

# REDES NEURAIS

**Thaís de Almeida Ratis Ramos**

**Baseado nas aulas de Thaís Gaudencio e Yuri Malheiros – UFPB**

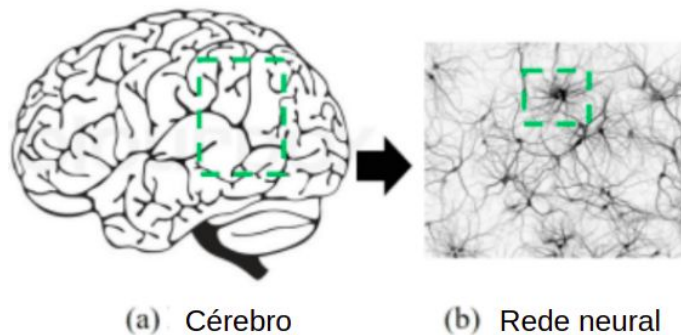
# INSPIRAÇÃO BIOLÓGICA

- O cérebro é o principal órgão associado à inteligência e aprendizagem
- O cérebro é composto por uma rede complexa de aproximadamente 100 bilhões de neurônios interconectados
- Existem mais de 500 trilhões de conexões entre neurônios no cérebro humano
- Mesmo as maiores redes neurais artificiais de hoje não chegam perto do cérebro humano

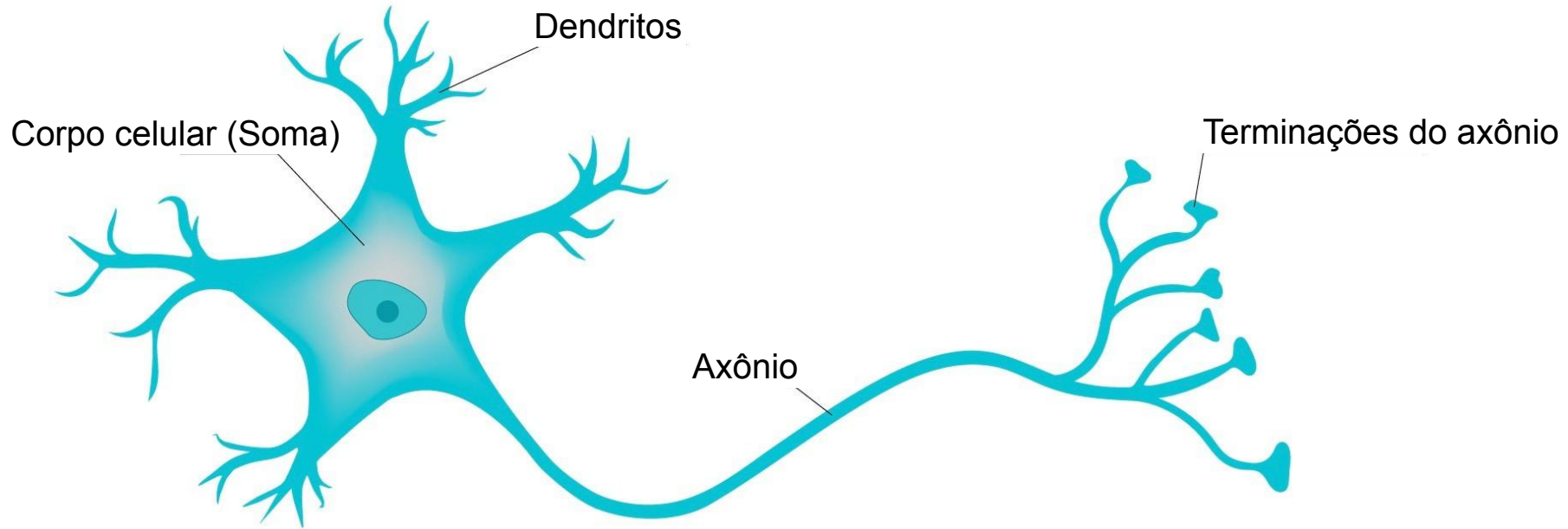
# INSPIRAÇÃO BIOLÓGICA

A unidade básica de uma rede neural é o neurônio artificial

Neurônios artificiais são modelados como neurônios biológicos do cérebro, nos quais são estimulados por entradas



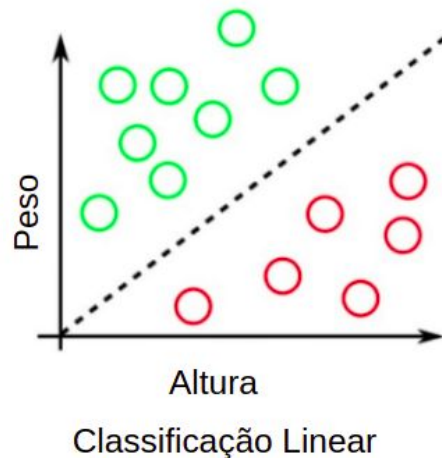
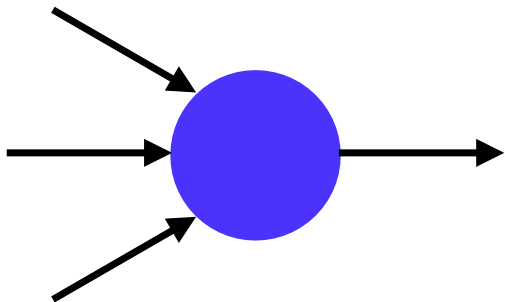
# NEURÔNIO



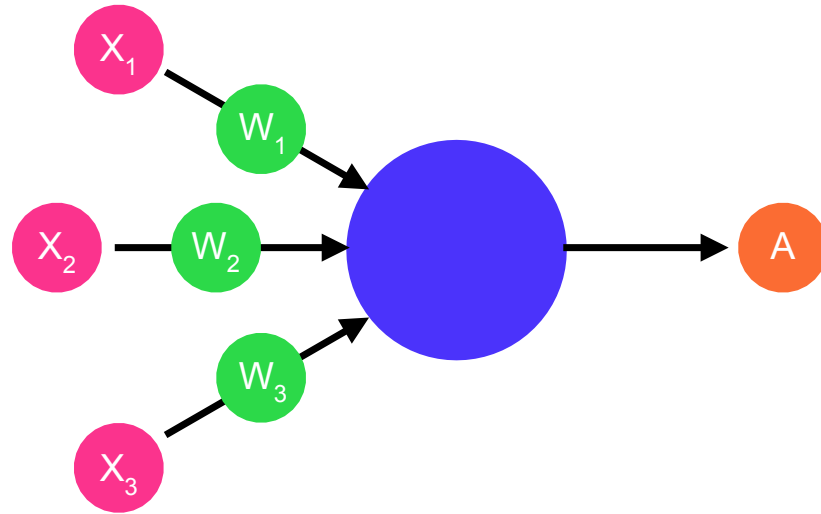
# PERCEPTRON

Em 1943, McCullock e Pitts propuseram um modelo de neurônio artificial: **perceptron**

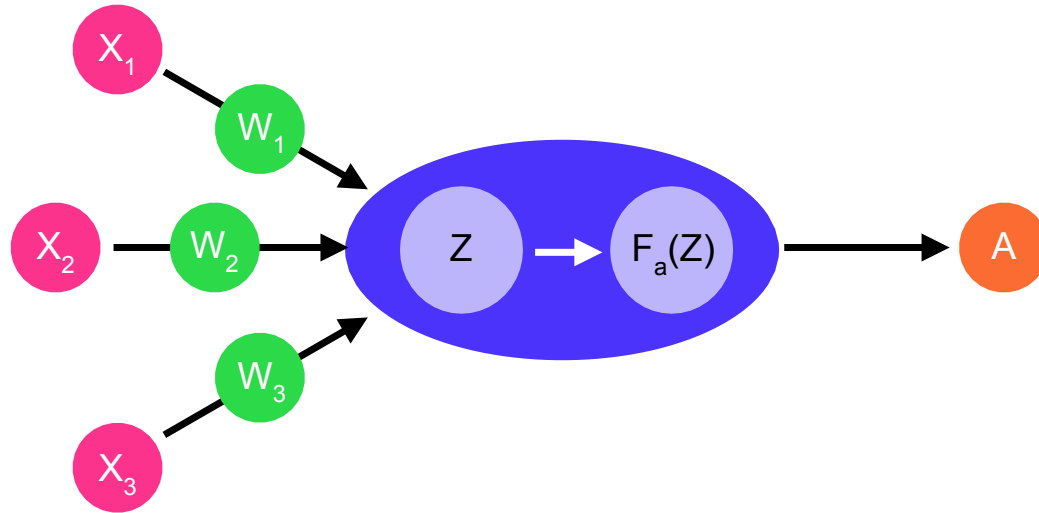
Modelo linear utilizado para classificação binária



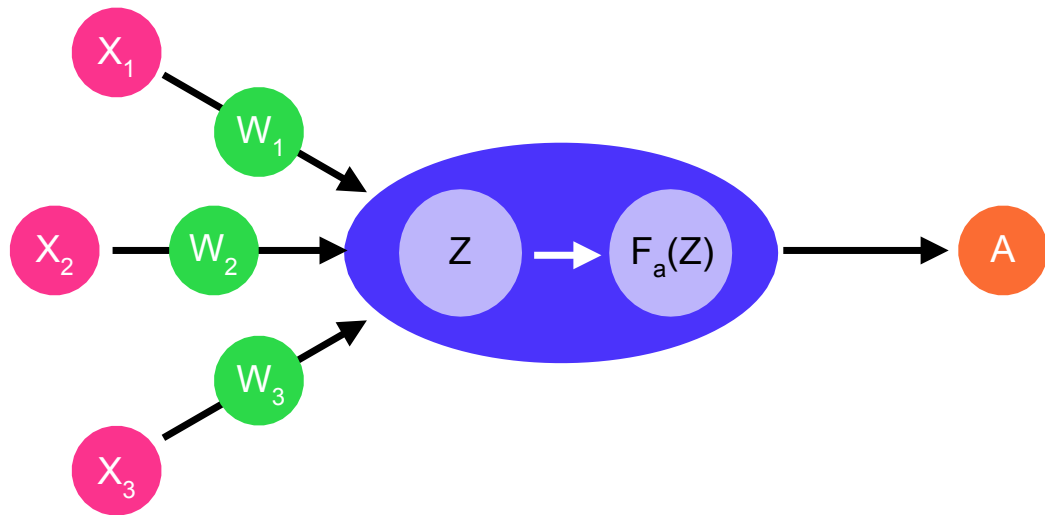
# PERCEPTRON



# PERCEPTRON



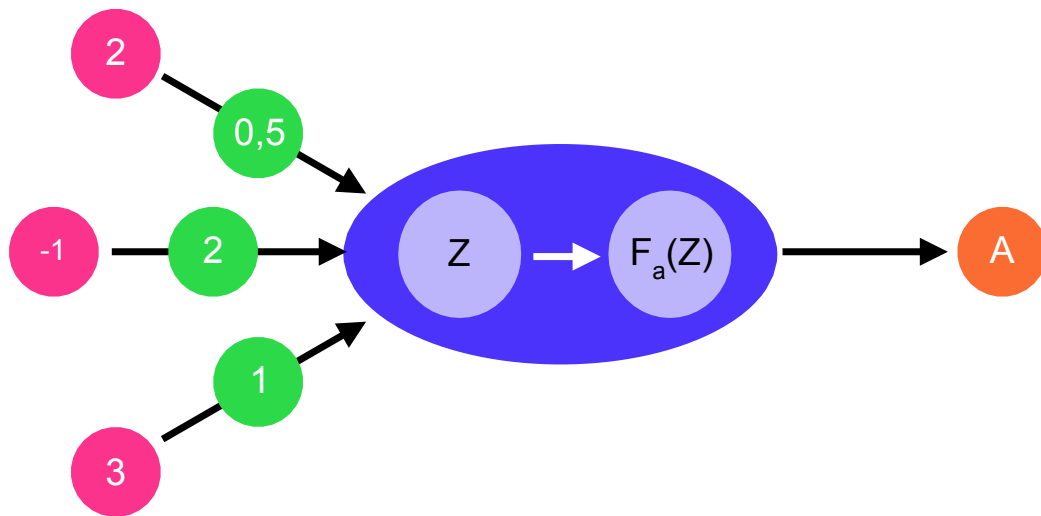
# PERCEPTRON



$$Z = w_1 \times x_1 + w_2 \times x_2 + w_3 \times x_3$$

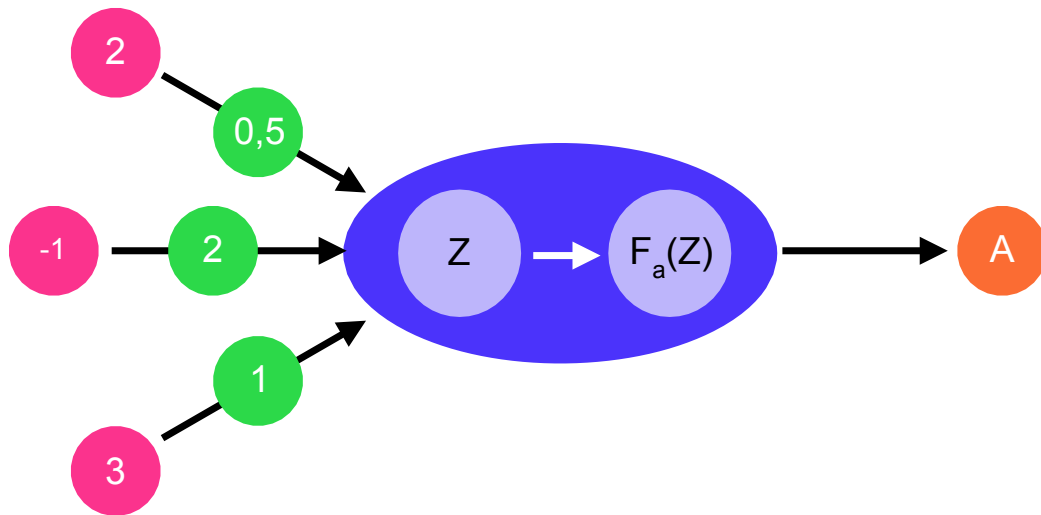


# PERCEPTRON



$$Z = 0,5 \times 2 + 2 \times -1 + 1 \times 3 = 2$$

# PERCEPTRON



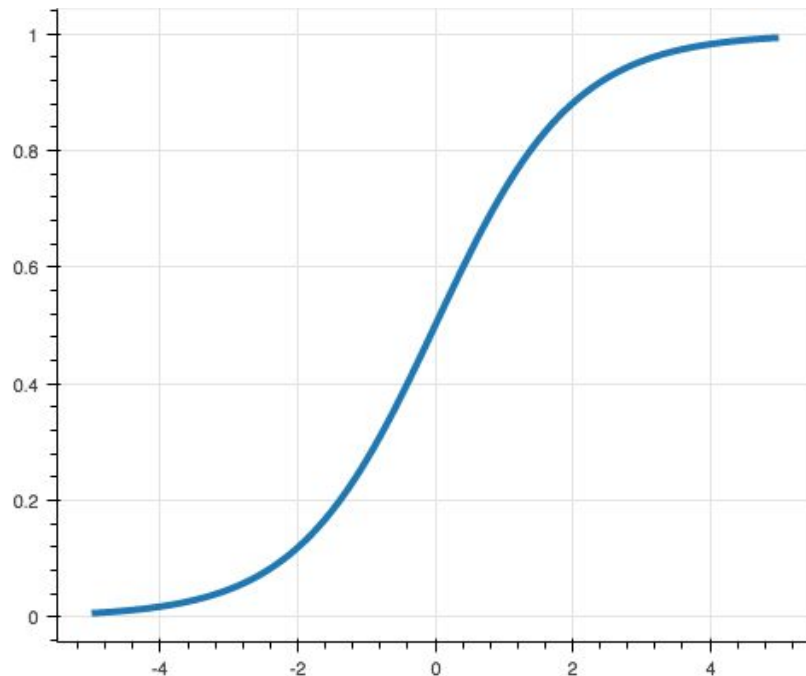
$$A = F_a(2)$$

# FUNÇÕES DE ATIVAÇÃO

Várias funções podem ser usadas como função de ativação ( $F_a$ )

Função Sigmóide:

$$f_a(z) = \frac{1}{1 + e^{-z}}$$

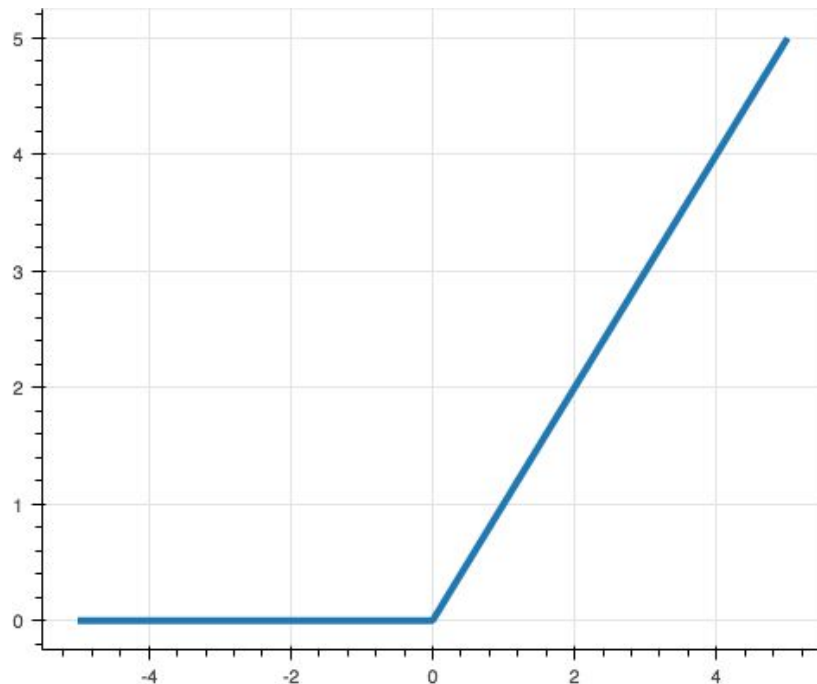


# FUNÇÕES DE ATIVAÇÃO

Várias funções podem ser usadas como função de ativação ( $F_a$ )

Função ReLU:

$$f_a(z) = \max(0, z)$$



# PERCEPTRON

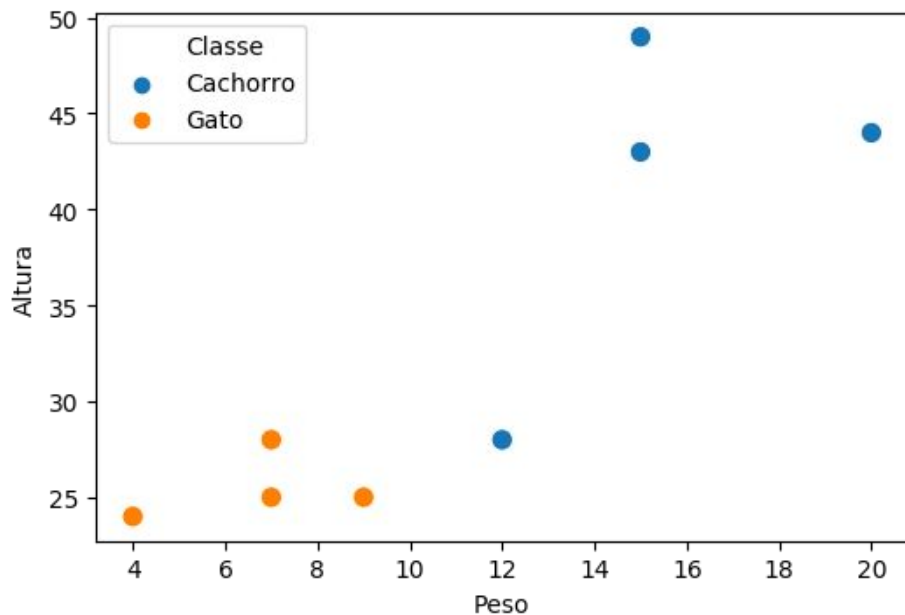
Podemos usar o neurônio artificial para diferentes problemas

Ele nos fornece uma saída de acordo com as entradas

O resultado do processamento das entradas para fornecer uma saída é determinado pela função de ativação e pelos pesos

# EXEMPLO

## Gatos e cachorros:



	Peso	Altura	Classe
0	20	44	Cachorro
1	15	43	Cachorro
2	12	28	Cachorro
3	15	49	Cachorro
4	7	28	Gato
5	7	25	Gato
6	9	25	Gato
7	4	24	Gato

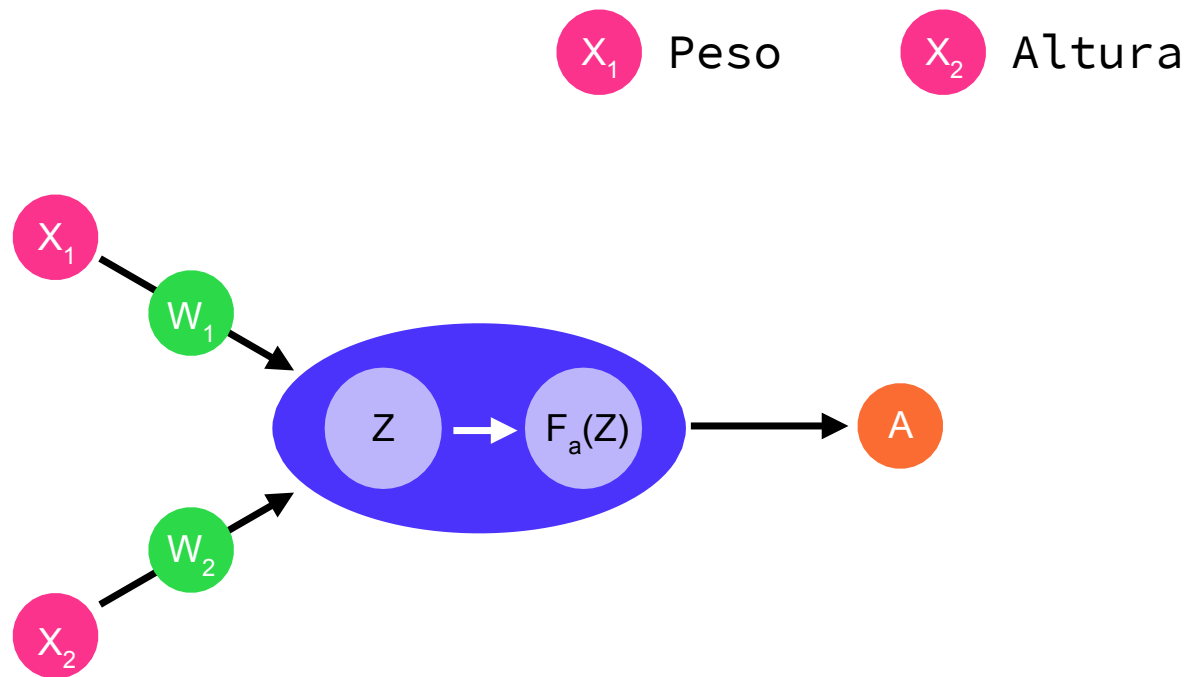
# EXEMPLO

Dada a altura e o peso, vamos fazer o neurônio ter:

Saída **0** se for um gato

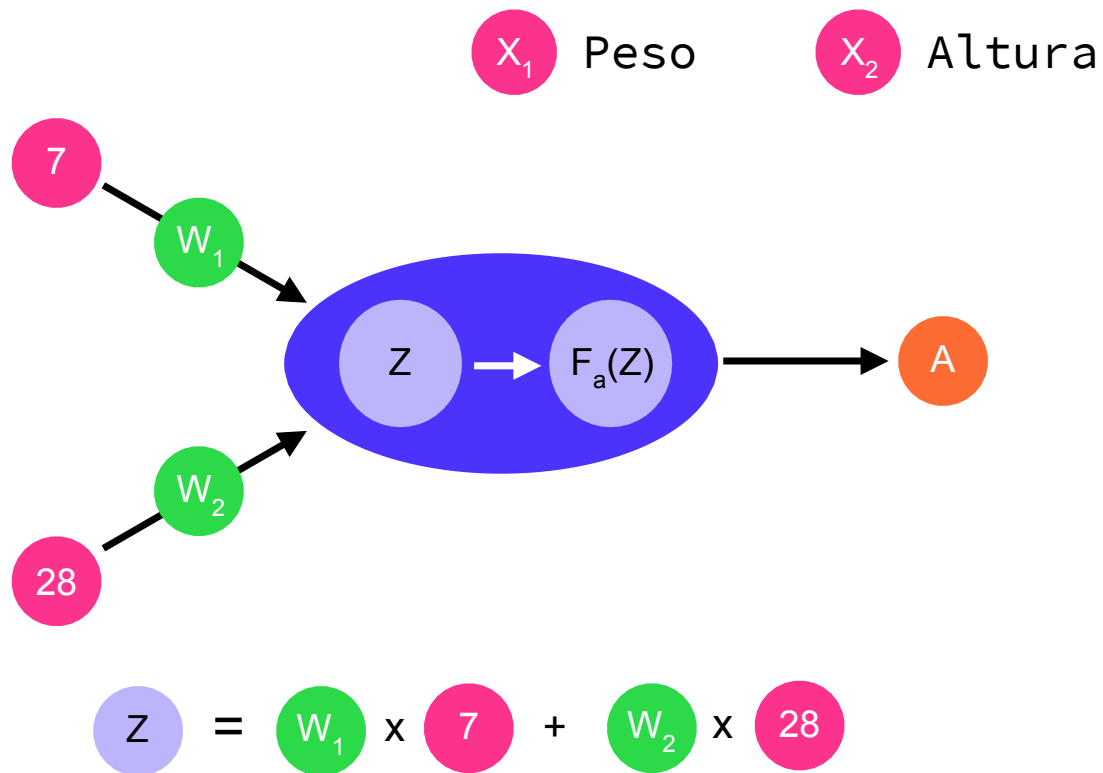
Saída **1** se for um cachorro

# EXEMPLO





# EXEMPLO



# TREINANDO UMA REDE NEURAL

Quais são os valores de  $w$ ?

Podemos tentar encontrá-los manualmente, porém este é um processo complexo e muitas vezes inviável

# TREINANDO UMA REDE NEURAL

Como não sabemos os pesos iniciais, vamos iniciá-los aleatoriamente

$$w_1 = 0,25$$

$$w_2 = 0,66$$

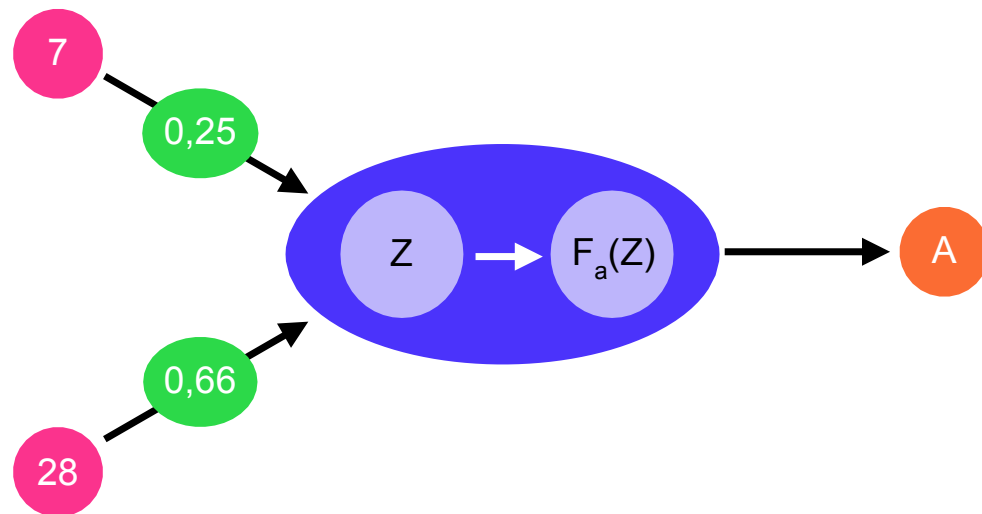
Em seguida, vamos calcular a saída para

$$x_1 = 7$$

$$x_2 = 28$$

# TREINANDO UMA REDE NEURAL

Saída 0: gato  
Saída 1: cachorro



$$Z = 0,25 \times 7 + 0,66 \times 28 = 20,23$$

$$A = F_a(20,23) = 0,99$$

# TREINANDO UMA REDE NEURAL

Esperado	Obtido	Ajuste	Esperado - Obtido
0	0	-	0
0	1	Diminuir W	-1
1	0	Aumentar W	1
1	1	-	0

# TREINANDO UMA REDE NEURAL

Para isso, eu preciso saber o erro, ou seja a diferença entre o valor esperado e o valor obtido

Quanto maior o erro, maior precisa ser o ajuste dos valores de  $w$

Então, vamos atualizar  $w_1$  da seguinte forma:

$$w_1 \leftarrow w_1 + \eta (\textit{esperado} - \textit{obtido}) x_1$$

 Taxa de aprendizagem: controla o tamanho da atualização

# TREINANDO UMA REDE NEURAL

Cada ajuste vai melhorando um pouco a rede neural

Vamos ajustar os pesos para cada um dos exemplos dos dados de treinamento

$$w_1 \leftarrow w_1 + \eta(\textit{esperado} - \textit{obtido})x_1$$

$$w_2 \leftarrow w_2 + \eta(\textit{esperado} - \textit{obtido})x_2$$

# TREINANDO UMA REDE NEURAL

Em muitos casos, fazer uma atualização para cada exemplo pode não ser suficiente

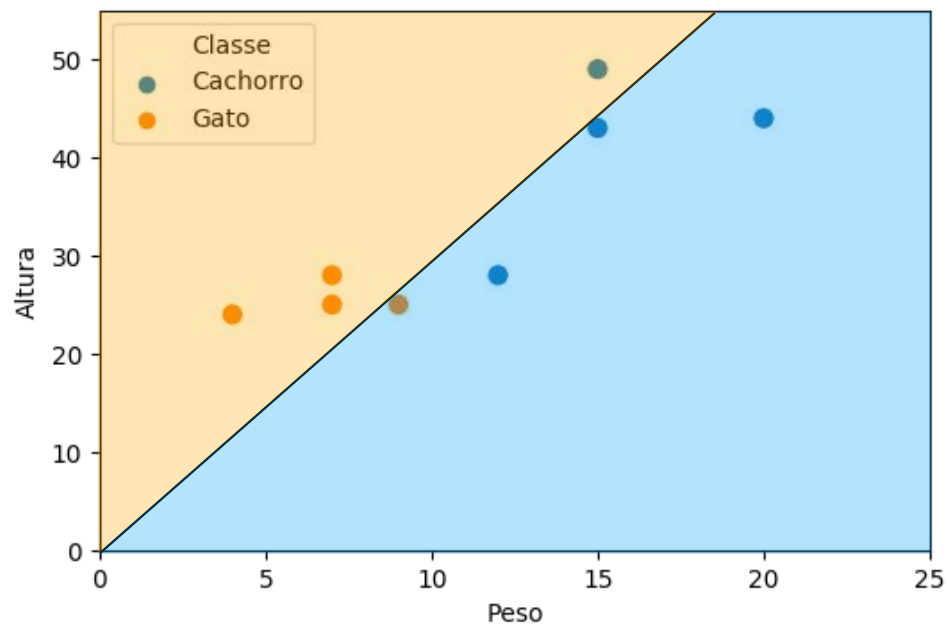
Por isso, costumamos atualizar para cada exemplo e em seguida refazer o processo várias vezes

Esse número de vezes é chamado de **épocas**



# EXEMPLO

Temos como resultado:



# TREINANDO UMA REDE NEURAL

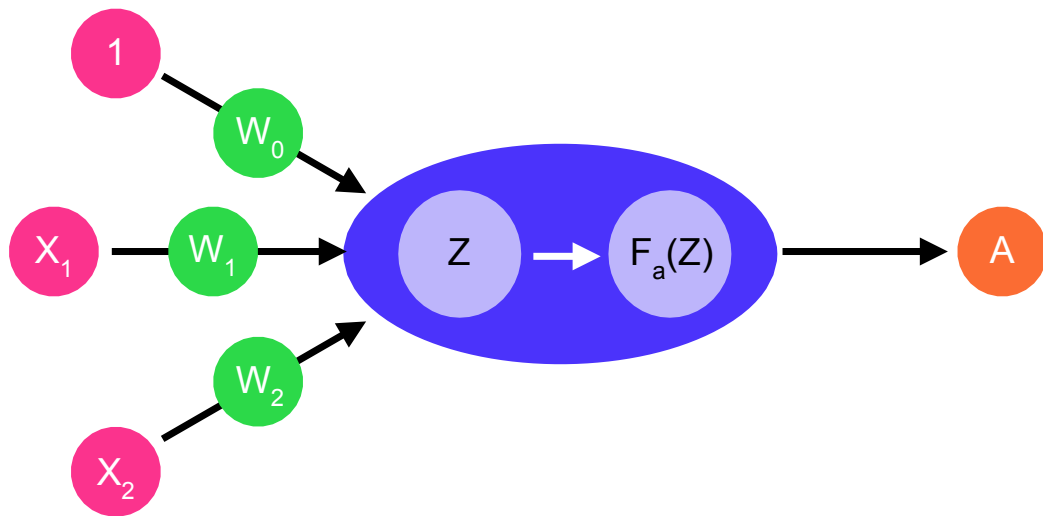
Usando o neurônio com a estrutura que fizemos, sempre teremos uma fronteira linear partindo da origem

Para permitir que a fronteira possa se deslocar no eixo X, vamos adicionar uma entrada extra no neurônio chamada de **viés**

O viés geralmente recebe 1 como entrada

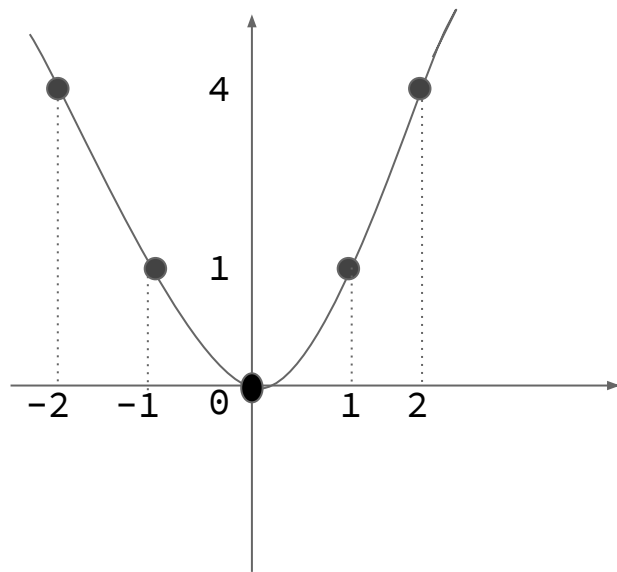
O peso dessa entrada segue o comportamento das outras

# TREINANDO UMA REDE NEURAL



# VIÉS (BIAS)

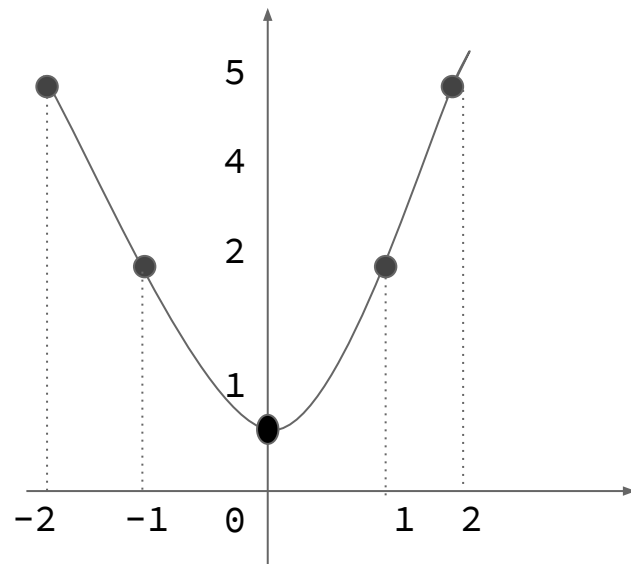
$$b = 0$$



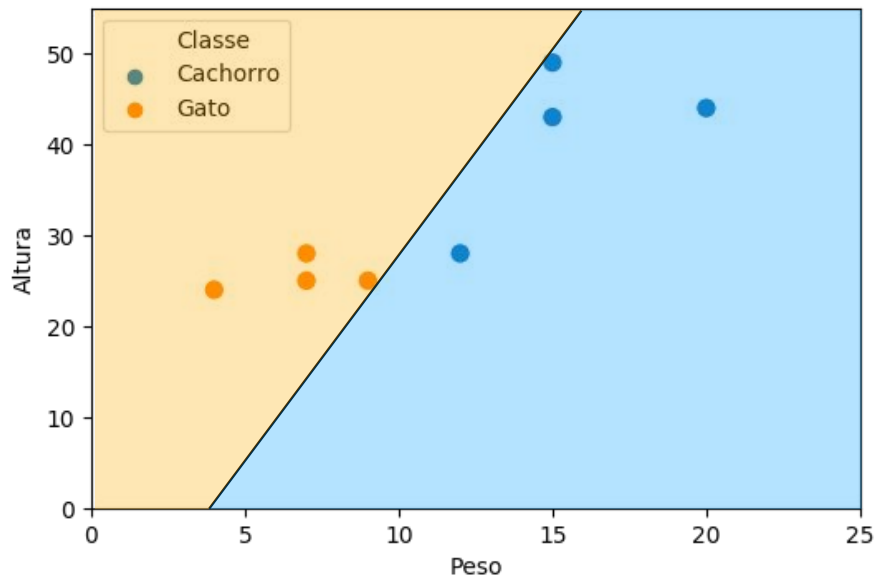
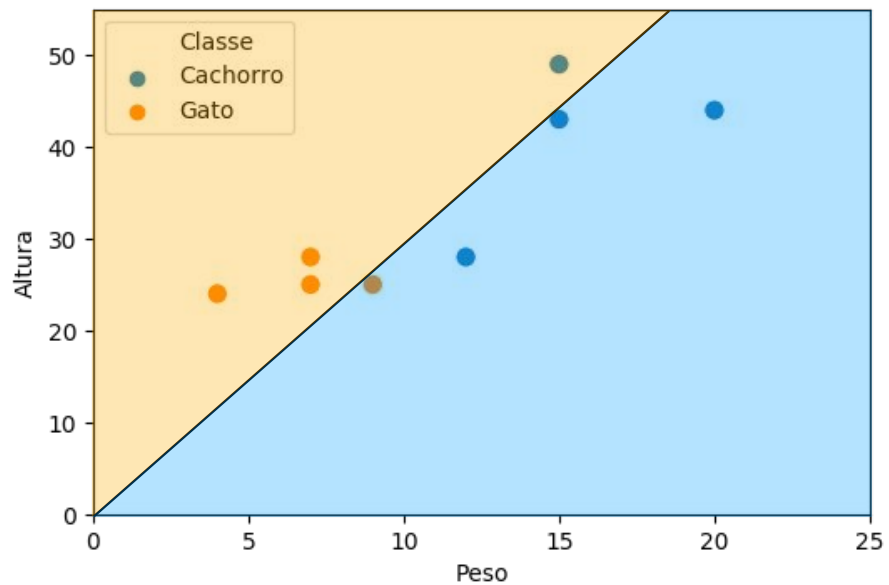
$$f(x) = x^2 + b$$

x	f(x) b=0	f(x) b=1
-2	4	5
-1	1	2
0	0	1
1	1	2
2	4	5

$$b = 1$$

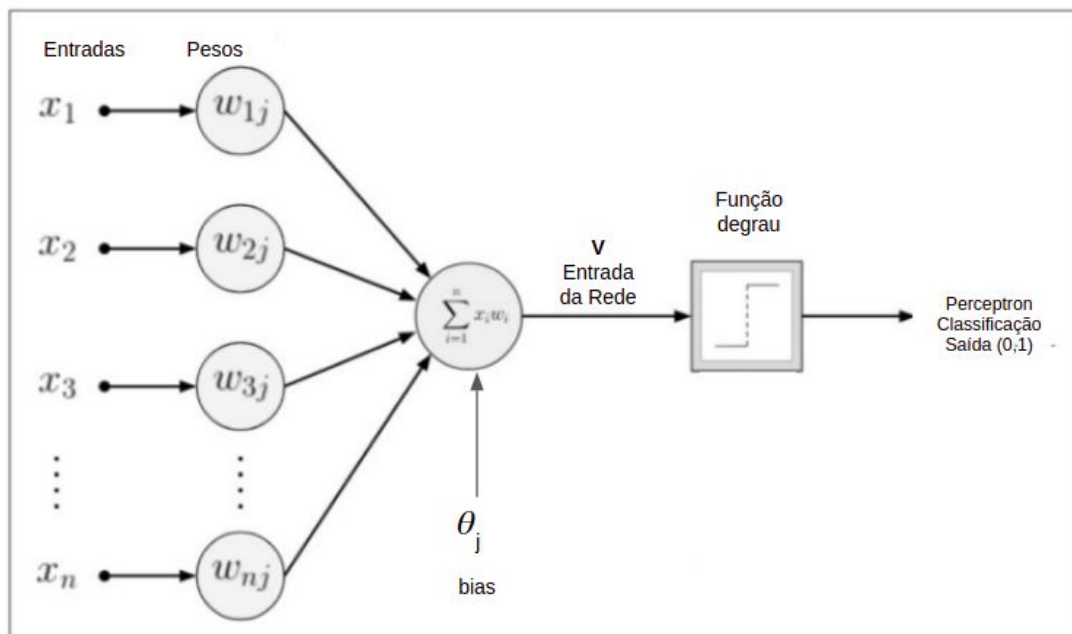


# VIÉS (BIAS)



# ESTRUTURA DO PERCEPTRON

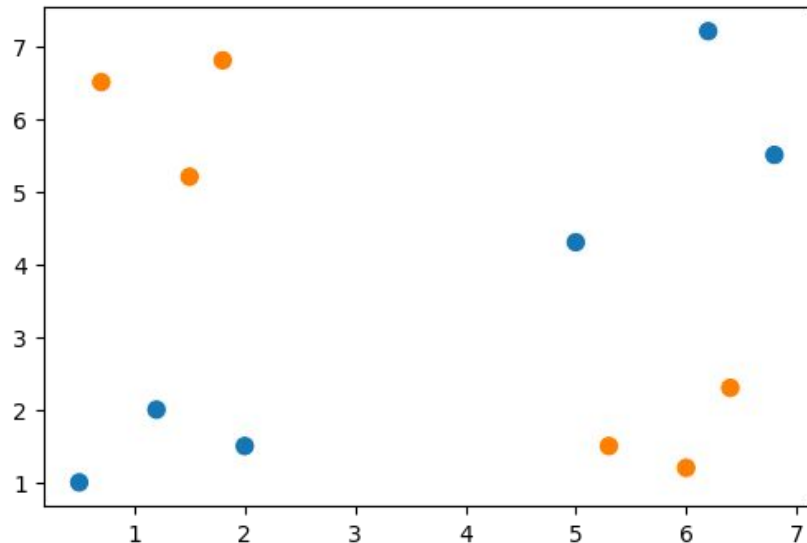
Função de ativação de Heaviside: 
$$f(v) = \begin{cases} 0 & v < 0 \\ 1 & v \geq 0 \end{cases}$$



# TREINANDO UMA REDE NEURAL

Como vimos, o perceptron possui a limitação de criar apenas fronteiras lineares

Ele não consegue classificar bem um conjunto de dados dessa forma:

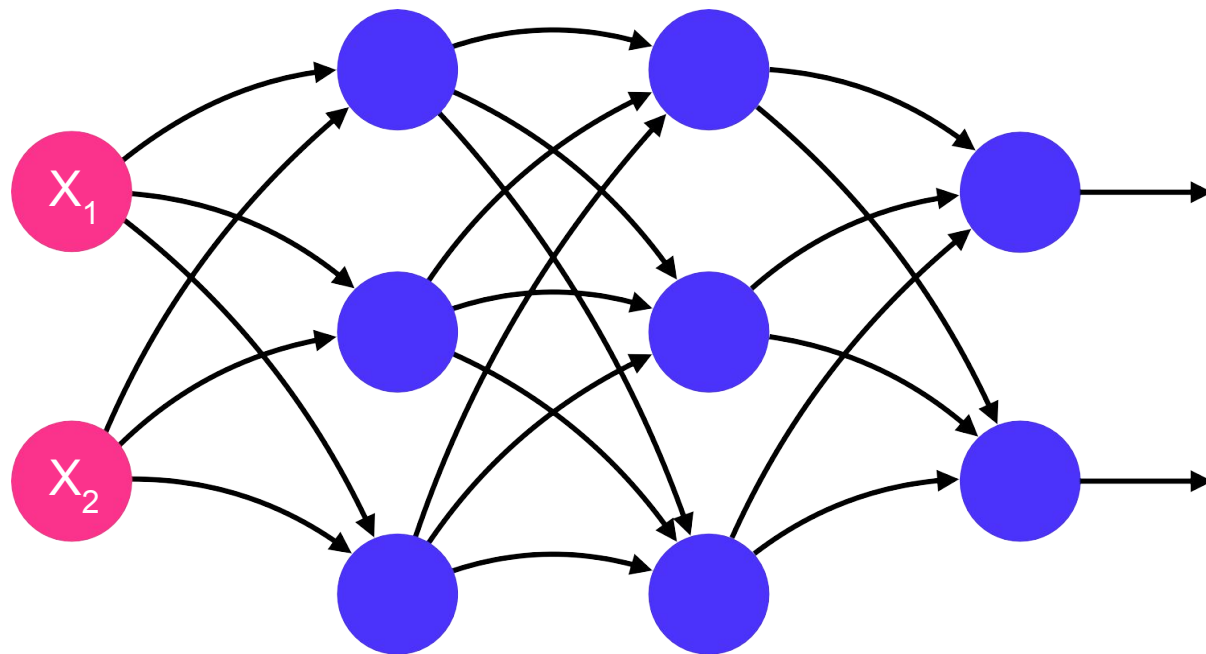


# REDE NEURAL

Camada de  
entrada

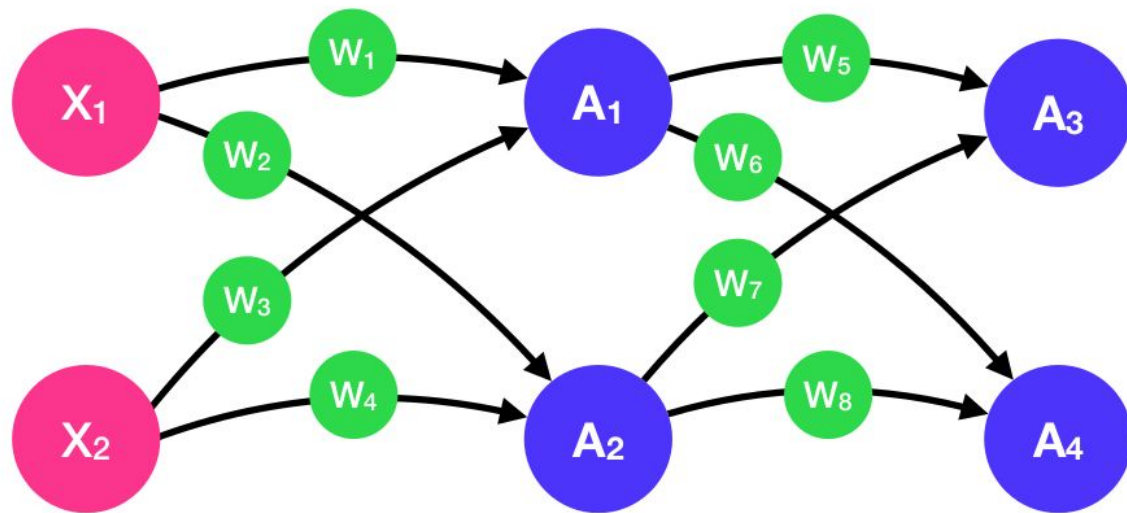
Camada  
intermediária

Camada de  
saída

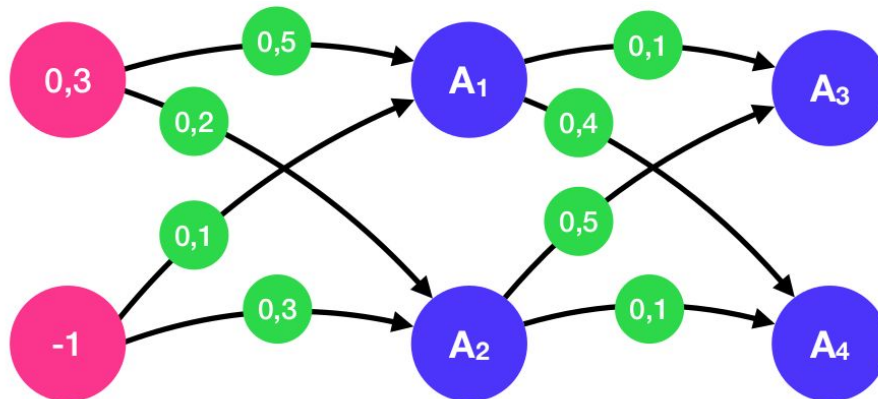




# REDE NEURAL



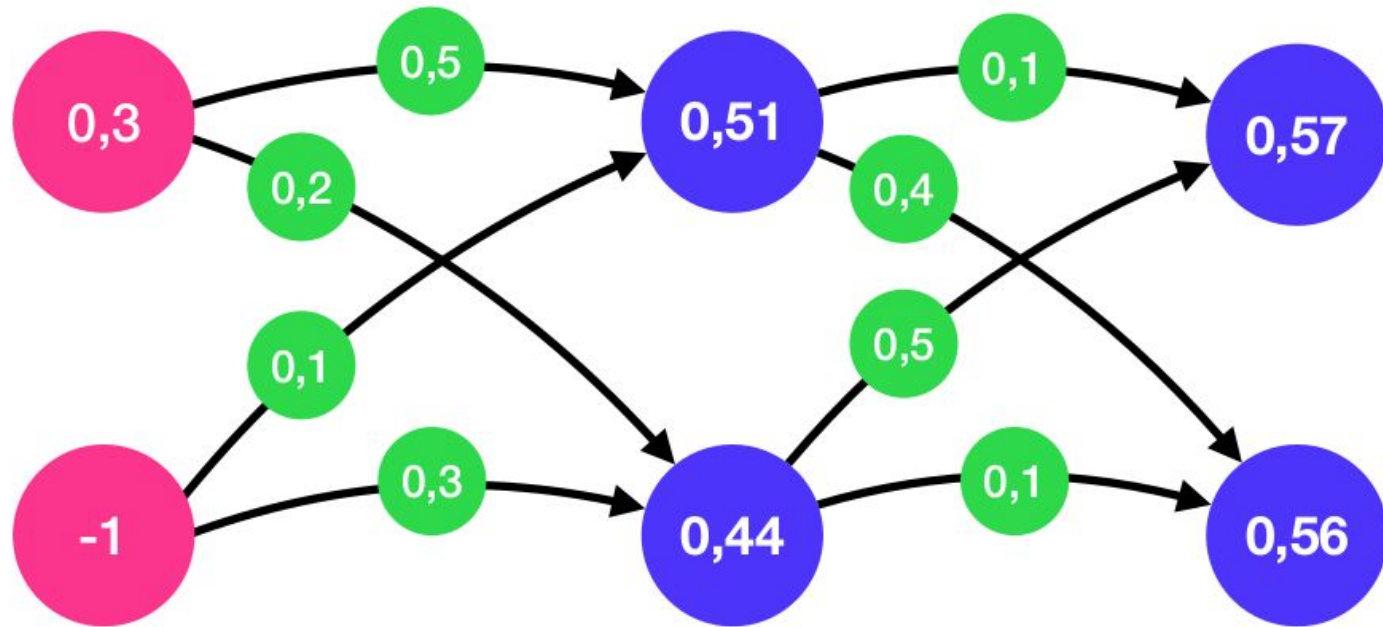
# REDE NEURAL



$$Z_1 = 0,5 \times 0,3 + 0,1 \times -1 = 0,049$$

$$A_1 = F_a(0,049) = 0,51$$

# REDE NEURAL



# TREINANDO UMA REDE NEURAL

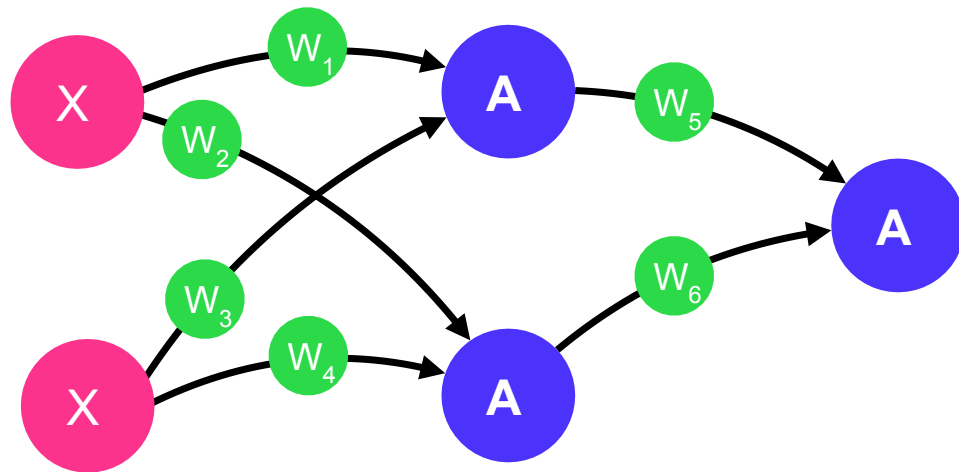
O algoritmo é semelhante ao que fizemos para apenas um neurônio

Iniciamos os pesos de forma aleatória

Entramos com os dados de treinamento e comparamos o resultado da rede com o resultado esperado

Ajustamos os pesos de acordo com o erro

# TREINANDO UMA REDE NEURAL



O valor de  $A_3$  é a saída da rede. Podemos compará-lo diretamente com o resultado dos dados de treinamento

E para  $A_1$  e  $A_2$ ? Quais são os valores esperados? Não temos esses dados

# TREINANDO UMA REDE NEURAL

Para treinar uma rede com multilaminadas usamos o algoritmo backpropagation

Nele, vamos definir uma função que mede o erro da saída da rede

Também conseguimos calcular através do gradiente da função de erro, em qual sentido devemos ajustar os pesos para que o erro diminua

# CORREÇÃO DE ERROS

Calculamos o erro para um exemplo através da função:

$$e_k = d_k - y_k$$

Onde:

e - Sinal de erro

d - Saída desejada apresentada durante o treinamento

y - Saída real da rede

# APRENDIZADO POR CORREÇÃO DE ERROS

- O processo de aprendizado por correção de erros utiliza algoritmos para caminhar sobre a curva de erros, com o intuito de alcançar o menor valor de erro possível, o mínimo global
- Muitas vezes o algoritmo não alcança este mínimo global, atingindo o que chamamos de mínimo local. Caso este erro alcançado seja desfavorável, é necessário recomeçar o processo de aprendizagem

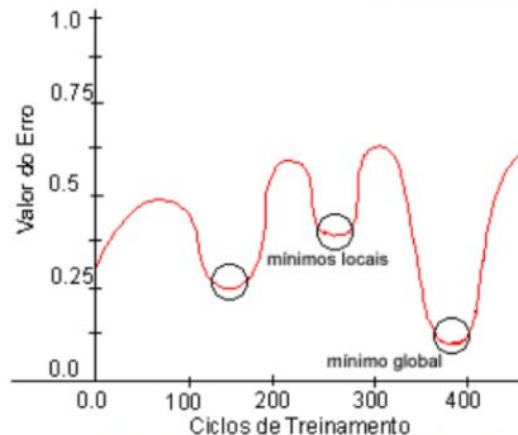


Gráfico de uma possível superfície de erro mostrando os mínimos locais e o mínimo global



# APRENDIZADO POR CORREÇÃO DE ERROS

Para correção do erro, os pesos devem ser ajustados de forma a aproximar a saída real à desejada

$$\Delta w_i(n) = \eta e(n) x_i(n)$$

Onde:

$\Delta w_i(n)$  - Valor de ajuste a ser acrescentado ao peso  $w_i$ ;

$\eta$  - Taxa de aprendizagem (constante positiva);

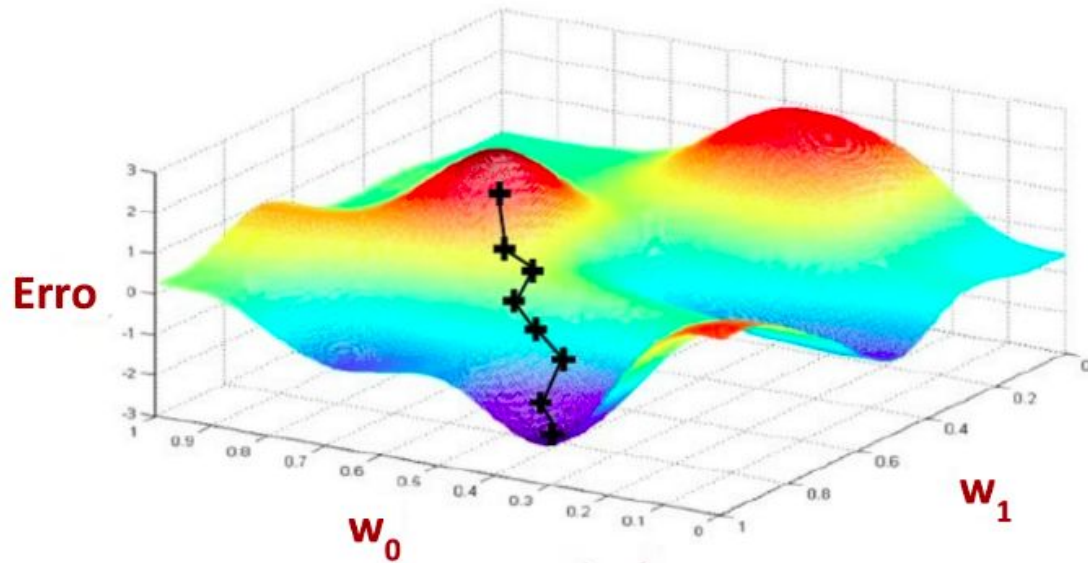
$e(n)$  - Valor do erro;

$x_i$  - Valor da entrada

ENCONTRAR O VETOR (MATRIZ) DE PESOS SINÁPTICOS ( $w^*$ )  
QUE MINIMIZEM O ERRO (E) ENTRE A SAÍDA DA REDE  
NEURAL E A SAÍDA DESEJADA

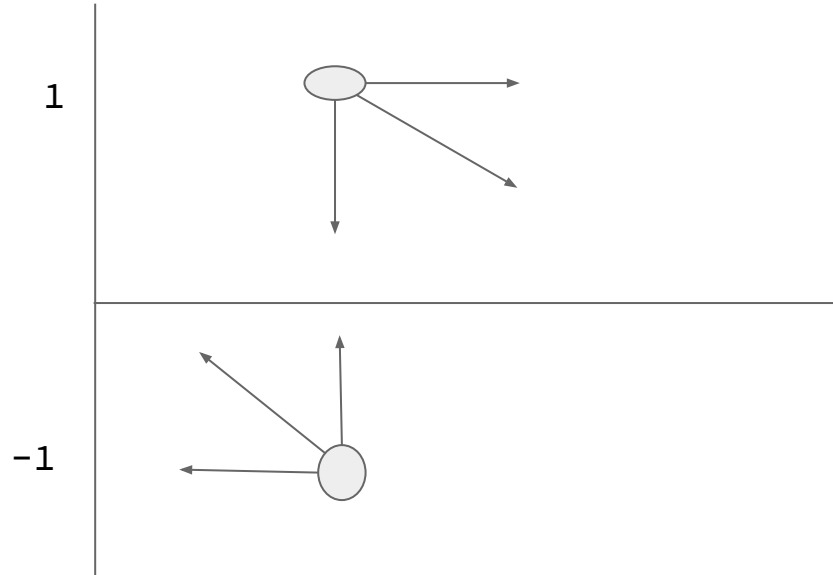
# MÉTODO DO GRADIENTE

Vetor que indica o sentido e a direção onde a função crescerá mais rápido



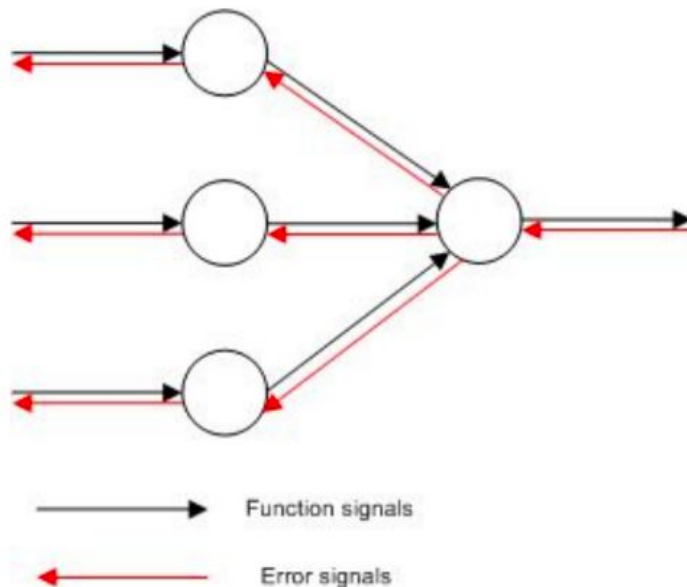
# MÉTODO DO GRADIENTE

Vetor que indica o sentido e a direção onde a função crescerá mais rápido



# BACKPROPAGATION

Sinal de erro se origina em um neurônio de saída e se propaga para trás (camada por camada) através da rede



# CORREÇÃO DOS PESOS - REGRA DELTA

$$\Delta w_{ji}(n) = \eta * \delta j(n) * y_i(n)$$

Correção  
de peso

Taxa de  
aprendizagem

Gradiente  
local

Sinal de  
entrada do  
neurônio  $j$

# GRADIENTE LOCAL

Camada de saída

$$\delta_j(n) = e_j * \varphi'(v_j(n))$$

Gradiente  
Local

Sinal de erro  
na saída de  $j$

Derivada da  
função de  
ativação no  
neurônio  $j$

Camada oculta

$$\delta_j(n) = \varphi'(v_j(n)) \sum_k \delta_k(n) w_{kj}(n)$$

Gradiente  
Local

Derivada  
associada ao  
neurônio  $j$

Soma ponderada  
dos gradientes  
locais da camada  
seguinte  $k$

# BACKPROPAGATION

```
function backpropagation-algorithm
  (network, training_examples, learning_rate)
  network <- initialize_weights(randomly)
  start loop
    for each example in training-examples do
      // Computação para frente (Propagação)
      network_out = rna_output(network, example)

      // Computacao para trás (Retropropagação)
      example_err = actual_out - network_out
      local_gradient = calc_gradient(example_err)
```



# BACKPROPAGATION

Atualiza os pesos dos neurônios da camada de saída  $j$  (regra delta)

**for each** previous-layer **in** network **do**

    Compute o erro em cada nó

    Atualiza os pesos dos neurônios da camada  
    (regra delta)

**end for**

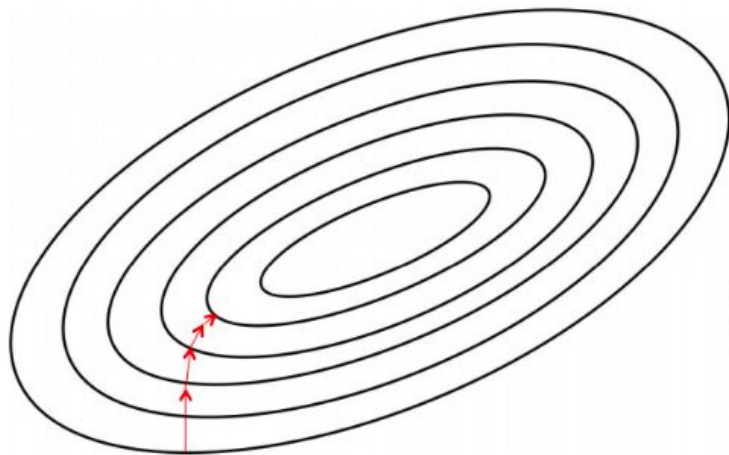
**end for**

**end loop** quando rede convergir

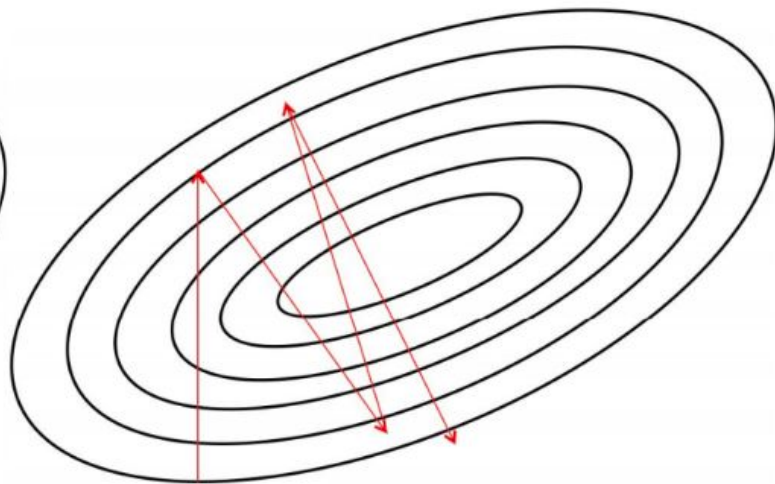
**return** network

# TAXA DE APRENDIZAGEM

**$\eta$  muito pequeno**  
Aprendizagem e  
convergência lenta



**$\eta$  muito grande**  
Rede instável



# REGRA DELTA COM TERMO DO MOMENTUM

Permite aumentar a taxa de aprendizado, evitando risco de instabilidade

Leva em consideração a variação dos pesos da iteração anterior:

- Variação de pesos com o mesmo sinal: Acelera a descida
- Variação de pesos com sinais opostos: Reduz o ajuste nos pesos (estabiliza o algoritmo)

Pode evitar que o processo de aprendizagem termine em um mínimo local raso

$$\Delta w_{ji}(n) = \alpha \Delta w_{ji}(n-1) + \eta \delta_j(n) y_i(n)$$

The equation is annotated with red boxes and arrows. A red box encloses the term  $\alpha \Delta w_{ji}(n-1)$ , with an arrow pointing to it from the label 'Constante do momento'. Another red box encloses the term  $\eta \delta_j(n) y_i(n)$ , with an arrow pointing to it from the label 'Regra delta'. A third red box encloses the entire right-hand side of the equation, with an arrow pointing to it from the label 'Variação de pesos da última iteração'.

$$0 \leq \alpha < 1$$

MODOS DE TREINAMENTO: ESTOCÁSTICO, POR LOTE E  
MINI-LOTE

# MODOS DE TREINAMENTO

Baseado na forma de atualização dos pesos, temos 3 modos de treinamento:

- Estocástico - Ajuste de pesos é realizado após a apresentação de cada exemplo
- Por lote (batch) - Ajuste de pesos é realizado após a apresentação de todos os exemplos à rede (fim da época)
- Mini-lote (mini-batch) - Ajuste de pesos é realizado após a apresentação de um subconjunto de exemplos

# ESTOCÁSTICO

## **Vantagens:**

- Ajuda a escapar de mínimos locais
- Mais simples de implementar
- Pode tirar vantagens de dados (exemplos) redundantes

## **Desvantagens:**

- Estimar o erro baseado em um único exemplo não é uma boa aproximação do erro real
- Treinamento muito lento
- Mais difícil provar teoricamente que o algoritmo converge

# POR LOTE (BATCH)

## **Vantagens:**

- Estimativa precisa do gradiente
- Convergência mais rápida sob condições simples
- Mais fácil de paralelizar

## **Desvantagens:**

- Pode ficar preso em mínimos locais

# POR MINI-LOTE (MINI-BATCH)

Bom balanço entre o modo estocástico e o modo por lote:

- Convergência mais rápida (modo por lote)
- Evita mínimos locais



# CRITÉRIOS DE PARADA

Existem alguns critérios razoáveis para a convergência:

- Vetor gradiente alcançar um limiar suficiente pequeno
- Taxa de variação do erro muito pequena entre as épocas (ex: menor que 1%)
- Rede apresenta um bom desempenho de generalização, ou seja, funciona bem com um outro conjunto de exemplos (conjunto de validação)