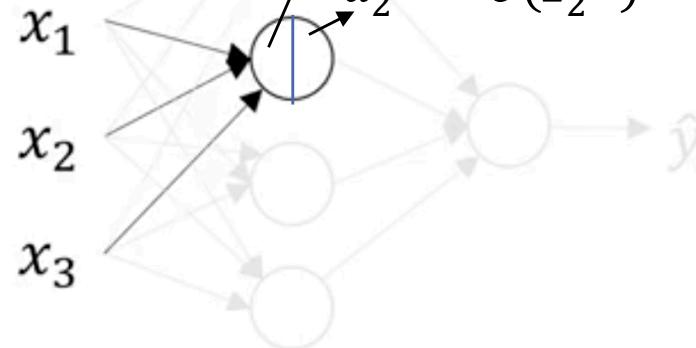
$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$$

$$a_1^{[1]} = \sigma(z_1^{[1]})$$

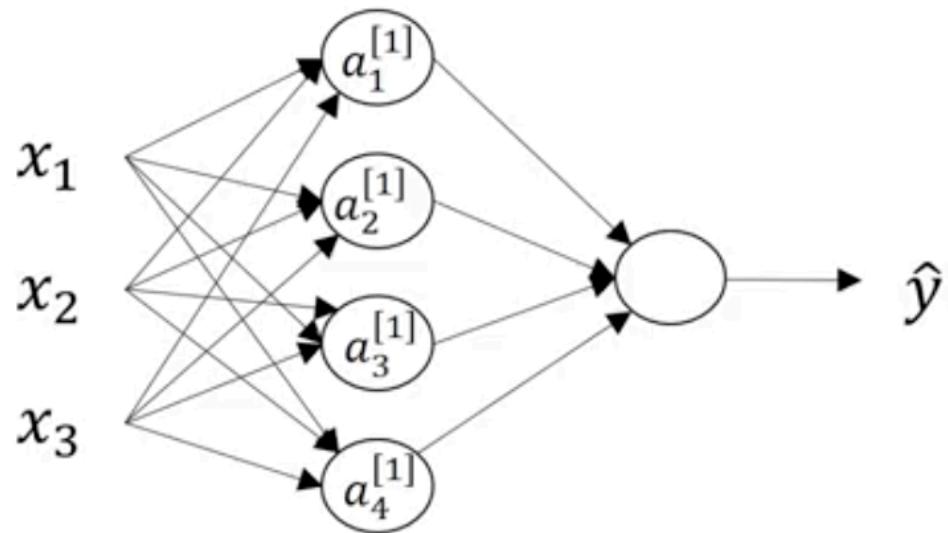


$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$

$$a_2^{[1]} = \sigma(z_2^{[1]})$$

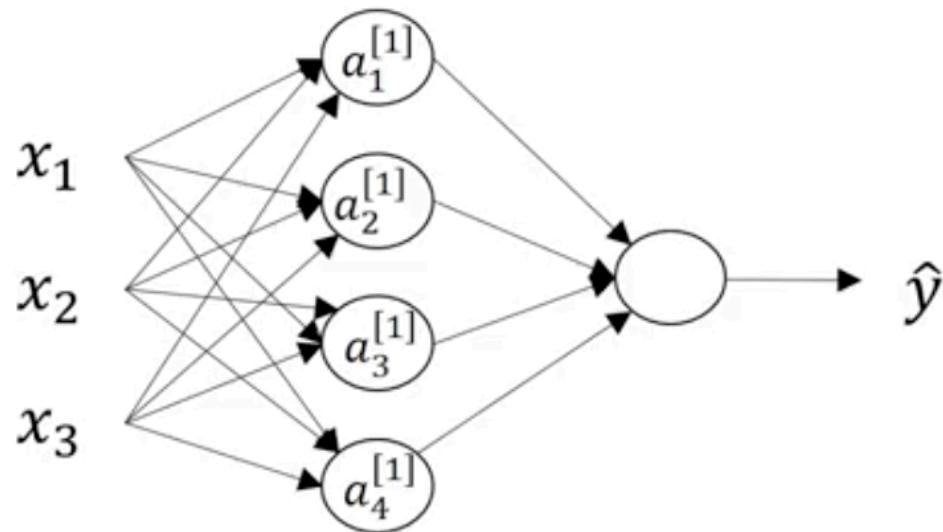


Representação



$$\begin{aligned}z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]}) \\z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]}) \\z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]}) \\z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})\end{aligned}$$

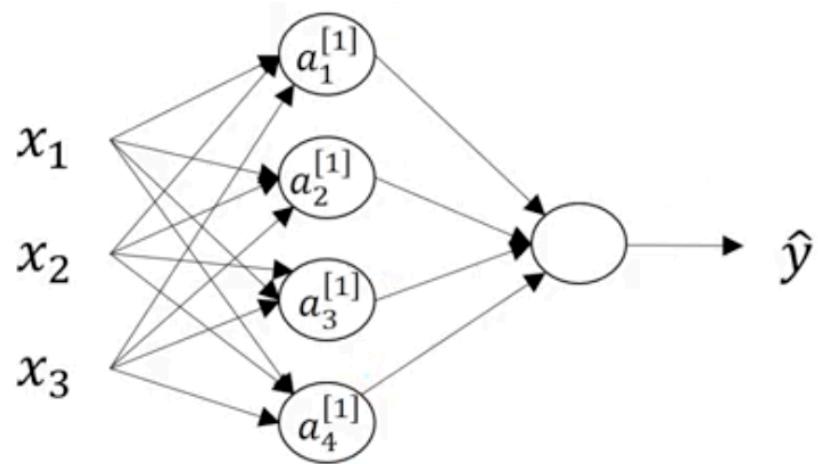
Representação



$$\begin{aligned}
 z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, a_1^{[1]} = \sigma(z_1^{[1]}) \\
 z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, a_2^{[1]} = \sigma(z_2^{[1]}) \\
 z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, a_3^{[1]} = \sigma(z_3^{[1]}) \\
 z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, a_4^{[1]} = \sigma(z_4^{[1]})
 \end{aligned}$$

$$\begin{aligned}
 W^{[1]} & & b^{[1]} & \\
 Z^{[1]} &= \begin{bmatrix} -w_1^{[1]T} & - \\ -w_2^{[1]T} & - \\ -w_3^{[1]T} & - \\ -w_4^{[1]T} & - \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} & = & \begin{bmatrix} w_1^{[1]T} x + b_1^{[1]} \\ w_2^{[1]T} x + b_2^{[1]} \\ w_3^{[1]T} x + b_3^{[1]} \\ w_4^{[1]T} x + b_4^{[1]} \end{bmatrix} & = & \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix}
 \end{aligned}$$

Representação



Dada a entrada x :

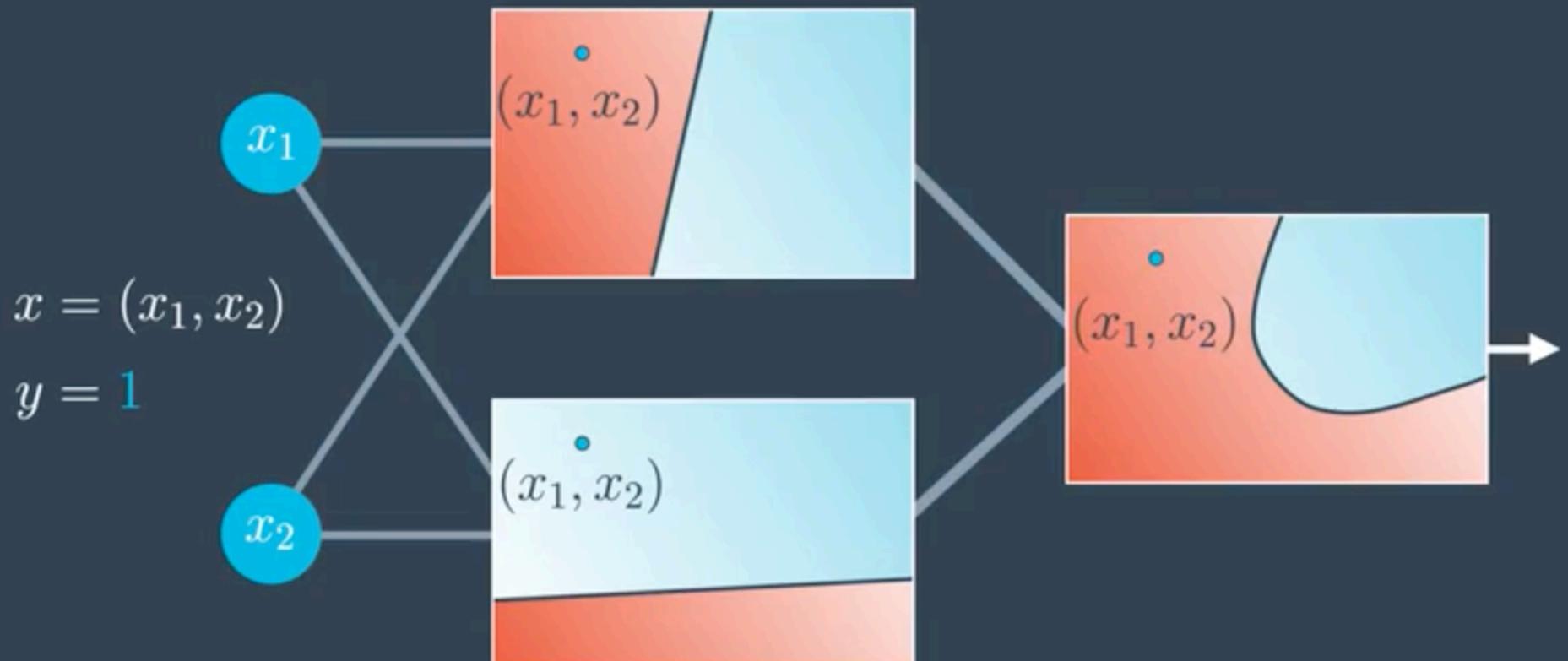
$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

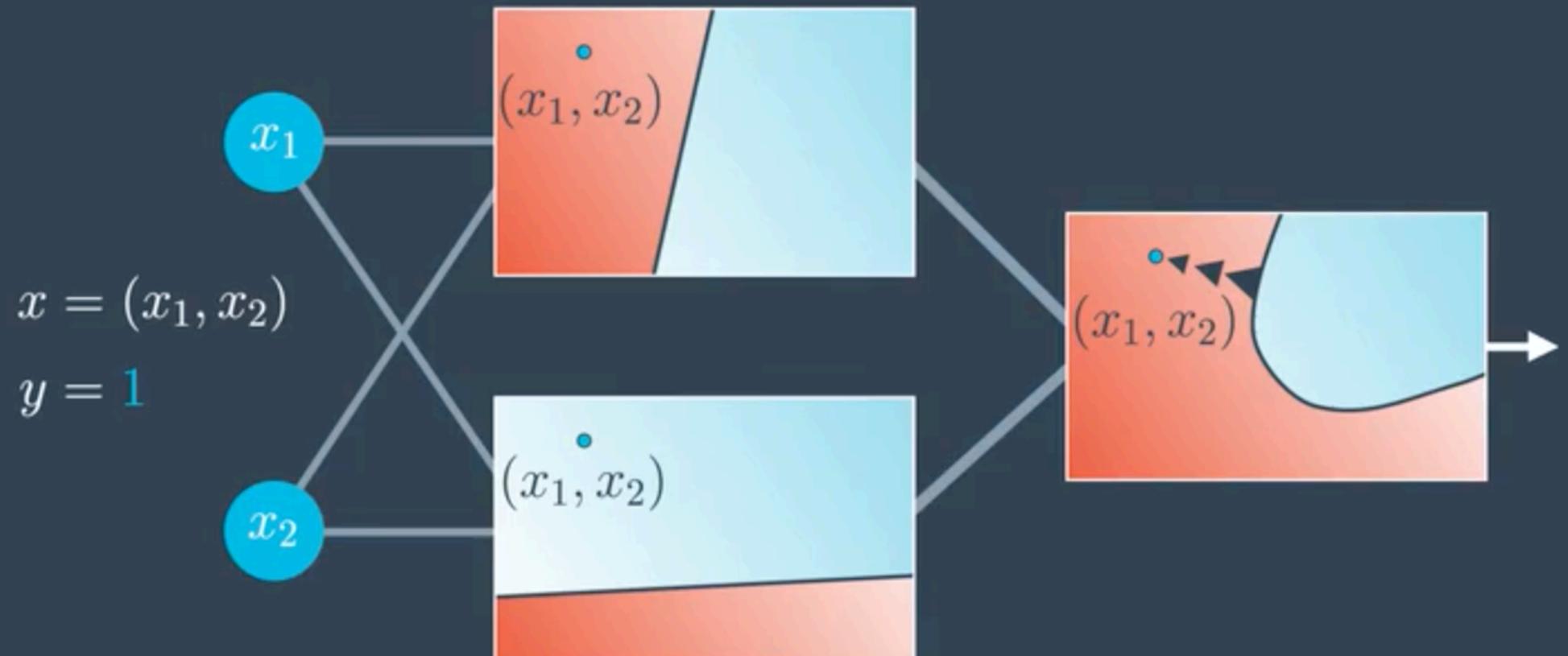
$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[2]} = \sigma(z^{[2]})$$

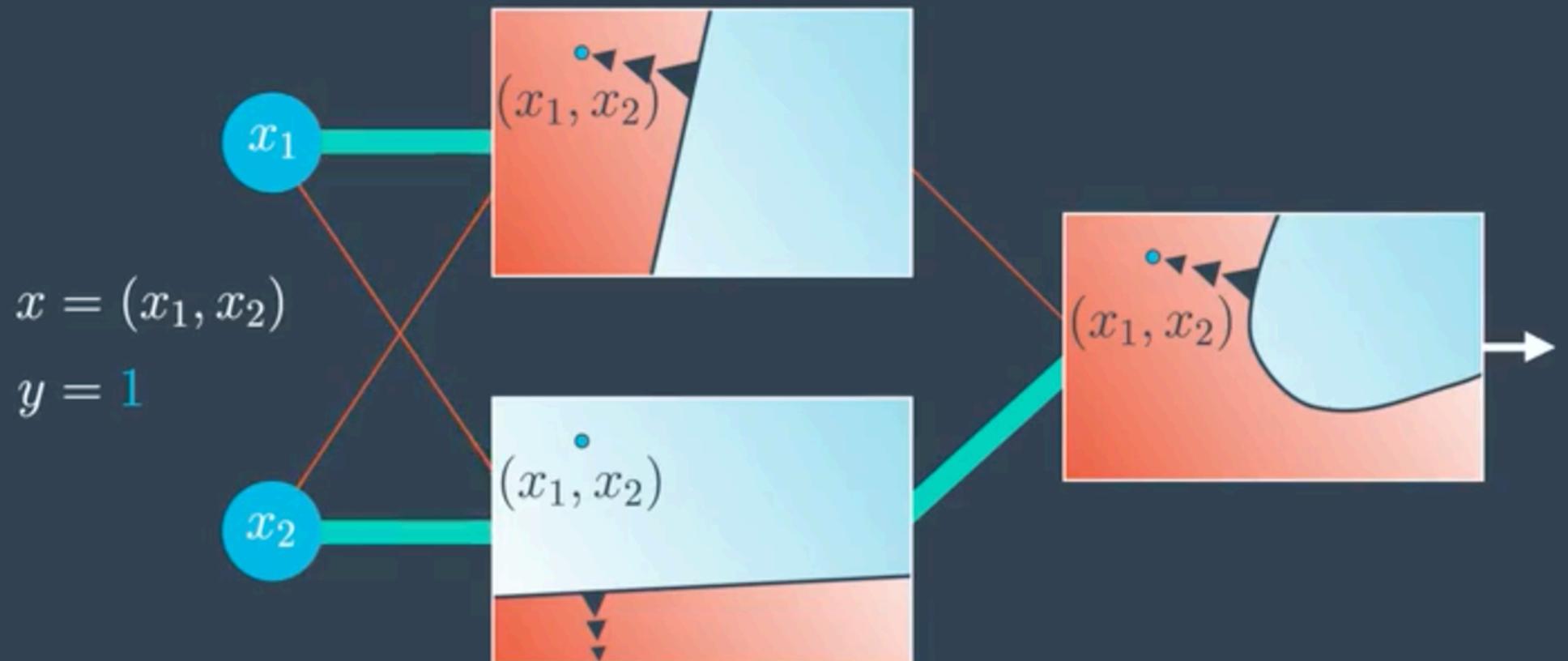
FeedForward



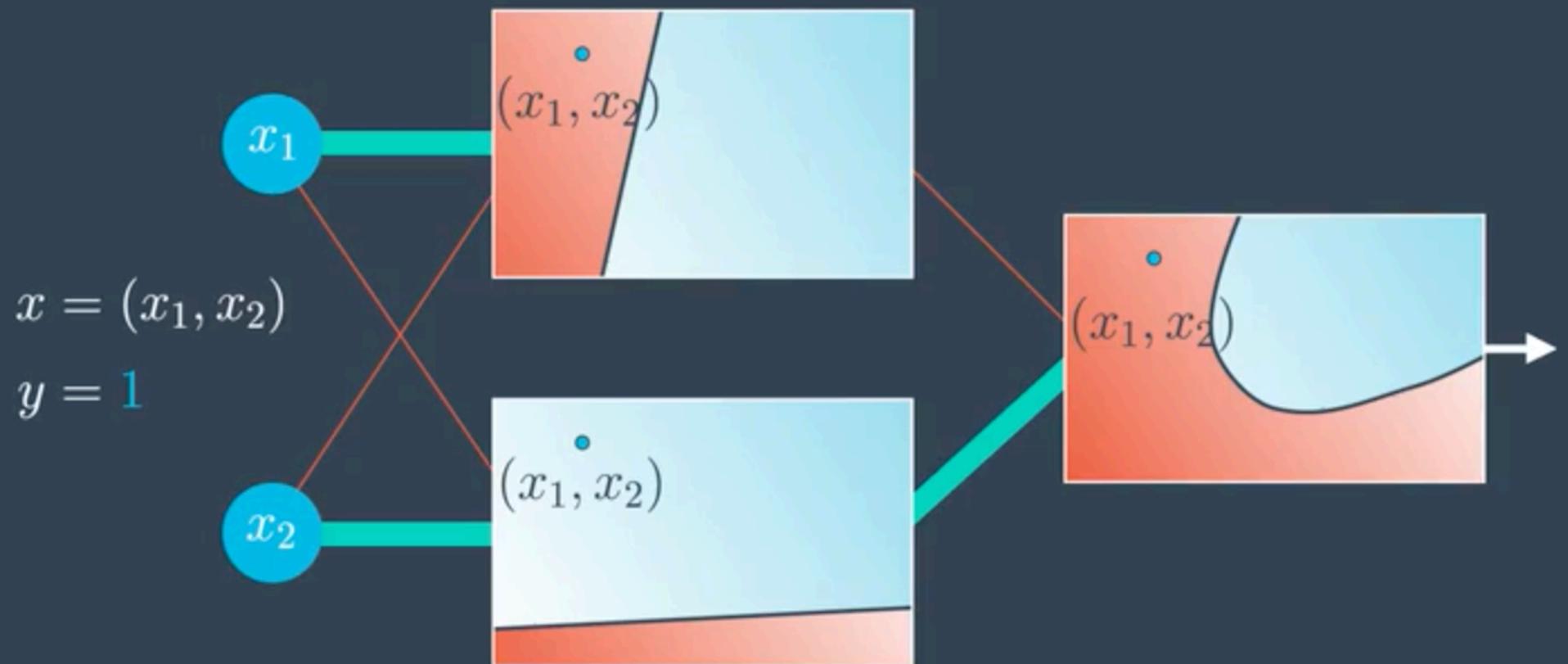
Backpropagation



Backpropagation



Backpropagation



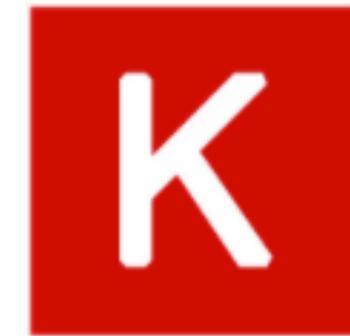
Exercício

▷ MLP_code.ipynb

MLP com Keras



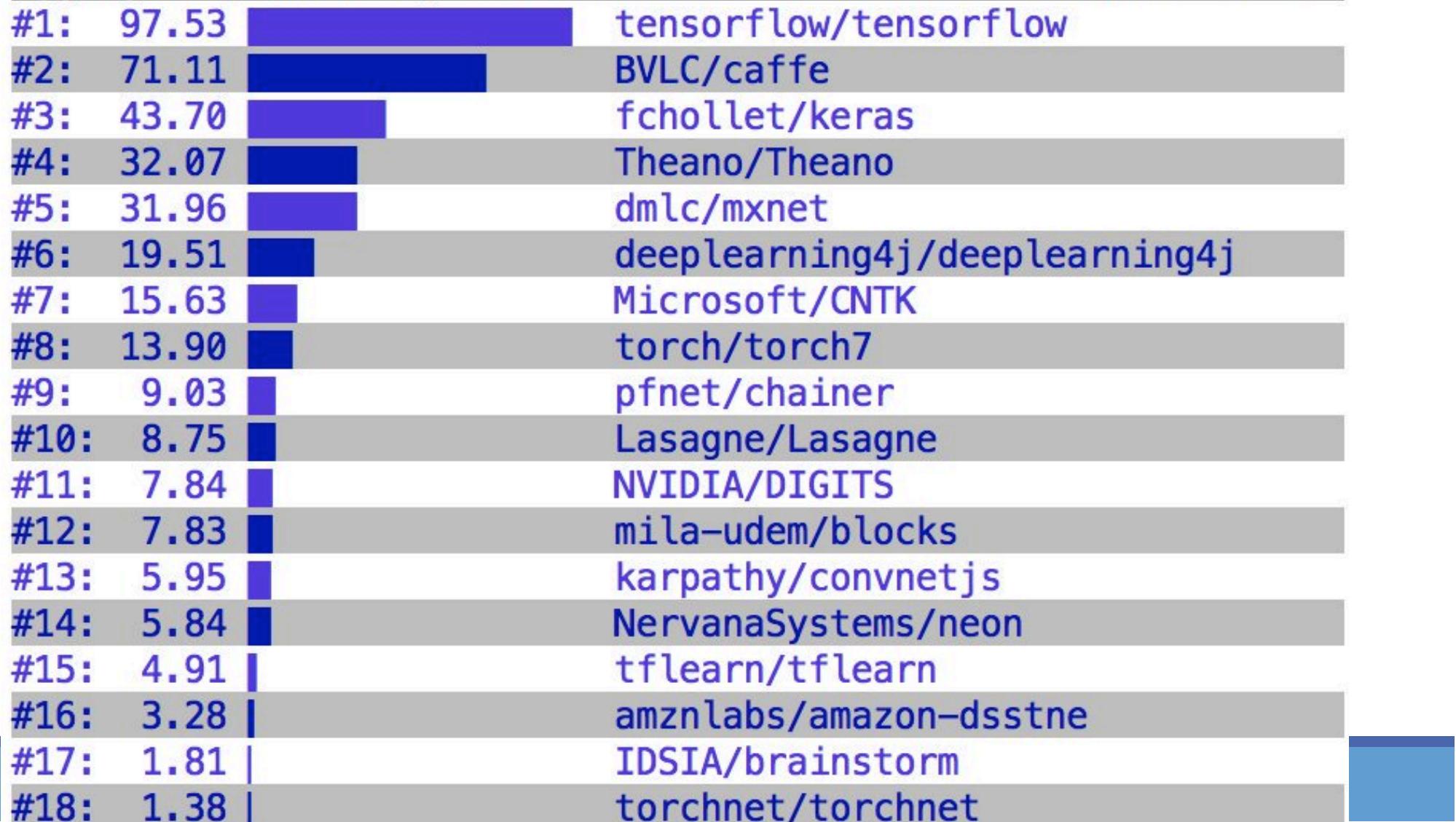
Frameworks



DEEPMLEARNING4J

Popularidade

Aggregate popularity ($30 \cdot \text{contrib} + 10 \cdot \text{issues} + 5 \cdot \text{forks}) \cdot 1e-3$



	Languages	Tutorials and training materials	CNN modeling capability	RNN modeling capability	Architecture: easy-to-use and modular front end	Speed	Multiple GPU support	Keras compatible
Theano	Python, C++	++	++	++	+	++	+	+
Tensor-Flow	Python	+++	+++	++	+++	++	++	+
Torch	Lua, Python (new)	+	+++	++	++	+++	++	
Caffe	C++	+	++		+	+	+	
MXNet	R, Python, Julia, Scala	++	++	+	++	++	+++	
Neon	Python	+	++	+	+	++	+	
CNTK	C++	+	+	+++	+	++	+	

Keras: High Level Wrapper

Keras é uma camada em cima do *Tensorflow*, que facilita alguns procedimentos

Também suporta *Theano* como backend

```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD

N, D, H = 64, 1000, 100

model = Sequential()
model.add(Dense(input_dim=D, output_dim=H))
model.add(Activation('relu'))
model.add(Dense(input_dim=H, output_dim=D))

optimizer = SGD(lr=1e0)
model.compile(loss='mean_squared_error',
              optimizer=optimizer)

x = np.random.randn(N, D)
y = np.random.randn(N, D)
history = model.fit(x, y, nb_epoch=50,
                     batch_size=N, verbose=0)
```

Keras: High Level Wrapper

Define modelo como uma sequência de camadas



```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD

N, D, H = 64, 1000, 100

model = Sequential()
model.add(Dense(input_dim=D, output_dim=H))
model.add(Activation('relu'))
model.add(Dense(input_dim=H, output_dim=D))

optimizer = SGD(lr=1e-0)
model.compile(loss='mean_squared_error',
              optimizer=optimizer)

x = np.random.randn(N, D)
y = np.random.randn(N, D)
history = model.fit(x, y, nb_epoch=50,
                     batch_size=N, verbose=0)
```

Keras: High Level Wrapper

Define optimizador



```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD

N, D, H = 64, 1000, 100

model = Sequential()
model.add(Dense(input_dim=D, output_dim=H))
model.add(Activation('relu'))
model.add(Dense(input_dim=H, output_dim=D))

optimizer = SGD(lr=1e-0)
model.compile(loss='mean_squared_error',
              optimizer=optimizer)

x = np.random.randn(N, D)
y = np.random.randn(N, D)
history = model.fit(x, y, nb_epoch=50,
                     batch_size=N, verbose=0)
```

Keras: High Level Wrapper

```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD

N, D, H = 64, 1000, 100

model = Sequential()
model.add(Dense(input_dim=D, output_dim=H))
model.add(Activation('relu'))
model.add(Dense(input_dim=H, output_dim=D))

optimizer = SGD(lr=1e-0)
model.compile(loss='mean_squared_error',
              optimizer=optimizer)

x = np.random.randn(N, D)
y = np.random.randn(N, D)
history = model.fit(x, y, nb_epoch=50,
                     batch_size=N, verbose=0)
```

Constrói modelo,
especifica função de perda



Keras: High Level Wrapper

Treina o modelo

```
from keras.models import Sequential
from keras.layers.core import Dense, Activation
from keras.optimizers import SGD

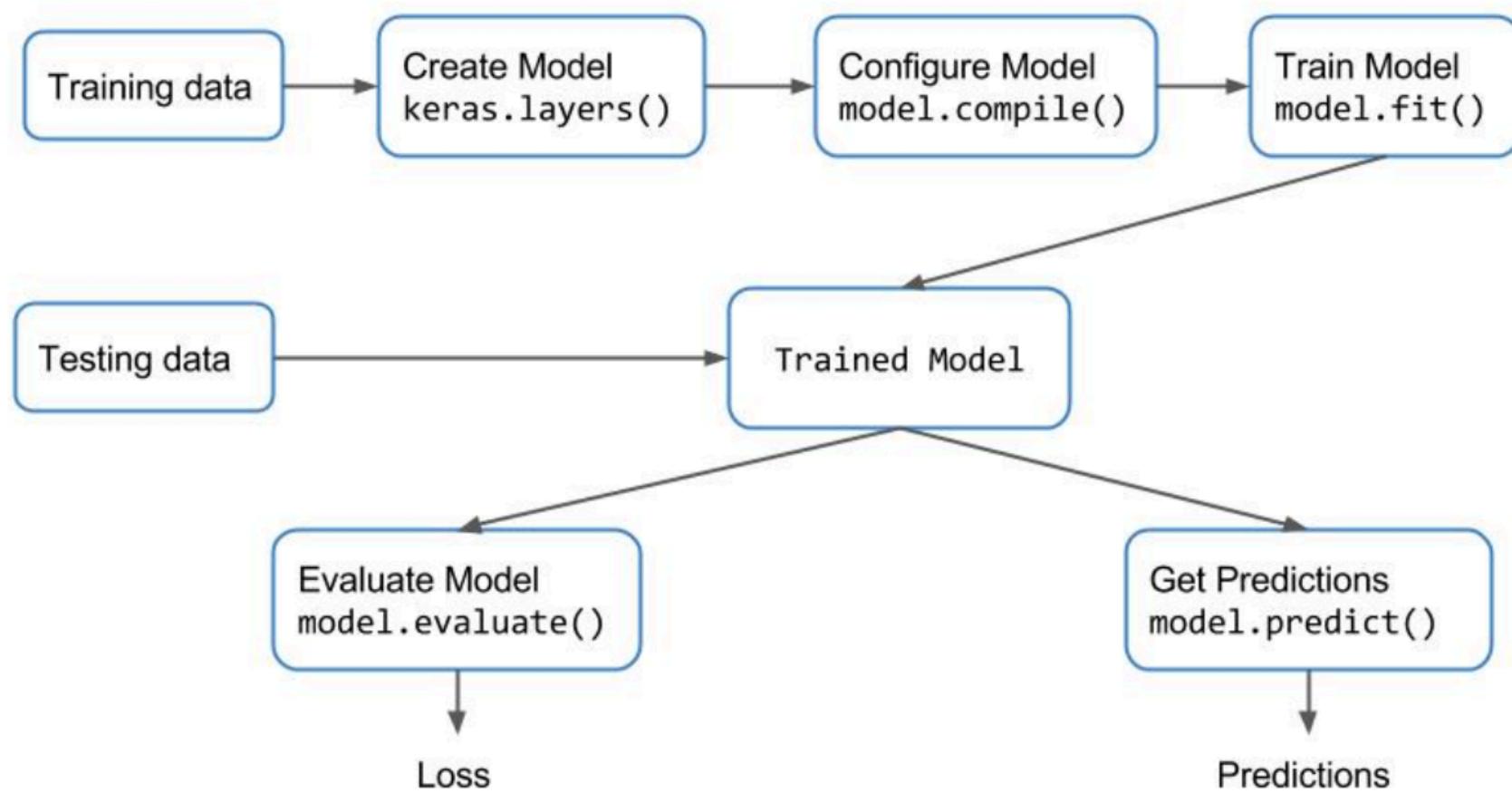
N, D, H = 64, 1000, 100

model = Sequential()
model.add(Dense(input_dim=D, output_dim=H))
model.add(Activation('relu'))
model.add(Dense(input_dim=H, output_dim=D))

optimizer = SGD(lr=1e-0)
model.compile(loss='mean_squared_error',
              optimizer=optimizer)

x = np.random.randn(N, D)
y = np.random.randn(N, D)
history = model.fit(x, y, nb_epoch=50,
                     batch_size=N, verbose=0)
```

Fluxograma Keras



MLP com Keras

▷ código minist_keras.ipynb

One-hot Encoding

ANIMAL	VALOR
	?
	?
	?
	?
	?



ANIMAL	PATO	CASTOR	MORSA?
	1	0	0
	0	1	0
	1	0	0
	0	0	1
	0	1	0

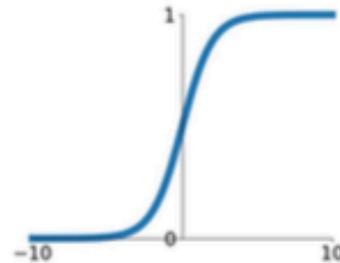
TODO

▷ Cifar10_keras.ipynb

Funções de Ativação

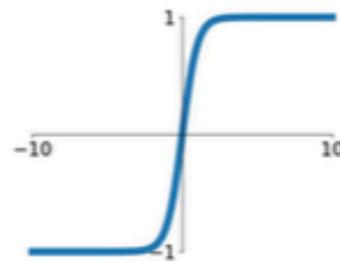
Sigmoid

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



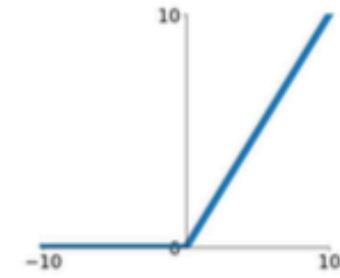
tanh

$$\tanh(x)$$



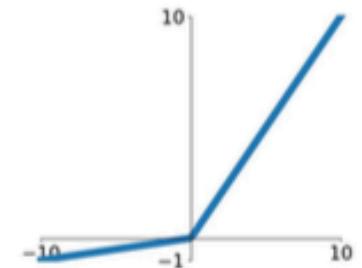
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

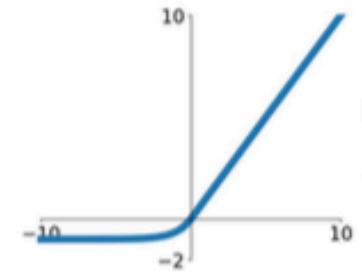


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Histórico das funções de ativação

Name	Formula	Year
none	$y = x$	-
sigmoid	$y = \frac{1}{1+e^{-x}}$	1986
tanh	$y = \frac{e^{2x}-1}{e^{2x}+1}$	1986
ReLU	$y = \max(x, 0)$	2010
(centered) SoftPlus	$y = \ln(e^x + 1) - \ln 2$	2011
LReLU	$y = \max(x, \alpha x), \alpha \approx 0.01$	2011
maxout	$y = \max(W_1x + b_1, W_2x + b_2)$	2013
APL	$y = \max(x, 0) + \sum_{s=1}^S a_i^s \max(0, -x + b_i^s)$	2014
VLReLU	$y = \max(x, \alpha x), \alpha \in 0.1, 0.5$	2014
RReLU	$y = \max(x, \alpha x), \alpha = \text{random}(0.1, 0.5)$	2015
PReLU	$y = \max(x, \alpha x), \alpha \text{ is learnable}$	2015
ELU	$y = x, \text{ if } x \geq 0, \text{ else } \alpha(e^x - 1)$	2015

