

Hands-on Deep Learning - Aula 2

Introdução às Redes Neurais

Camila Laranjeira¹², Hugo Oliveira¹², Pedro H. Barros¹³

¹Programa de Pós-Graduação em Ciência da Computação (PPGCC)
Universidade Federal de Minas Gerais

²Interest Group in Pattern Recognition and Earth Observation (PATREO)
Universidade Federal de Minas Gerais

³Wireless Informational Sensing Embedded systems Modeling Algorithms and Protocols (WISEMAP)
Universidade Federal de Minas Gerais

08 de Fevereiro, 2020





Agenda

1 Origem

2 Fundamentos

- Classificação Linear
- Perceptron
- Otimização
- Redes Neurais Artificiais
- Treinamento

3 Framework

- Comparação Entre Frameworks
- Pytorch

Agenda



1 Origem

2 Fundamentos

- Classificação Linear
- Perceptron
- Otimização
- Redes Neurais Artificiais
- Treinamento

3 Framework

- Comparação Entre Frameworks
- Pytorch



O Cérebro

- Redes Neurais foram originalmente inspiradas no cérebro humano
- Sua unidade básica é o neurônio, composto por 3 partes principais:
 - Dendritos: Recebem sinais de outros neurônios conectados a ele
 - Corpo: Responsável por “processar” a informação
 - Axônio: Transmite o sinal processado para outros neurônios. A ativação do axônio depende da força sináptica do sinal.

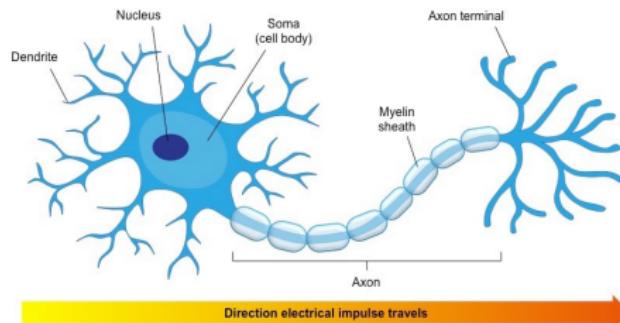


Figura: Representação de um neurônio [1]

O Cérebro

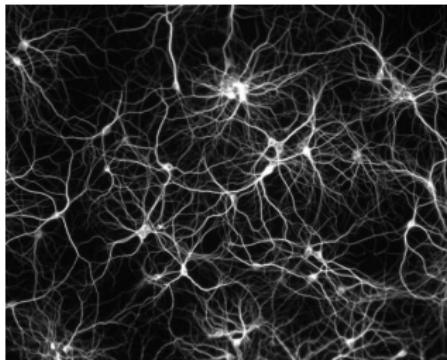


Figura: Representação da malha neural [2]

- Os neurônios se conectam uns aos outros, criando uma grande malha.
- Responsável pelo processamento de todas as funções do corpo.

Representação Artificial do Cérebro



- Redes neurais artificiais são sistemas compostos de neurônios artificiais **inspirados** no neurônio biológico.

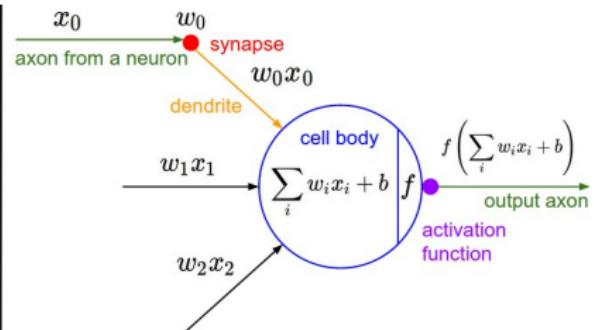
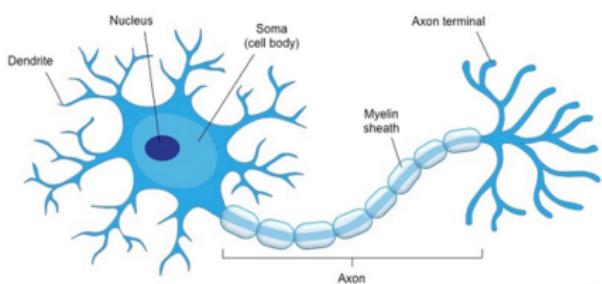
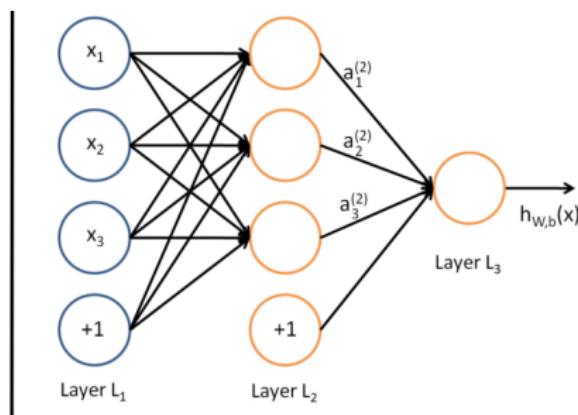
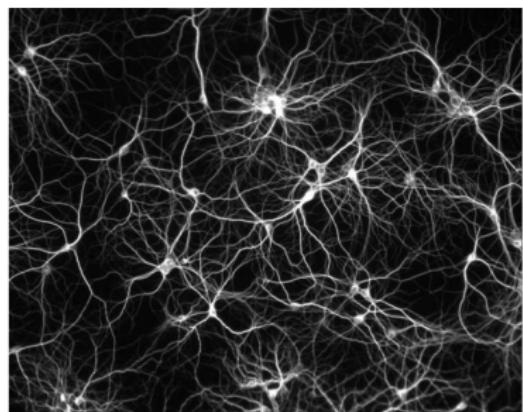


Figura: Comparando os neurônios biológico e artificial (Perceptron).

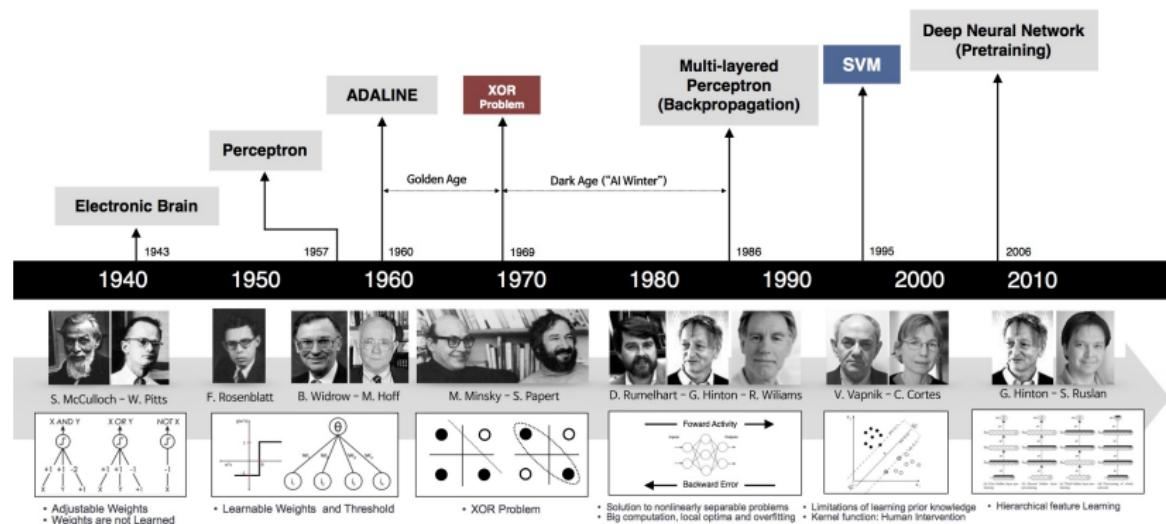
Representação Artificial do Cérebro



- Redes neurais artificiais são sistemas compostos de neurônios artificiais **inspirados** no neurônio biológico.
 - Similarmente combinados na forma de uma malha neural.
 - Esse modelo é **muito limitado** se comparado com o cérebro real.



Representação Artificial do Cérebro



Agenda



1 Origem

2 Fundamentos

- Classificação Linear
- Perceptron
- Otimização
- Redes Neurais Artificiais
- Treinamento

3 Framework

- Comparação Entre Frameworks
- Pytorch

Classificação Linear



- **Mapeamento:** Uma regra que, para cada elemento do conjunto X , provê o elemento correspondente do conjunto Y .

$$f : X \rightarrow Y$$

- Um mapeamento nada mais é do que uma **função**.

$$y = f(x)$$

Classificação Linear



- **Mapeamento:** Uma regra que, para cada elemento do conjunto X , provê o elemento correspondente do conjunto Y .

$$f : X \rightarrow Y$$

- Um mapeamento nada mais é do que uma **função**.

$$y = f(x)$$

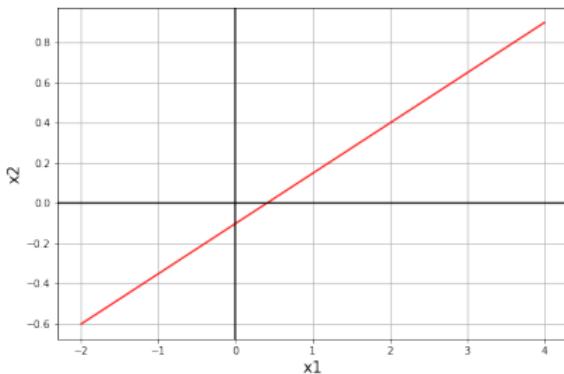
Classificação Linear



- Equação da Reta

$$w_1x_1 + w_2x_2 + b = 0$$

- Um peso w_i por dimensão
- Um viés b



Classificação Linear



- Equação da Reta: Pode servir como uma função de mapeamento.

$$f(x; W, b) = w_1x_1 + w_2x_2 + b$$

- Em outras palavras: Dados os **parâmetros** $W = \{w_1, w_2\}$ e b de uma reta, é possível mapear uma entrada $X = \{x_1, x_2\}$ para a saída $f(x; W, b)$.
- Representação matricial

$$\underbrace{\begin{bmatrix} w_1 & w_2 \end{bmatrix}}_W \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_{X^i} + \begin{bmatrix} b \end{bmatrix} = \begin{bmatrix} \text{score} \end{bmatrix}$$
$$f(X^i; W, b)$$

Classificação Linear



- Equação da Reta:

$$f(x; W, b) = w_1 x_1 + w_2 x_2 + b$$

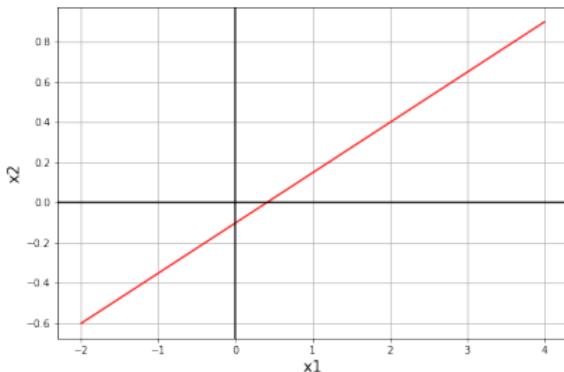
- Vamos supor uma reta fixa:

$$f(x; W, b) = -1x_1 + 4x_2 + 0.4$$

- Representação matricial

$$\underbrace{\begin{bmatrix} -1 & 4 \end{bmatrix}}_{W} \underbrace{\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}}_{X^i} + \begin{bmatrix} 0.4 \end{bmatrix} = \begin{bmatrix} \text{score} \end{bmatrix}$$

$$f(X^i; W, b)$$



Classificação Linear



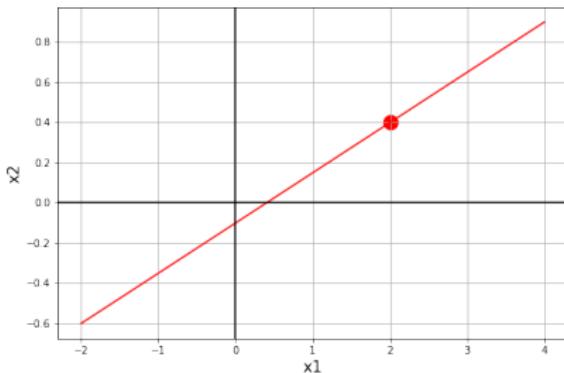
- Vamos supor uma reta fixa:

$$f(x; W, b) = -1x_1 + 4x_2 + 0.4$$

- Para $X = \{2.0, 0.4\}$

$$\underbrace{\begin{bmatrix} -1 & 4 \end{bmatrix}}_{W} + \underbrace{\begin{bmatrix} 2 \\ 0.4 \end{bmatrix}}_{X^i} = \begin{bmatrix} 0 \end{bmatrix}$$

$f(X^i; W, b)$



- $f(x; W, b) = 0$: ponto sobre a reta.

Classificação Linear

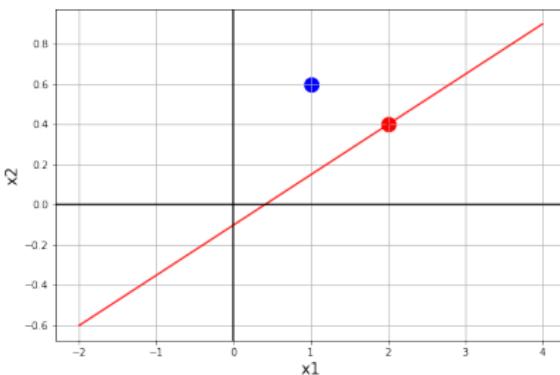


- Vamos supor uma reta fixa:

$$f(x; W, b) = -1x_1 + 4x_2 + 0.4$$

- Para $X = \{1.0, 0.6\}$

$$\underbrace{\begin{bmatrix} -1 & 4 \end{bmatrix}}_{W} + \underbrace{\begin{bmatrix} 1 \\ 0.6 \end{bmatrix}}_{X^i} = \begin{bmatrix} 1.8 \end{bmatrix} \quad f(X^i; W, b)$$



- $f(x; W, b) > 0$: ponto a cima da reta.

Classificação Linear

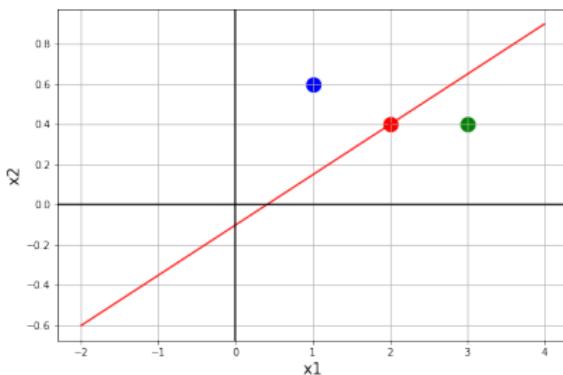


- Vamos supor uma reta fixa:

$$f(x; W, b) = -1x_1 + 4x_2 + 0.4$$

- Para $X = \{3.0, 0.4\}$

$$\underbrace{\begin{bmatrix} -1 & 4 \\ 3 & 0.4 \end{bmatrix}}_{W} + \underbrace{\begin{bmatrix} 0.4 \\ -1 \end{bmatrix}}_{x^i} = f(x^i; W, b)$$



- $f(x; W, b) < 0$: ponto abaixo da reta.

Classificação Linear



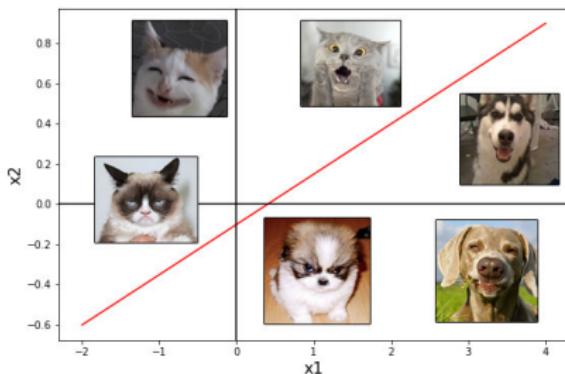
- Equação da Reta

$$f(x; W, b) = w_1 x_1 + w_2 x_2 + b$$

- Representação matricial

$$\begin{array}{c} \begin{matrix} w_1 & w_2 \end{matrix} \\ \underbrace{\phantom{\begin{matrix} w_1 & w_2 \end{matrix}}}_{W} \end{array} + \begin{matrix} x_1 \\ x_2 \end{matrix} + \boxed{b} = \boxed{\text{score}}$$

$$f(X^i; W, b)$$

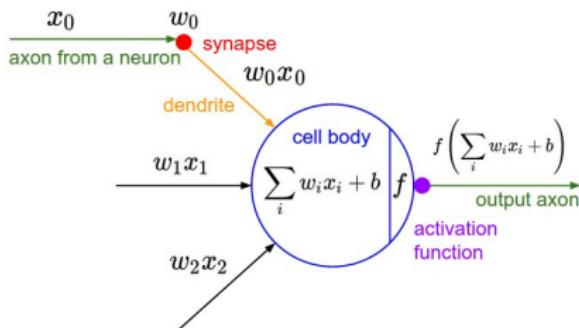


- Em termos gerais, o treinamento de um classificador linear busca os **parâmetros W e b** para compor a **fronteira de decisão** entre as classes.

Perceptron



- Unidade elementar das redes neurais.
- Modelo matemático que mapeia suas entradas para uma saída.
- Exemplo: um ponto em 3 dimensões $X = \{x_0, x_1, x_2\}$ vai compor um perceptron com 3 entradas e 1 saída.



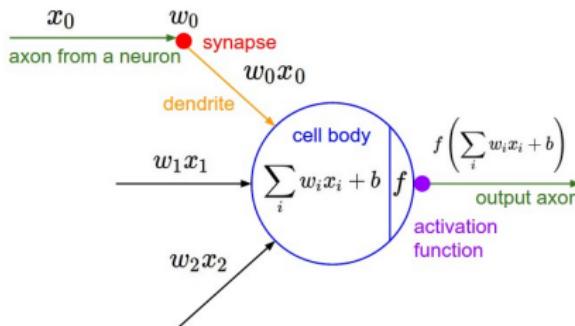
Perceptron



- Mapeamento linear

$$z = \sum_{i=0}^N w_i x_i + b$$

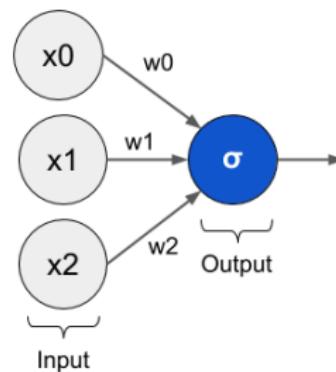
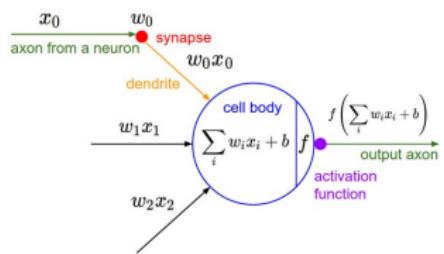
- Seguido de uma função de ativação f



Perceptron



- Pode ser representado como uma rede neural de 2 camadas
 - Camada 1: Entrada. A dimensão dessa camada é definida de acordo com o dado.
 - Camada 2: Saída. Dimensão diretamente ligada ao problema. Por exemplo, para um problema de classificação a dimensão equivale ao número de classes.



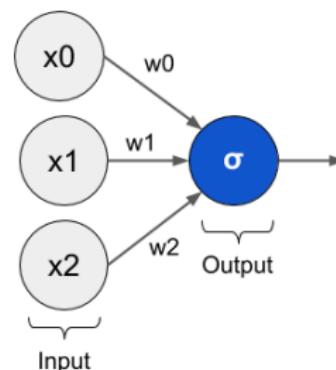
Perceptron



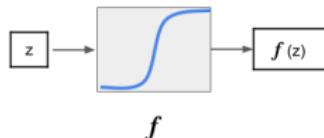
$$z = \sum_{i=1}^N w_i x_i + b$$

- W é o peso de cada ligação
- b é o bias

$$\underbrace{\begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}}_W \quad \underbrace{\begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix}}_X + \begin{bmatrix} b \end{bmatrix} = \begin{bmatrix} z \end{bmatrix}$$

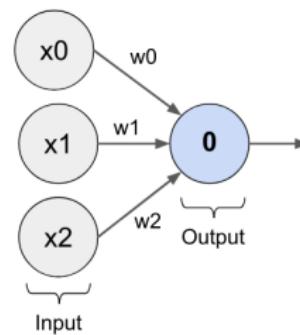
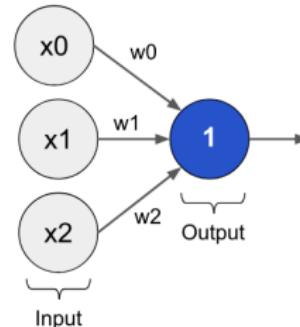


Ativações



- f é a função de ativação
- Por exemplo:

$$f(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{otherwise} \end{cases}$$

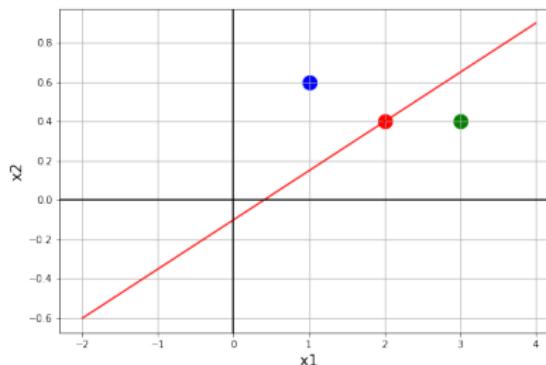


Ativações



- $f(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{otherwise} \end{cases}$

- $z = -1x_1 + 4x_2 + 0.4$



- $z = -1 \times 1 + 4 \times 0.6 + 0.4 = 1.8$

- $f(z) = 1$

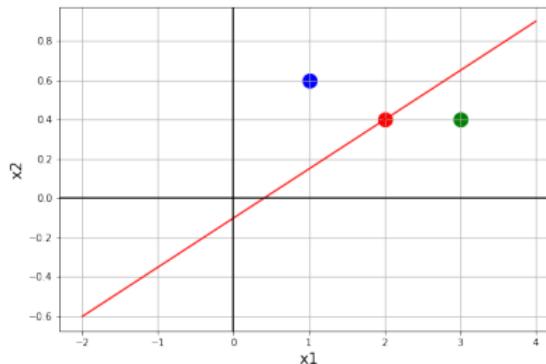
- $z = -1 \times 3 + 4 \times 0.4 + 0.4 = -1$

- $f(z) = 0$

Ativações

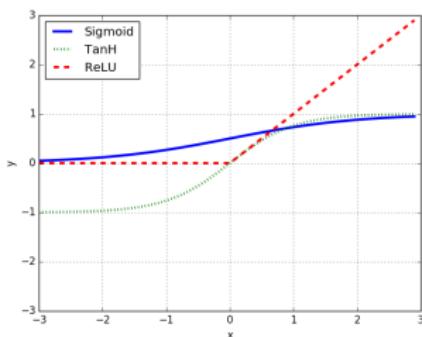


- $f(z) = \begin{cases} 1, & \text{if } z \geq 0 \\ 0, & \text{otherwise} \end{cases}$
- $z = -1x_1 + 4x_2 + 0.4$



- Essa função de ativação permite maior controle sobre a saída do perceptron. A saída z que antes poderia assumir qualquer valor $z = \{-\infty, \dots, \infty\}$ agora está limitada a $z = \{0, 1\}$

Ativações



- Sigmoid (σ)

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

$$f(x) = 2\sigma(2x) - 1$$
- ReLU

$$f(x) = \max(0, x)$$

- Entre outras
 - Threshold, Leaky ReLU, Maxout, Linear, etc.
- Em redes com mais de uma camada escondida, a adição de uma função de ativação não linear permite que a rede modele dependências não lineares.

Perceptron como Classificador Linear



- Função de mapeamento do neurônio

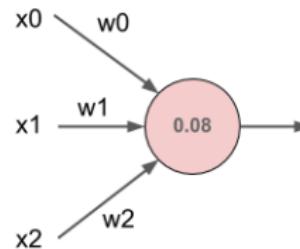
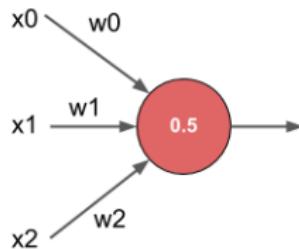
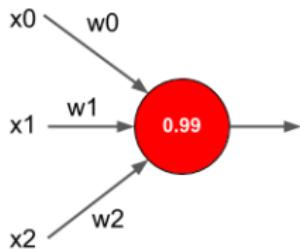
$$f\left(\sum_{i=0}^N w_i x_i + b\right)$$

- O neurônio faz um mapeamento linear, permitindo que seja possível treiná-lo para encontrar os parâmetros da fronteira de decisão tal qual um classificador linear para um problema de classificação binária $Y = \{0, 1\}$.

Perceptron como Classificador Linear



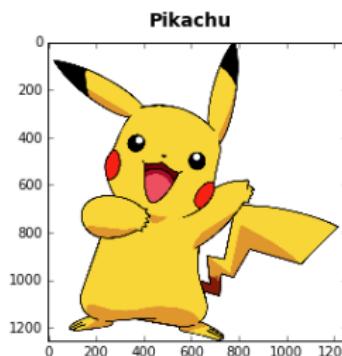
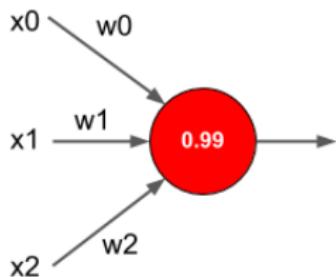
- De forma abstrata, um perceptron interpreta as suas entradas, e libera uma ativação com uma determinada força.



Perceptron como Classificador Linear



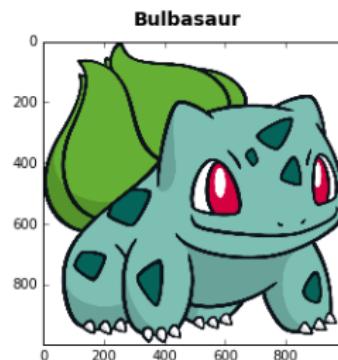
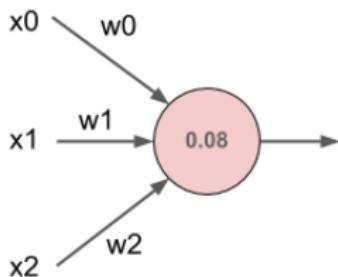
- Exemplo: você treinou um perceptron para fazer uma classificação binária, associando um rótulo numérico para cada classe
 $Y = \{1 : \text{Pikachu}, 0 : \text{Bulbasaur}\}$. Quanto mais próxima a ativação for de 1, mais certeza o perceptron tem da classe "Pikachu".



Perceptron como Classificador Linear



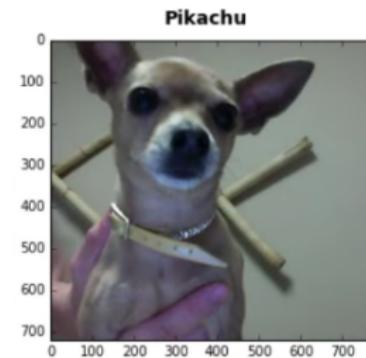
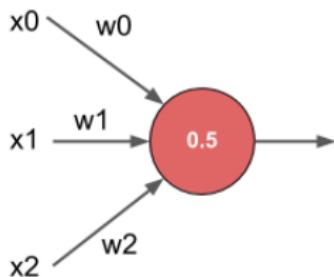
- Exemplo: você treinou um perceptron para fazer uma classificação binária, associando um rótulo numérico para cada classe
 $Y = \{1 : \text{Pikachu}, 0 : \text{Bulbasaur}\}$. Quanto mais próxima a ativação for de 0, mais certeza o perceptron tem da classe "Bulbasaur".



Perceptron como Classificador Linear



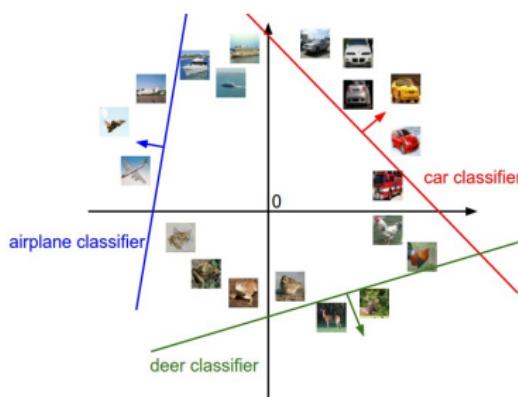
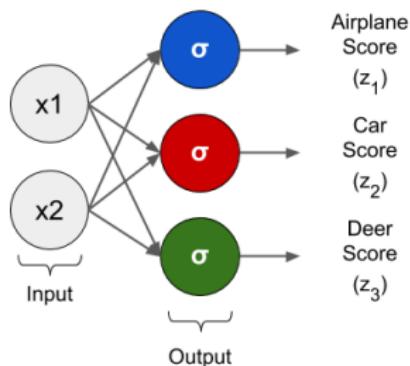
- Exemplo: você treinou um perceptron para fazer uma classificação binária, associando um rótulo numérico para cada classe
 $Y = \{1 : \text{Pikachu}, 0 : \text{Bulbasaur}\}$. Valores intermediários indicam que a imagem está na **zona de dúvida** da classificação.



Perceptron como Classificador Linear



- Para problemas com múltiplas classes, é necessário especializar cada perceptron em uma classe específica.
- Cada perceptron vai otimizar os parâmetros da fronteira de decisão da classe que ele deve se especializar.

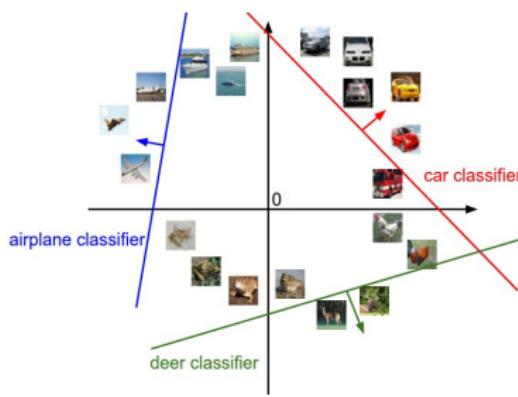


Perceptron como Classificador Linear



- Na representação matricial, isso significa que a matriz de pesos W tem a forma $K \times D$, sendo K o número de classes e D a dimensão do dado. E o vetor de bias b possui a forma $K \times 1$, apenas um bias por classe.
- Em outras palavras, cada linha das matrizes W e b produzem o score de uma classe específica.

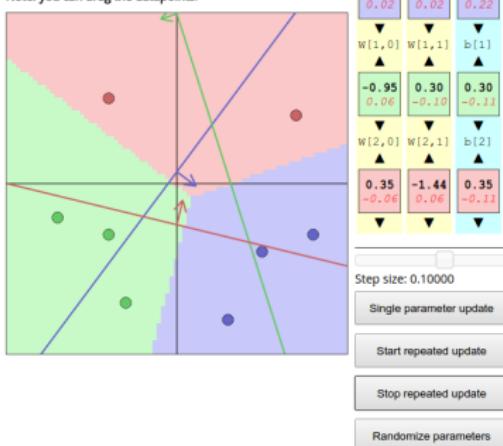
$$\begin{matrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{matrix} \quad \begin{matrix} x_1 \\ x_2 \end{matrix} + \begin{matrix} b_1 \\ b_2 \\ b_3 \end{matrix} = \begin{matrix} z_1 \\ z_2 \\ z_3 \end{matrix}$$



Perceptron como Classificador Linear

Datapoints are shown as circles colored by their class (red/green/blue). The background regions are colored by whichever class is most likely at any point according to the current weights. Each classifier is visualized by a line that indicates its zero score level set. For example, the blue line shows the set of points (x_0, x_1) that give score of zero. The blue arrow draws the vector $(W_{0,0}, W_{0,1})$, which shows the direction of score increase and its length is proportional to how steep the increase is.

Note: you can drag the datapoints.



Parameters W, b are shown below. The value is in **bold** and its gradient (computed with backprop) is in **red, italic** below. Click the triangles to control the parameters.

$W[0,0]$	$W[0,1]$	$b[0]$
1.12 0.02	0.84 0.02	0.06 0.22
$W[1,0]$	$W[1,1]$	$b[1]$
-0.95 0.06	0.30 -0.10	0.30 -0.13
$W[2,0]$	$W[2,1]$	$b[2]$
0.35 -0.08	-1.44 0.06	0.35 -0.13

Visualization of the data loss computation. Each row is loss due to one datapoint. The first three columns are the 2D data x_i and the label y . The next three columns are the three class scores from each classifier $f(x_i; W, b) = Wx_i + b$ (E.g. $s[0] = x[0] * W[0,0] + x[1] * W[0,1] + b[0]$). The last column is the data loss for a single example, L_i .

$x[0]$	$x[1]$	y	$s[0]$	$s[1]$	$s[2]$	L_i
0.50	0.40	0	0.95	-0.05	-0.05	0.00
0.80	0.30	0	1.21	-0.37	0.20	0.00
0.30	0.80	0	1.06	0.26	-0.70	0.19
-0.40	0.30	1	-0.14	0.77	-0.22	0.09
-0.30	0.70	1	0.31	0.80	-0.76	0.51
-0.70	0.20	1	-0.56	1.03	-0.18	0.00
0.70	-0.40	2	0.51	-0.49	1.17	0.34
0.80	-1.15	2	-0.00	-0.81	2.28	0.00
-0.40	-0.50	2	-0.81	0.53	0.93	0.60

mean:

Total data loss: 0.19

Regularization loss: 0.52

Total loss: 0.71

L2 Regularization strength: 0.10000

Multi-class SVM loss formulation:

- Weston Watkins 1999
- One vs. All
- Structured SVM
- Softmax

Figura: Demo interativa da Universidade de Stanford (CS231n):

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>



Funções de Perda

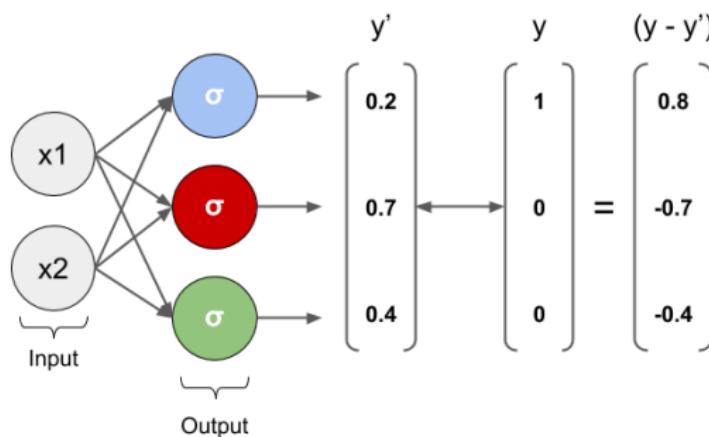


- Primeiro precisamos definir uma medida de qualidade
- Isso é obtido através da função de perda (*loss function*)
 - Também conhecida como função de custo
- A escolha da *loss function* vai depender do problema sendo abordado

Funções de Perda (Exemplo)



- \hat{y} representa a saída da rede (predição) dada a entrada x
 - $f(x; \mathbf{W}, \mathbf{b})$
- y é o rótulo da amostra x
- $(y - \hat{y})$ é a função de custo



Funções de Perda



- \hat{y} representa a saída da rede (predição) dada a entrada x
 - $f(x; \mathbf{W}, \mathbf{b})$
- y é o rótulo da amostra x

1 Erro quadrático

- $\mathcal{L}(W, b; x, y) = \frac{1}{2} \|\hat{y} - y\|^2$

2 Log-Loss (Cross-Entropy)

- $\mathcal{L}(W, b; x, y) = -\sum_x y \log \hat{y}$

3 Entre (inúmeras) outras

- Kullback-Leibler, Hinge, Huber, Mean Absolute Error, ...
- Você pode customizar a sua *loss* (os frameworks dão suporte a isso)

Otimização



- O nosso objetivo então é minimizar essa função de custo
- Como pôde ser visto, a função de custo $\mathcal{L}(W, b; x, y)$ é baseada nas variáveis da rede:
 - pesos W
 - bias b
- Logo, a minização dessa função se dá ao mudar essas variáveis baseado no erro que ela gera

Otimização

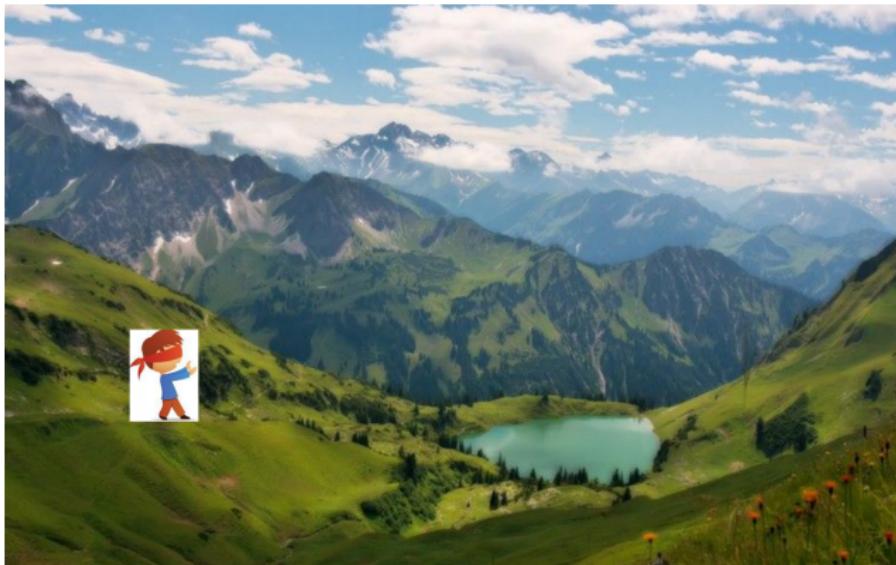


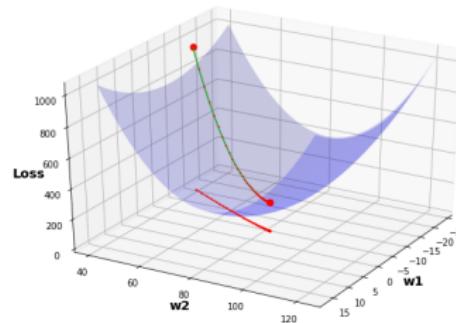
Figura: Analogia do processo de otimização¹.

¹ http://cs231n.stanford.edu/slides/2016/winter1516_lecture3.pdf

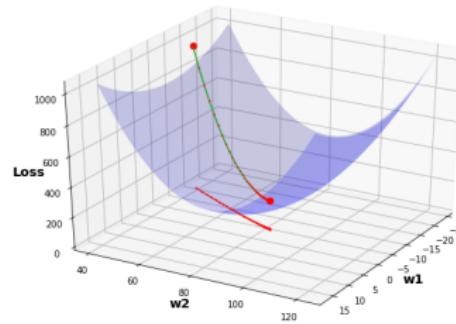
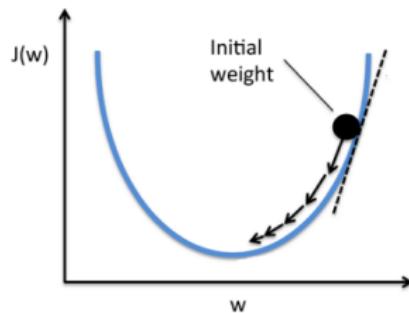
Otimização



- Intuitivamente, o nosso "medidor de altura" é uma função da *loss* em relação aos pesos.



Otimização



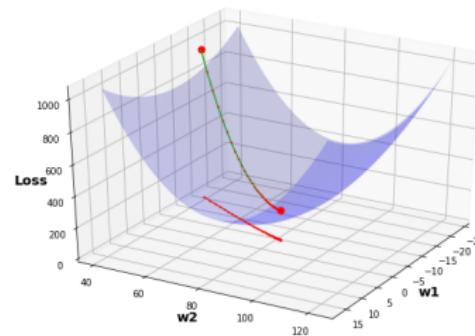
- A direção é dada pela derivada do custo em relação ao peso.
- Para múltiplas dimensões o vetor de derivadas parciais é chamado de **gradiente**.

Otimização



- Em outras palavras, eu consigo alterar os pesos e detectar como a função de perda foi impactada pela minha alteração.

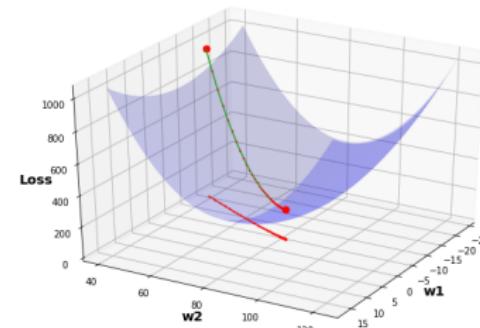
- $\frac{\partial \mathcal{L}(W, b)}{\partial W_{ij}^{(L)}}$



Otimização



- Isso vai indicar:
 - O sinal: se o peso deve aumentar ou diminuir
 - A magnitude: o quanto aquele peso deve ser ajustado.



- $\frac{\partial \mathcal{L}(W, b)}{\partial W_{ij}^{(L)}}$

Otimização



- Descida do gradiente "vanilla"

$$W_{ij}^{(L)} = W_{ij}^{(L)} - \alpha \frac{\partial \mathcal{L}(W,b)}{\partial W_{ij}^{(L)}}$$

$$b_{ij}^{(L)} = b_{ij}^{(L)} - \alpha \frac{\partial \mathcal{L}(W,b)}{\partial b_{ij}^{(L)}}$$

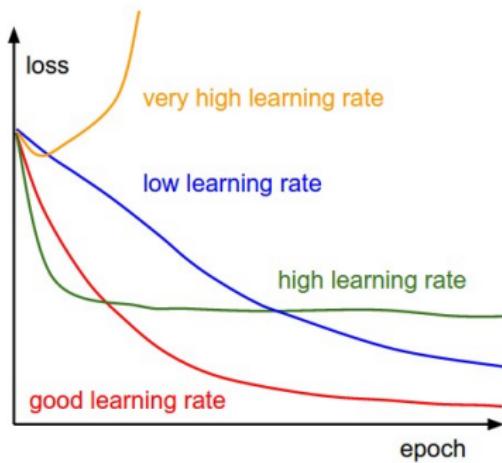
- Sendo α a taxa de aprendizado (*learning rate*)
 - Um fator que interfere na magnitude da alteração dos pesos.

```
for i in range(nb_epochs):
    params_grad = evaluate_gradient(loss_function, data, params)
    params = params - learning_rate * params_grad
```

Otimização



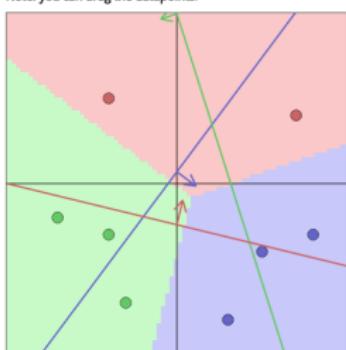
- Taxa de Aprendizado: Define o tamanho do "passo" de otimização



Perceptron como Classificador Linear

Datapoints are shown as circles colored by their class (red/green/blue). The background regions are colored by whichever class is most likely at any point according to the current weights. Each classifier is visualized by a line that indicates its zero score level set. For example, the blue line shows the set of points (x_0, x_1) that give score of zero. The blue arrow draws the vector $(W_{0,0}, W_{0,1})$, which shows the direction of score increase and its length is proportional to how steep the increase is.

Note: you can drag the datapoints.



Parameters W, b are shown below. The value is in **bold** and its gradient (computed with backprop) is in **red, italic** below. Click the triangles to control the parameters.

$w[0,0]$	$w[0,1]$	$b[0]$
1.12 0.02	0.84 0.02	0.06 0.22
-0.95 -0.06	0.30 -0.10	0.30 -0.13
0.35 -0.08	-1.44 0.06	0.35 -0.11

Step size: 0.10000

Single parameter update

Start repeated update

Stop repeated update

Randomize parameters

Visualization of the data loss computation. Each row is loss due to one datapoint. The first three columns are the 2D data x_i and the label y . The next three columns are the three class scores from each classifier $f(x_i; W, b) = Wx_i + b$ (E.g. $s[0] = x[0] * W[0,0] + x[1] * W[0,1] + b[0]$). The last column is the data loss for a single example, L_i .

$x[0]$	$x[1]$	y	$s[0]$	$s[1]$	$s[2]$	L_i
0.50	0.40	0	0.95	-0.05	-0.05	0.00
0.80	0.30	0	1.21	-0.37	0.20	0.00
0.30	0.80	0	1.06	0.26	-0.70	0.19
-0.40	0.30	1	-0.14	0.77	-0.22	0.09
-0.30	0.70	1	0.31	0.80	-0.76	0.51
-0.70	0.20	1	-0.56	1.03	-0.18	0.00
0.70	-0.40	2	0.51	-0.49	1.17	0.34
0.80	-1.15	2	-0.00	-0.81	2.28	0.00
-0.40	-0.50	2	-0.81	0.53	0.93	0.60

mean:

0.19

Total data loss: 0.19

Regularization loss: 0.52

Total loss: 0.71

L2 Regularization strength: 0.10000

Multi-class SVM loss formulation:

Weston Watkins 1999

One vs. All

Structured SVM

Softmax

Figura: Demo interativa da Universidade de Stanford (CS231n):

<http://vision.stanford.edu/teaching/cs231n-demos/linear-classify/>



Otimização



- Métodos de minimização mais populares
 - Gradiente descendente (*Stochastic Gradient Descent (SGD)*) [3]
 - Adam [4]
 - RMSProp [5]
 - Adadelta [6]
 - Entre outros

Otimização



Gif - Algoritmos de Minimização

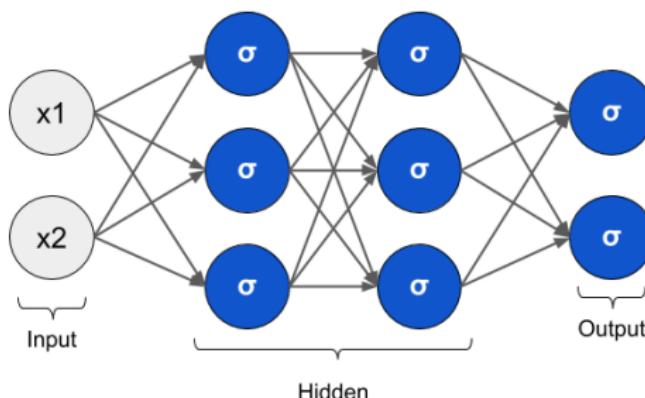
gifs/optimization-2D.gif

gifs/optimization-3D.gif

Multi-Layer Perceptron



- Além das camadas de entrada e saída, temos também as camadas escondidas (*hidden layers*).

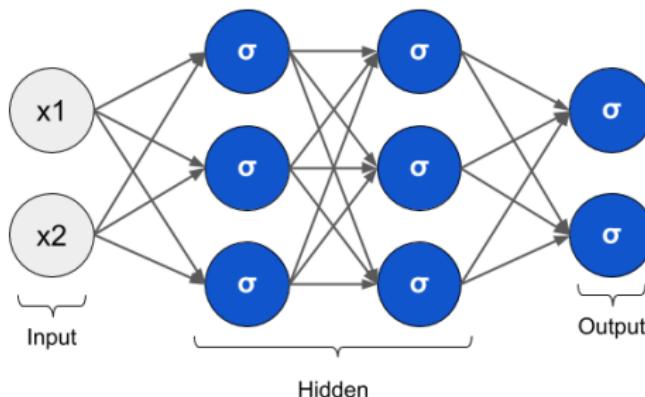


Multi-Layer Perceptron



- Considerando o que já sabemos sobre o perceptron, uma arquitetura multicamada é uma função composta

$$f(f(f(x; W^0, b^0); W^1, b^1); W^2, b^2)$$



Multi-Layer Perceptron



- Considerando o que já sabemos sobre o perceptron, uma arquitetura multicamada é uma função composta

$$f(f(f(x; W^0, b^0); W^1, b^1); W^2, b^2)$$

$$\begin{cases} a^0 = f(x; W^0, b^0) \\ a^1 = f(a^0; W^1, b^1) \\ a^2 = f(a^1; W^2, b^2) \end{cases}$$

Multi-Layer Perceptron



- Teorema de Aproximação Universal²
 - "*Uma rede neural feed forward com apenas uma camada (escondida) é suficiente para representar qualquer função, mas a camada pode ser inviavelmente grande e pode falhar em aprender e generalizar corretamente.*"

²Ian Goodfellow, Deep Learning Book. <http://www.deeplearningbook.org/contents/mlp.html>

Multi-Layer Perceptron



- Teorema de Aproximação Universal ²
 - "*Uma rede neural feed forward com apenas uma camada (escondida) é suficiente para representar qualquer função, mas a camada pode ser inviavelmente grande e pode falhar em aprender e generalizar corretamente.*"
- **Inclusive não linearidades!** (considerando ativações não lineares na camada escondida)

² Ian Goodfellow, Deep Learning Book. <http://www.deeplearningbook.org/contents/mlp.html>

Multi-Layer Perceptron



- Teorema de Aproximação Universal

Qual desses modelos você acha melhor?

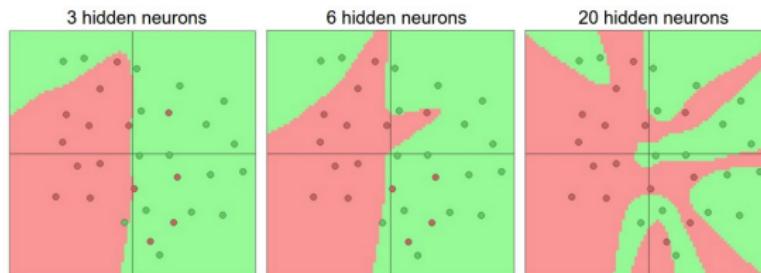


Figura: ConvnetJS demo³

³ConvnetJS demo <https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Multi-Layer Perceptron



- Teorema de Aproximação Universal
- A capacidade das redes neurais de aproximar qualquer função contínua (inclusive as mais complexas), a torna muito vulnerável a **overfitting**.

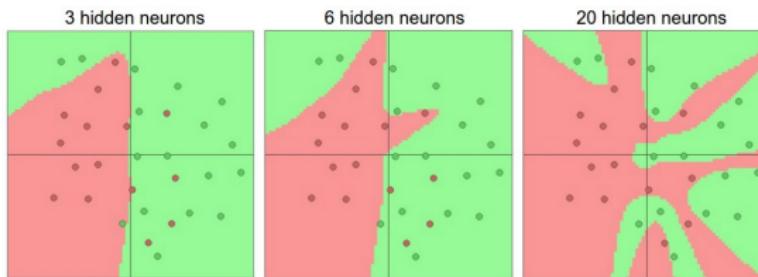


Figura: ConvnetJS demo ³

³ConvnetJS demo <https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Regularização



- Existem artifícios para controlar a complexidade do modelo.

1 Regularização L2

- Adiciona o termo $\frac{1}{2}\lambda W^2$ à função objetivo para todos os pesos da rede
- λ controla a força da regularização
- Pesos decaem em direção a 0

2 Dropout

- Durante o treinamento, cada neurônio tem uma probabilidade p de estar ativo.
- Neurônios inativos não são atualizados naquela iteração
- Não há dropout em tempo de teste

Regularização



- Efeito de aplicar regularização L2 no modelo com 20 neurônios

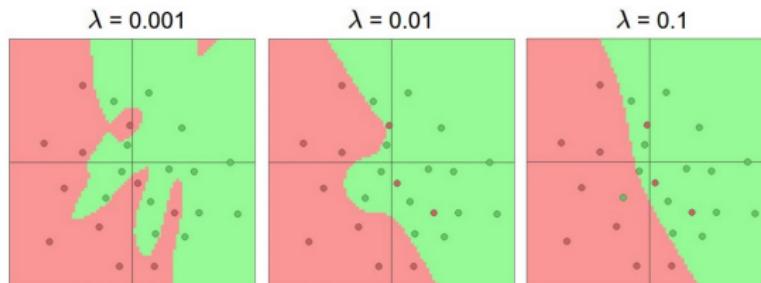


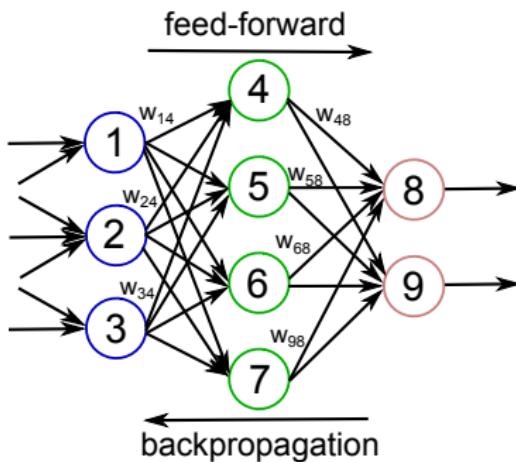
Figura: ConvnetJS demo ³

³ConvnetJS demo <https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

Redes Neurais Artificiais



- O processo de aprendizado pode ser dividido em dois passos:
 - Feed-forward: recebe os dados e gera a saída final da rede
 - Backpropagation: ajusta os pesos de acordo com o gradiente

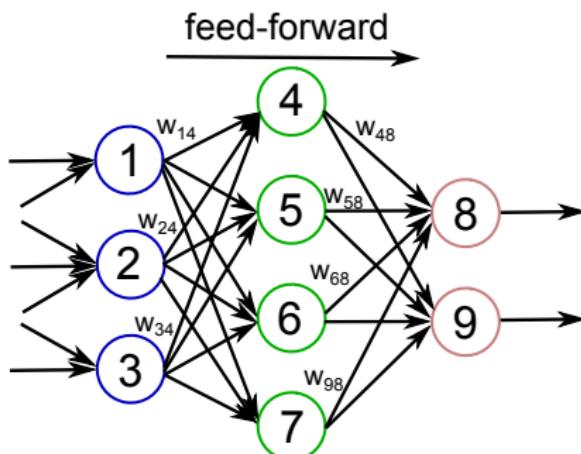


Feed-forward



$$z = \sum_i w_i x_i + b$$

$$a = f(z)$$

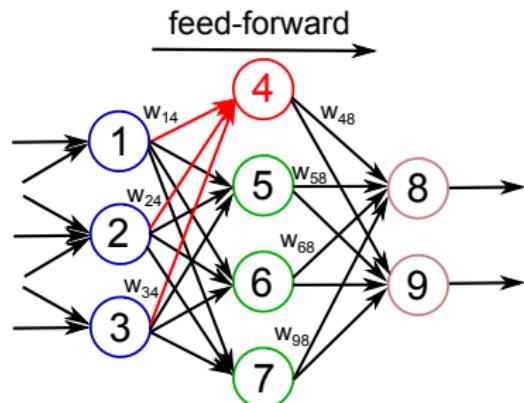


Feed-forward



$$z^4 = W^{14}a^1 + W^{24}a^2 + W^{34}a^3 + b^4$$

$$a^4 = f(z^4)$$

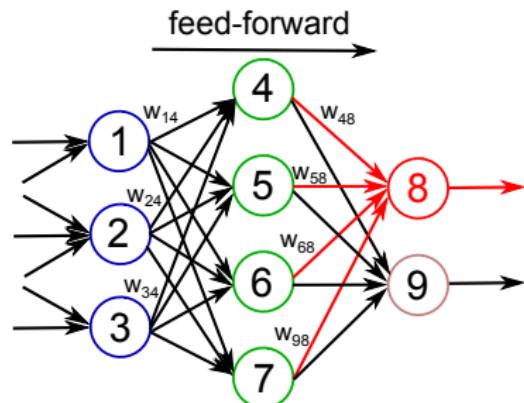


Feed-forward



$$z^8 = W^{48}a^4 + W^{58}a^5 + W^{68}a^6 + W^{78}a^7 + b^8$$

$$a^8 = f(z^8)$$

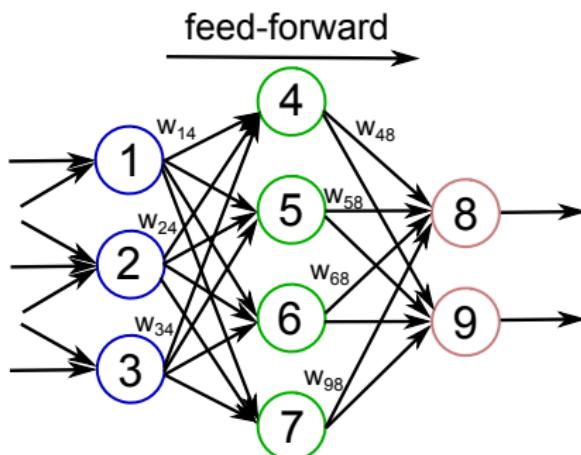


Feed-forward



$$z^{(L+1)} = W^{(L)}a^{(L)} + b^{(L)}$$

$$a^{(L+1)} = f(z^{(L+1)})$$



Backpropagation



- Podemos calcular a velocidade com que o erro muda (usando a *loss function*) à medida que mudamos as saídas dos neurônio
- Para mudar as saídas, precisamos mudar os pesos, o que é feito através do algoritmo de otimização
- Como vimos, precisamos das derivadas parciais
- Backpropagation calcula essas derivadas parciais!

Backpropagation



- Suponhamos:

- Erro quadrático: $\mathcal{L}(W, b; x, y) = \frac{1}{2} \|y - a^{(L)}\|^2$
- Sigmóide como função de ativação: $f(z) = \sigma(z) = \frac{1}{1 + \exp^z}$

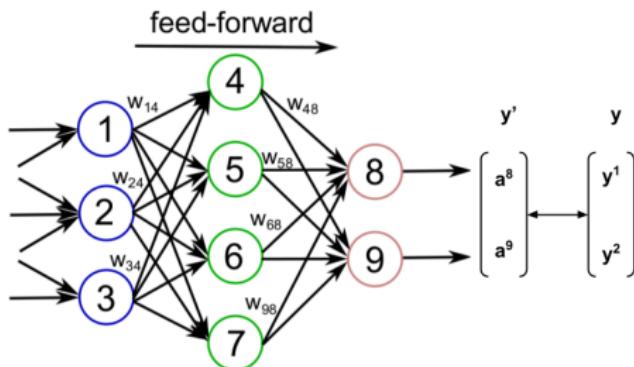
Backpropagation



- Passo a passo, temos:

1 Calcula o erro da rede

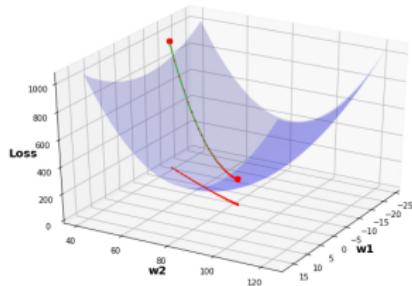
$$\mathcal{L}(W, b; x, y) = \frac{1}{2} \|y - a^{(L)}\|^2$$



Backpropagation



2 Calcula o erro dos neurônios da última camada ⁴



- $\frac{\partial \mathcal{L}}{\partial W^{(L)}} = a^{(L-1)} a^{(L)} (1 - a^{(L)}) (y - a^{(L)})$
- $\delta^{(L)} = a^{(L)} (1 - a^{(L)}) (y - a^{(L)})$
- $\frac{\partial \mathcal{L}}{\partial W^{(L)}} = a^{(L-1)} \delta^{(L)}$

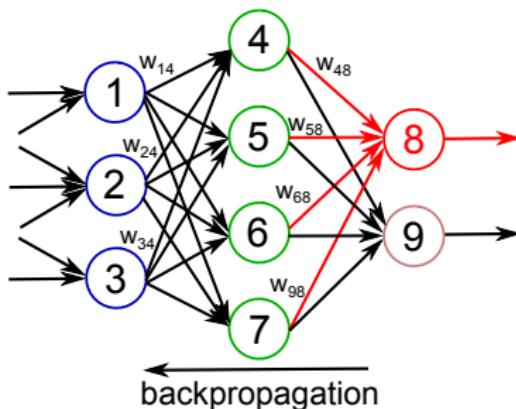
⁴ Backpropagation calculus | Appendix to deep learning chapter 3 <https://youtu.be/tleHLnjs5U8>

Backpropagation



- Propagando o erro para neurônios na última camada

- $\delta^8 = a^8(1 - a^8)(y - a^8)$



Backpropagation



- Propagando o erro para neurônios na última camada

- $W^{(L)} = W^{(L)} - \alpha \frac{\partial \mathcal{L}}{\partial W^{(L)}}$

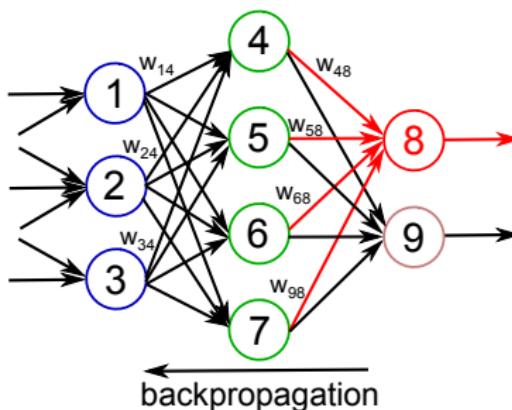
- $\frac{\partial \mathcal{L}}{\partial W^{(L)}} = a^{(L-1)} \delta^{(L)}$

- $W_{48} = W_{48} - \alpha a^4 \delta^8$

- $W_{58} = W_{58} - \alpha a^5 \delta^8$

- $W_{68} = W_{68} - \alpha a^6 \delta^8$

- $W_{78} = W_{78} - \alpha a^7 \delta^8$



Backpropagation



- 3 Calcula os erros das camadas escondidas
- Mais regra da cadeia!

$$\frac{\partial \mathcal{L}}{\partial W^{(L-1)}}$$

- Iterando em todos os neurônios da camada L , o $\delta^{(L-1)}$ de qualquer neurônio da camada $L-1$ é dado por:

$$\delta^{(L-1)} = (\sum_i \delta_i^{(L)} W_i^{(L)}) f'(z^{(L-1)})$$

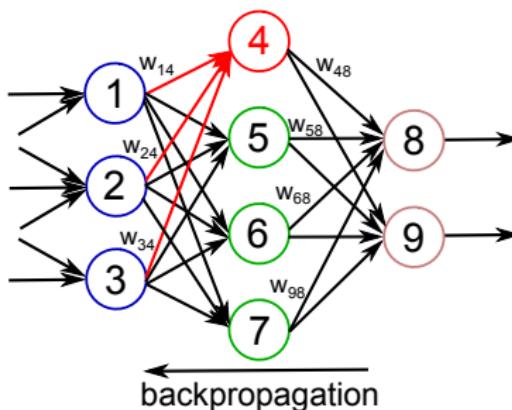
Backpropagation



- Propagando para um neurônio de camada escondida

- $\delta^4 = (\delta^8 W_{48} + \delta^9 W_{49}) a^4 (1 - a^4)$

- $W_{14} = W_{14} - \alpha \delta^4 a^1$
- $W_{24} = W_{24} - \alpha \delta^4 a^2$
- $W_{34} = W_{34} - \alpha \delta^4 a^3$

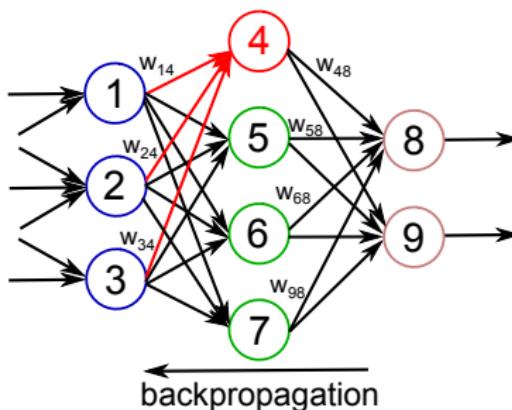


Backpropagation



- Propagando para um neurônio de camada escondida em relação ao bias

- $\delta^4 = (\delta^8 W_{48} + \delta^9 W_{49}) a^4 (1 - a^4)$
- $b_4 = b_4 - \alpha \delta^4 * 1$



Procedimento de Treinamento



- **Configurações Iniciais** (O que sabemos até então)
 - Arquitetura
 - Função de Perda
 - Regularização
 - Otimizador
- Tudo isso é definido experimentalmente (conhecimento prévio pode ajudar bastante) e é diretamente ligado ao problema em questão

Procedimento de Treinamento



- Dinâmica **iterativa** de treinamento
 - A otimização é realizada a cada **passo** de treinamento (analogia da montanha).
 - De forma prática que é um passo de treinamento?

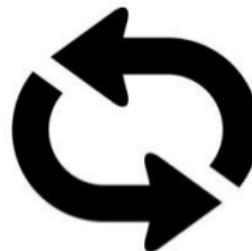
Procedimento de Treinamento



Batch

Iteração

Época



Batch

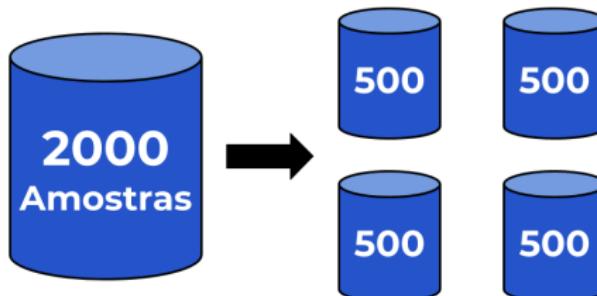


- Todas as amostras do dataset.
- Quantidade de amostras em um passo de treinamento.
 - Um passo: Forward + Backward
 - Em outras palavras, quantas amostras o modelo vai ver antes de calcular a loss e otimizar.

Batch



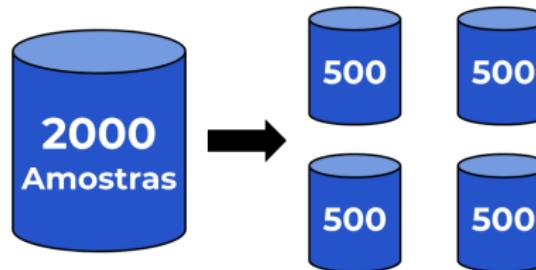
- Por exemplo, dado um dataset de 2000 amostras, poderíamos dividi-lo em mini-batches de 500 amostras.
 - *batch_size* é um hiperparâmetro do treinamento (e geralmente se refere a mini-batches)
 - Nesse caso, $batch_size = 500$



Época



- Uma época é completada quando todas as amostras do dataset foram vistas pelo modelo.
- No nosso exemplo, 1 época se passa a cada 4 iterações.



Época



- *num_epochs* também é hiperparâmetro do treinamento.
- **Mas se ao final de uma época o modelo viu todas as amostras do dataset, por que precisamos de mais de uma época?**

Época



- Modelos profundos são *data hungry* (famintos por dados).
- Na maioria dos casos, uma época não é suficiente para o modelo convergir.
 - Quantidade insuficiente de amostras.
- Nessas condições, é necessário apresentar as amostras mais de uma vez ao modelo.

Agenda



1 Origem

2 Fundamentos

- Classificação Linear
- Perceptron
- Otimização
- Redes Neurais Artificiais
- Treinamento

3 Framework

- Comparação Entre Frameworks
- Pytorch

Framework

Comparação Entre Frameworks

Frameworks



Software	Creator	Software License	Open source	Platform	Written in	Interface	OpenMP support	OpenCL support	CUDA support	Parallel executor (host-node)	Automatic differentiation	Has pretrained	Requires net	Convolutional	RNN/Echo	Model support
nnlib (old)	Karim Lak	MIT	Yes	Linux, mac OS, Windows	Python	Python			Yes		Yes	Yes	Yes	Yes		
BogDL	Jacques Ertl	Apache 2.0	Yes	Apache Software Foundation	Scala	Scala, Python			No			Yes	Yes	Yes	Yes	
Caffe	Berkeley Vision and Learning Center	BSD license	Yes	Linux, mac OS, Windows	C++	Python, MATLAB, C++	Yes	Under development[1]	Yes	Yes	Yes	Yes	Yes	Yes	Yes	?
DeepLearning4J	Stanford engineering team, DeepLearning4J community, originally Alex Krizhevsky	Apache 2.0	Yes	Linux, mac OS, Windows, Android (Cross-platform)	C++, Java	Java, Scala, Cognos, Python (Others, Keras)	Yes	On roadmap[1]	Yes[2]	Computational Graph	Yes[2]	Yes	Yes	Yes	Yes	Yes[2]
Chainer	Preferred Networks	MIT license	Yes	Linux, mac OS, Windows	Python	Python	No	No[3]	Yes[4]		Yes	Yes	Yes	Yes		
Stackel	Joseph Redden	Public Domain, Apache License	Yes	Cross-Platform	C	C, Python	Yes	No[3]	Yes	Yes	Yes	Yes	Yes	Yes		
DLB	Daniel King	Boost Software License	Yes	Cross-Platform	C++	C++	Yes	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
TensorFlow (Old)	D. Chen et al.	Apache 2.0	Yes	TensorFlow	Java	Java	No	No	No	No	No	No	No	No	No	No
PyTorch	Facebook AI Research	Apache 2.0	Yes	Linux, mac OS, Windows	C++, Python	Python	No	No[3]	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
Intel Data Analytics Acceleration Library	Intel	Apache License 2.0	Yes	Linux, mac OS, Windows	C++, Python, Java	Python, Java	Yes	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Realistic Kernel Library	Intel	Proprietary	No	Linux, mac OS, Windows, as Intel Driver[1]	C/C++	C/C++	No	No[3]	No	Yes	Yes	Yes	Yes	Yes	Yes	No
Kaldi	François Fleuret	MIT license	Yes	Linux, mac OS, Windows	Python	Python, R	Only using TensorFlow backend	Under development for TensorFlow (and on roadmap to Kaldi)	Yes	Yes	Yes[5]	Yes	Yes	Yes	Yes[5]	
LightGBM (Hierarchical TreeBoost)	Matthew Wistow	Proprietary	No	Linux, mac OS, Windows	C, C++, Java, MATLAB	MATLAB	No	No								
Microsoft Cognitive Toolkit	Microsoft Research	MIT License[6]	Yes	Windows, Linux[7] (macOS, Linux via Docker on roadmap)	C++	Python (TensorFlow), C++, Cuda, PyTorch, Go, R, Scala, Perl, Small C++, Java, Python, Julia, Matlab, Cuda, Go, R, Scala, Perl	Yes[8]	No								
Apache MLLib	Apache Software Foundation	Apache 2.0	Yes	Linux, mac OS, Windows, MPFR, Java, Python, JavaScript	C/C++	Python (TensorFlow), C++, Cuda, PyTorch, Go, R, Scala, Perl	Yes	On roadmap[1]	Yes	Yes[9]	Yes[9]	Yes	Yes	Yes	Yes[9]	
TensorFlow	Google Brain Team	Apache 2.0	Yes	Linux, mac OS, Windows	C/C++	Graphical user interface	Yes	No	No	?	?	No	No	No	No	?
TensorLayer	Hao Dang	Apache 2.0	Yes	Cross-platform	C/C++	Python	Yes	No	Yes	?	?	No	No	No	No	?
Theano	University of Montréal	BSD license	Yes	Linux, mac OS, Windows	C++, Python, Matlab, Python, C	Python (TensorFlow), C++, C	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Theano	Pierre又有, Sam Sri, Sundeep Chirikar, Gregory Charisin	BSD license	Yes	Linux, mac OS, Windows	Python, C, C++, C#	Python	Yes	Via repository maintained package[10]	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Apache Mahout	Apache Incubator	Apache 2.0	Yes	Linux, mac OS, Windows	Python, C++, Java	Python, C++, Java	No	No	?	Yes	Yes	Yes	Yes	Yes	Yes	Yes
TensorFlow	Google Brain Team	Apache 2.0	Yes	Windows, Linux, Android, iOS	C/C++	C/C++	No	On roadmap[1] but already with GPU support	Yes	Yes[11]	Yes[11]	Yes	Yes	Yes	Yes	Yes
TensorLayer	Hao Dang	Apache 2.0	Yes	Linux, mac OS, Windows, Android	C/C++	Python	No	On roadmap[1] but already with GPU support	Yes	Yes[12]	Yes[12]	Yes	Yes	Yes	Yes	Yes
Theano	Université de Montréal	BSD license	Yes	Cross-platform	Python	Python (TensorFlow)	Yes	Under development[1]	Yes	Yes[13]	Through Language model API[13]	Yes	Yes	Yes	Yes	Yes[13]
Theano	Pierre又有, Sam Sri, Sundeep Chirikar, Gregory Charisin	BSD license	Yes	Windows, Linux, C/C++, Java, Python	C, C++, Python, Go, Angular	Python (TensorFlow), C, C++, Java, Python	Yes	Third party implementation[14]	Yes[15]	Through Tensorflow API[15]	Yes[15]	Yes	Yes	Yes	Yes	Yes[15]
TensorMathematics	Wolfram Research	Proprietary	No	Windows, Linux, Cloud computing	Mathematica Language	Mathematica Language	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Under Development
VisualDF	Wei Li	Proprietary	No	Linux, Web-based	Graphical user interface, C/C++	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes

Figura: Comparação das características dos principais frameworks de Deep Learning¹.

¹ https://en.wikipedia.org/wiki/Comparison_of_deep_learning_software

Frameworks



● Caffe²

- Berkeley AI Research (BAIR)
- Linguagens: **Python**, C++, Matlab
- Pouca customização

² <http://caffe.berkeleyvision.org/>

Frameworks



- Torch³

- Ronan Collobert, Koray Kavukcuoglu, Clement Farabet
- Linguagem: Lua
- Problemas com o suporte a bibliotecas externas no Lua

³<http://torch.ch/>

Frameworks



- Tensorflow⁴
 - Google Brain team
 - Linguagens: **Python**, C/C++, Java, Go, R, Julia
 - Grafos Estáticos
 - Comunidade vasta e participativa

⁴<https://www.tensorflow.org/>

Frameworks



- Pytorch⁵
 - Facebook
 - Linguagem: Python
 - Grafos Dinâmicos
 - Vasta gama de modelos “off-the-shelf”
 - Eficiente, simples e com várias funcionalidades além do treinamento de NNs

⁵ <https://pytorch.org/>

Grafos Dinâmicos



A graph is created on the fly



```
from torch.autograd import Variable  
  
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))
```

Figura: Grafos Dinâmicos.

Grafos Dinâmicos



A graph is created on the fly

```
from torch.autograd import Variable  
  
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))  
  
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())
```

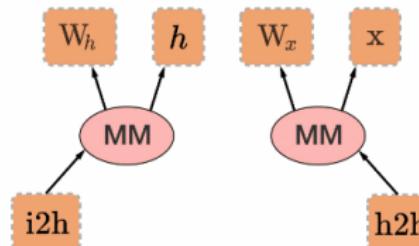


Figura: Grafos Dinâmicos.

Grafos Dinâmicos



A graph is created on the fly

```
from torch.autograd import Variable  
  
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))  
  
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h
```

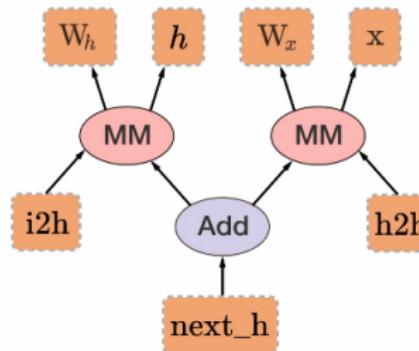


Figura: Grafos Dinâmicos.

Grafos Dinâmicos



A graph is created on the fly

```
from torch.autograd import Variable  
  
x = Variable(torch.randn(1, 10))  
prev_h = Variable(torch.randn(1, 20))  
W_h = Variable(torch.randn(20, 20))  
W_x = Variable(torch.randn(20, 10))  
  
i2h = torch.mm(W_x, x.t())  
h2h = torch.mm(W_h, prev_h.t())  
next_h = i2h + h2h  
next_h = next_h.tanh()
```

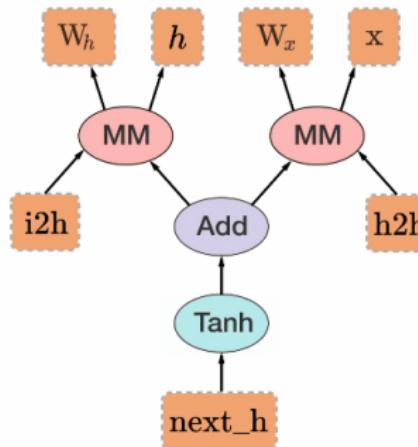


Figura: Grafos Dinâmicos.

Grafos Dinâmicos



Back-propagation
uses the dynamically built graph

```
from torch.autograd import Variable

x = Variable(torch.randn(1, 10))
prev_h = Variable(torch.randn(1, 20))
W_h = Variable(torch.randn(20, 20))
W_x = Variable(torch.randn(20, 10))

i2h = torch.mm(W_x, x.t())
h2h = torch.mm(W_h, prev_h.t())
next_h = i2h + h2h
next_h = next_h.tanh()

next_h.backward(torch.ones(1, 20))
```

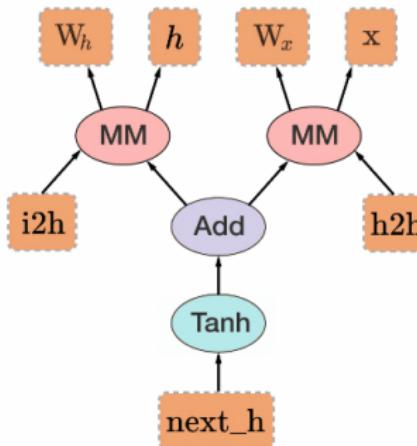


Figura: Grafos Dinâmicos.

Pytorch



- Tensores
- Data Loader
- Definindo um modelo
- Funções de Perda
- Procedimento de treinamento

Tensores



- Sintaxe parecida com ndarrays (Numpy)
- Armazena tanto os inputs e labels quanto os pesos e viéses das NNs
- Funções para tradução de ndarrays para tensores e vice-versa

Tensores



Data type	dtype	CPU tensor
32-bit floating point	<code>torch.float32</code> or <code>torch.float</code>	<code>torch.FloatTensor</code>
64-bit floating point	<code>torch.float64</code> or <code>torch.double</code>	<code>torch.DoubleTensor</code>
16-bit floating point	<code>torch.float16</code> or <code>torch.half</code>	<code>torch.HalfTensor</code>
8-bit integer (unsigned)	<code>torch.uint8</code>	<code>torch.ByteTensor</code>
8-bit integer (signed)	<code>torch.int8</code>	<code>torch.CharTensor</code>
16-bit integer (signed)	<code>torch.int16</code> or <code>torch.short</code>	<code>torch.ShortTensor</code>
32-bit integer (signed)	<code>torch.int32</code> or <code>torch.int</code>	<code>torch.IntTensor</code>
64-bit integer (signed)	<code>torch.int64</code> or <code>torch.long</code>	<code>torch.LongTensor</code>

Figura: Tipos de Tensores.

Tensores



Demo - Tensors

Tensors.ipynb

Data Loader



- Classes *Dataset* e *DataLoader*
 - *DataLoader*: normalmente usada a versão original
 - *Dataset*: normalmente customizada
- Subclasses da classe *Dataset* possuem duas funções principais:
 - `make_dataset()`
 - `__getitem__()`

Data Loader



Demo - Data Loader

Data_Loader.ipynb

Definindo uma Arquitetura



- O pacote *torch.nn* contém as principais camadas
 - **Convolução/Deconvolução:** Conv1d, Conv2d, Conv3d, ConvTranspose1d, ConvTranspose2d, ConvTranspose3d
 - **Pooling/Unpooling:** MaxPool1d, MaxPool2d, MaxPool3d, MaxUnpool1d, MaxUnpool2d, MaxUnpool3d
 - **Funções de Ativação:** ReLU, Sigmoid, Tanh, Softplus, Softmax, LogSoftmax
 - **Regularização:** Dropout, Dropout2d, Dropout3d
 - **Normalização:** BatchNorm1d, BatchNorm2d, BatchNorm3d, InstanceNorm1d, InstanceNorm2d, InstanceNorm3d
 - **Recorrência:** LSTM, GRU
 - **Densa:** Linear, Bilinear

Dimensões de Entrada e Saída



- Camadas Fully Connected: Linear, Softmax...
 - Input: $B \times * \times F_{in}$
 - B = batch size
 - $*$ = qualquer número de dimensões adicionais
 - F_{in} = features de entrada
 - Output: $B \times * \times F_{out}$
 - F_{out} = features de saída

Dimensões de Entrada e Saída



- Muitas camadas fazem operações em tensores de quaisquer dimensões e retornam um tensor do mesmo tamanho da entrada: ReLU, Dropout, Sigmoid, Tanh...
 - Input: *
 - * = qualquer número de dimensões
- Output: $B \times * \times F_{out}$

Definindo uma Arquitetura



Demo - Architecture

Architecture.ipynb

Funções de Perda



- CrossEntropyLoss – Classificação
- NLLLoss – Classificação
- MSELoss – Regressão/Reconstrução
- L1Loss – Regressão/Reconstrução
- KLDivLoss – Reconstrução
- HingeEmbeddingLoss – Aprendizado Semi-Supervisionado

Funções de Perda



Funções de Perda

Funções de Perda diferentes recebem inputs e targets com dimensões diferentes

Funções de Perda



Demo - Loss Functions

Loss_Functions.ipynb

Otimizadores



- Gradient Descent (GD)
- Stochastic Gradient Descent (SGD)
- Momentum
- AdaGrad (Adaptive Gradient)
- RMSprop
- Adam

Otimizadores



Comparação de Otimizadores

Para mais informações sobre otimizadores, visitar o blog post “An overview of gradient descent optimization algorithms”⁶

⁶ <http://ruder.io/optimizing-gradient-descent/index.html>

Introdução ao Pytorch



Demo - Introdução ao Pytorch

Introduction_to_Pytorch.ipynb

Procedimento de Treinamento



Exercício Prático - Classificação de Faces

MLP_Classification_Olivetti.ipynb

└ Framework

 └ Pytorch

Procedimento de Treinamento



Exercício Prático - Classificação de Dígitos

MLP_Classification_MNIST.ipynb

Procedimento de Treinamento



Exercício Prático - Classificação de Peças de Roupa

MLP_Classification_FashionMNIST.ipynb

└ References

- [1] Aleksandra C.
How is a neuron adapted to perform its function?, 2018.
- [2] Fractal Foundation.
Fractal neurons, 2018.
- [3] Herbert Robbins and Sutton Monro.
A stochastic approximation method.
In Herbert Robbins Selected Papers, pages 102–109. Springer, 1985.
- [4] Diederik P Kingma and Jimmy Ba.
Adam: A method for stochastic optimization.
arXiv preprint arXiv:1412.6980, 2014.
- [5] Geoffrey Hinton.
Overview of mini-batch gradient descent, 2018.
- [6] Matthew D Zeiler.
Adadelta: an adaptive learning rate method.
arXiv preprint arXiv:1212.5701, 2012.

