

# Introdução ao Aprendizado Profundo

## Deep Learning 101

---

**María Fernanda Rodríguez Ruiz**

Prof. Dr. José Mario De Martino

Maio de 2018

Departamento de Engenharia de Computação e Automação Industrial — LCA  
Faculdade de Engenharia Elétrica e de Computação — FEEC  
Universidade Estadual de Campinas — UNICAMP



# Agenda

1. Introdução
2. Ambiente de Desenvolvimento
3. Aprendizado de Máquina
4. Aprendizado Profundo
5. Revisão Final

# Introdução

---

## Objetivo

- Apresentar as **noções básicas** que envolvem o conceito de **Aprendizado Profundo** ou *Deep Learning*.

## Modelos

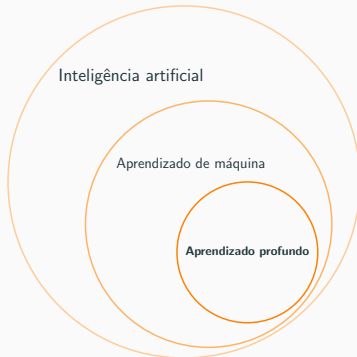
Apresentar alguns dos **modelos** do aprendizado de máquina:

- Regressão linear,
- Regressão logística,
- Redes neurais artificiais,
- Redes neurais artificiais profundas, e
- Redes neurais convolucionais.

## Implementação

- Breve introdução **teórica** e
- Implementações **práticas** através de um mesmo problema.
  - Os exemplos serão desenvolvidos fazendo uso do **Python** e **Keras/TensorFlow** em **Jupyter**.

## Disciplinas da Inteligência Artificial — IA



**Figura 1:** Diagrama de Venn da relação entre algumas das disciplinas da Inteligência Artificial [Goodfellow et al., 2016].

## Aprendizado de máquina — AM

Disciplina da IA que cria sistemas que **aprendem a partir de dados ou com sua própria experiência** a resolver tarefas específicas sem que estas precisem ser explicitamente programadas [Goodfellow et al., 2016].

- **Aprendizado profundo**

Técnicas de AM baseada em **múltiplos níveis de aprendizagem**, para análise de grandes volumes de dados [Goodfellow et al., 2016].



## Aplicações do AM

### Dados

- Binário
- Séries Temporais
- Texto
- Sons
- Imagens
- ...

### Áreas

- Agricultura
- Robótica
- Ciência
- Entretenimento
- Finanças
- ...

# Ambiente de Desenvolvimento

---

## Ferramentas

- *Frameworks* (aprendizado de máquina e profundo),
- *Python* e
- *Jupyter*.

## Frameworks para aprendizado de máquina e aprendizado profundo

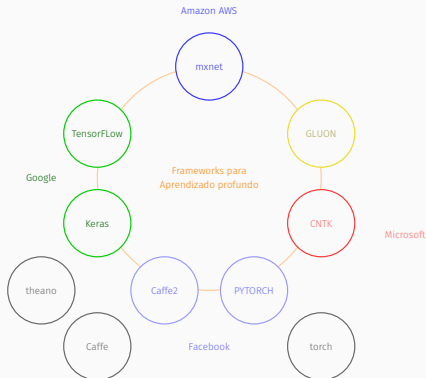


Figura 2: Frameworks<sup>1</sup> de aprendizado profundo de código aberto.

<sup>1</sup>The 5 Deep Learning Frameworks — <https://heartbeat.fritz.ai/>

## Frameworks para aprendizado de máquina e aprendizado profundo

Operam em 2 níveis de abstrações<sup>2</sup>:

- **Low Level** — operações matemáticas e primitivas de redes neurais são implementadas (**TensorFlow**, Theano, PyTorch.)
- **High Level** — primitivas de baixo nível são usados para implementar abstrações de redes neurais, como modelos e camadas (**Keras**).

---

<sup>2</sup>The 5 Deep Learning Frameworks — <https://heartbeat.fritz.ai/>

## TensorFlow<sup>3</sup>



- Biblioteca de software de **código aberto** para computação numérica de alto desempenho, do **Google**.
- **Arquitetura flexível** em várias plataformas CPUs, GPUs, TPUs, e desde *desktops* a dispositivos móveis.
- Forte suporte para **aprendizado de máquina** e **aprendizado profundo**.

---

<sup>3</sup><https://www.tensorflow.org/>

## Keras<sup>4</sup>



- API de redes neurais de **alto nível**, escrita em Python.
- Pode ser executado em cima do **TensorFlow**, CNTK ou Theano
- Foi desenvolvida para permitir a **experimentação rápida**.
- **Suporta** redes convolucionais e redes recorrentes.
- **Funciona** para CPU e GPU.

---

<sup>4</sup><https://keras.io/>

## Python<sup>5</sup>



- Linguagem de programação interpretada de **alto nível**.
- Pacotes:
  - **NumPy**<sup>6</sup>: suporta arrays e matrizes multidimensionais.
  - **Matplotlib**<sup>7</sup>: biblioteca para gerar gráficos a partir de dados contidos em listas ou matrizes.

---

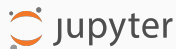
<sup>5</sup><https://www.python.org/>

<sup>6</sup><http://www.numpy.org/>

<sup>7</sup><https://matplotlib.org/>



## Jupyter notebook<sup>8</sup>



- **Linguagens (*kernels*):**
  - Julia, Python, R, MATLAB/Octave, Javascript, C, ....
- **Aplicação Web:**
  - Ambiente de computação interativo que permite criar (escrever e executar) documentos/código no navegador.
  - Markdown<sup>9</sup>, HTML, LaTeX, PNG, SVG, PDF
- **Documento interativo:**
  - Sequência linear de células

---

<sup>8</sup><http://jupyter.org/>

<sup>9</sup><https://www.markdownguide.org/>

## Fluxo de trabalho Keras

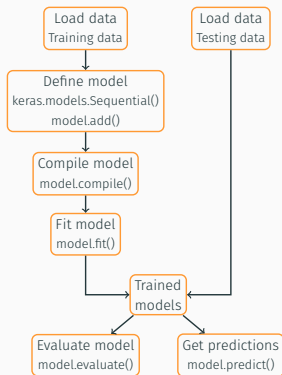


Figura 3: *Workflow* Keras

## Treinamento e avaliação dos modelos.

1. Load Data
2. Define Model
3. Compile Model
4. Fit Model
5. Evaluate Model

## 1. Load data

- Ler dados,
  - Conjunto de **treinamento** e
  - Conjunto de **teste**.
- Visualizar dados,
- Pre-processar.

```
X_train, y_train = load_data_train()  
X_test, y_test = load_data_test()
```

## 2. Define model

Dois modelos: **Sequential** e **Functional API**.

**Sequential**<sup>10</sup> usado para empilhar camadas (layers),  
**model.add()**<sup>11</sup> usada para adicionar as camadas,  
**input\_shape=()** especificar a forma de entrada.

```
model = keras.models.Sequential()  
model.add(layer1 ..., input_shape=(nFeatures))  
model.add(layer2 ... )
```

---

<sup>10</sup><https://keras.io/models>

<sup>11</sup><https://keras.io/layers>

## 3. Compile model

Configurar o **processo de aprendizado** especificando:

**Otimizador**<sup>12</sup> que determina como os pesos são atualizados,

**Função de custo**<sup>13</sup> ou função de perda,

**Métricas**<sup>14</sup> para avaliar durante o treinamento e o teste.

```
model.compile(optimizer='SGD',  
              loss='mse',  
              metrics=['accuracy'])
```

---

<sup>12</sup><https://keras.io/optimizers>

<sup>13</sup><https://keras.io/losses>

<sup>14</sup><https://keras.io/metrics>

## 4. Fit model

Iniciar o **processo de treinamento**.

**batch\_size**<sup>15</sup>: dividir o conjunto de dados em número de lotes.

**epochs**: número de vezes que é treinado o conjunto de dados completo.

```
model.fit(X_train, y_train,  
          batch_size=500,  
          epochs=1)
```

---

<sup>15</sup>Epoch vs Batch Size vs Iterations — <https://towardsdatascience.com>

## 5. Evaluate model

Avaliar o **desempenho** do modelo.

`model.evaluate()` encontra a perda e as métricas especificadas.

Fornece uma medida **quantitativa** da precisão.

`model.predict()` encontra a saída para os dados de teste fornecidos e é útil para verificar as saídas **qualitativamente**.

```
history = model.evaluate(X_test, y_test)
y_pred = model.predict(X_test)
```

# Aprendizado de Máquina

---



Para entender o **aprendizado profundo**, é preciso ter uma sólida compreensão dos princípios básicos do aprendizado de máquina [Goodfellow et al., 2016].

## Aprendizado de Máquina

É tudo sobre a criação de um **algoritmo que pode aprender** com os dados para fazer uma previsão [Goodfellow et al., 2016].

## Noções matemáticas

Quatro grandes **disciplinas matemáticas**<sup>16</sup> compõem o aprendizado de máquina:

- **Estatística** é um núcleo de tudo, nos diz qual é o nosso objetivo.
- **Cálculo** nos diz como aprender e otimizar nosso modelo.
- **Álgebra linear** torna a execução dos algoritmos viáveis em conjuntos de dados massivos.
- **Probabilidade** ajuda a prever a possibilidade de um evento ocorrer.

---

<sup>16</sup>[https://github.com/llSourcell/math\\_of\\_machine\\_learning](https://github.com/llSourcell/math_of_machine_learning)

## Aprendizado de Máquina

Pode ser classificada em [Goodfellow et al., 2016]:

- **Supervisionado:** conjunto de dados de características, associada a um rótulo ou alvo,
- **Não supervisionado:** conjunto de dados de muitas características que aprendem propriedades úteis da estrutura desse conjunto de dados, e
- **Por reforço:** algoritmos interagem com um ambiente, portanto há um ciclo de feedback entre o sistema de aprendizado e suas experiências.

## Regressão linear

- **Algoritmo simples** de aprendizado de máquina, resolve um problema de **regressão** [Goodfellow et al., 2016]:
  - Na **estatística**, é uma técnica para modelar a relação entre uma variável dependente  $y$ , e uma ou mais variáveis independentes  $x$ .
- A saída é uma **função linear** da entrada,

$$\hat{y} = Wx + b \quad (1)$$

- **Hipótese**  $\hat{y}$  o valor que o modelo prevê,
- **Parâmetros**  $W$  que determinam como cada característica afeta a previsão, e  $b$  que controla o deslocamento fixo da previsão [Goodfellow et al., 2016].

## Parâmetros

- **Parâmetros** são valores que controlam o comportamento do sistema.
  - Objetivo de **encontrar** o “melhor” conjunto de parâmetros possível  $W$  e  $b$ , para descrever os dados.
- Primeiro precisamos definir o **erro/custo**.

## Função de custo

- Definir uma **medida de desempenho** do modelo.
- Uma maneira de medir o desempenho do modelo é calcular o **erro quadrático médio** (MSE) no conjunto de teste.
  - Na **estatística**, o MSE mede a média dos quadrados dos erros ou desvios [Goodfellow et al., 2016].

$$L = MSE = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \quad (2)$$

- **Diferença** entre a hipótese  $\hat{y}$  a partir do valor real  $y$ .

## Problema de otimização

- Precisamos **minimizar o custo** da hipótese  $\hat{y}$ , em função dos parâmetros do modelo  $W$  e  $b$ :

$$\min_{W,b} L = \min_{W,b} \frac{1}{N} \sum_i (y_i - \hat{y}_i)^2 \quad (3)$$

## Possíveis soluções

- Analítica
- Numérica:
  - **Algoritmos de otimização** que iteram sobre o conjunto de dados.

## Gradiente descendente

- Em **cálculo**, é um algoritmo de otimização iterativa de primeira ordem para encontrar o **mínimo de uma função** [Goodfellow et al., 2016].
  - **Gradiente** é uma operação que assume uma função de múltiplas variáveis, e retorna um **vetor na direção da inclinação** máxima no gráfico da função original.

$$\nabla L = \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix} \quad (4)$$

- Se quisermos **descer**, o que temos a fazer é **andar na direção oposta** ao gradiente.
- Esta seria a **estratégia para minimizar** as funções de custo.



## Taxa de Aprendizagem $\epsilon$

- Escalar positivo que determina o **tamanho do passo** [Goodfellow et al., 2016].

$$\begin{bmatrix} W' \\ b' \end{bmatrix} = \begin{bmatrix} W \\ b \end{bmatrix} - \epsilon \nabla L \quad (5)$$

$$= \begin{bmatrix} W \\ b \end{bmatrix} - \epsilon \begin{bmatrix} \frac{\partial L}{\partial W} \\ \frac{\partial L}{\partial b} \end{bmatrix} \quad (6)$$

- Poderia ser um problema de duas maneiras:
  - Se o passo é **muito pequeno**, nos moveremos lentamente para o mínimo,
  - Se o passo é **muito grande**, podemos acabar pulando além do mínimo.

## Back-propagation

Método para **calcular** o gradiente.

Enquanto **outro** algoritmo, como gradiente descendente, é usado para realizar a **aprendizagem** (otimização) usando esse gradiente [Goodfellow et al., 2016].

Duas fases:

- Forward-propagation
- Back-propagation

## Forward-propagation<sup>17</sup>

Processo de alimentação dos valores.

A entrada  $x$  fornece a informação inicial que se **propaga**, e finalmente produz  $\hat{y}$  [Goodfellow et al., 2016].

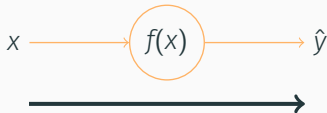


Figura 4: Forward-propagation.

---

<sup>17</sup>Everything you need know about neural networks — <https://hackernoon.com/>

## Back-propagation<sup>18</sup>

Usamos o valor da função de custo para calcular o erro. O valor do erro é propagado para atrás afim de calcular o gradiente com relação aos pesos.

Os gradientes são calculados usando a regra da cadeia.

$$\frac{dL}{dx} = \frac{dL}{dz} \frac{dz}{dx} \leftarrow df \leftarrow \frac{dL}{dz}$$

Figura 5: Back-propagation.

---

<sup>18</sup>Everything you need know about neural networks — <https://hackernoon.com/>

## Algoritmo Back-propagation<sup>19</sup>

- Após Forward-propagation, obtemos um valor de saída que é o valor previsto.
- Usamos uma função de perda  $L$  para calcular o valor do erro.
- Calculamos o gradiente do erro em relação a cada peso,
- Subtraímos o valor do gradiente do valor do peso.

Desta forma, nos aproximamos ao mínimo local.

---

<sup>19</sup>Everything you need know about neural networks — <https://hackernoon.com/>

## Problemas com gradiente descendente

O método de **gradiente descendente tradicional** calculará o gradiente de **todo o conjunto de dados**, mas executará apenas uma atualização.

Por isso pode ser muito lento e difícil de controlar para conjuntos de dados muito grandes e que não se encaixam na memória.

## Otimizadores<sup>20</sup>

Existem **extensões** (modificações ao método original) que tentam solucionar os problemas do gradiente descendente.

- Básico [Goodfellow et al., 2016]:
  - **SGD**: Stochastic Gradient Descent.
- Taxas de aprendizado adaptativas [Goodfellow et al., 2016]:
  - **AdaGrad** [Duchi et al., 2011].
  - **RMSprop** [Hinton, 2012].
  - **Adam** [Kingma and Ba, 2014].

---

<sup>20</sup>Types of Optimization Algorithms — <https://towardsdatascience.com/>

Regressão linear:  
Código GitHub



## Resumo parcial

- Regressão linear define a relação entre duas variáveis.
- Como **encontrar** a “melhor” combinação dos parâmetros?
  - **Função de custo:**
    - Erro quadrático médio
  - **Gradiente descendente**
    - Back-propagation para calcular o gradiente
    - Otimizadores (extensões do gradiente descendente)
- O processo de encontrar a melhor combinação de parâmetros (otimização) é chamado de **treinamento**.

## Regressão logística

- Na **estatística**, é um modelo de regressão em que a variável dependente é **categórica**.
  - Que a variável dependente seja categórica implica que estamos tratando um **problema de classificação**.
- A saída deve estar sempre entre 0 e 1 [Goodfellow et al., 2016].
  - A saída indica a probabilidade de uma entrada pertencer a uma classe.
- A **probabilidade** é modelada usando a relação

$$\hat{y} = \sigma(Wx + b) \quad (7)$$

- onde  $\sigma$  é chamada de **função de ativação** e limita a saída para o intervalo (0, 1).

## Função de ativação<sup>21</sup> $\sigma$

Deve ser uma função **sigmóide** para mapear qualquer número real da hipótese para o intervalo (0, 1) [Goodfellow et al., 2016]:

logistic

$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (8)$$

tanh

$$\sigma(z) = \frac{2}{1 + e^{-2z}} - 1 \quad (10)$$

softmax

$$\sigma(z)_i = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (9)$$

O nome sigmóide vem da forma em **S** do seu gráfico.

---

<sup>21</sup>Activation functions — <https://towardsdatascience.com/>

## Representação gráfica

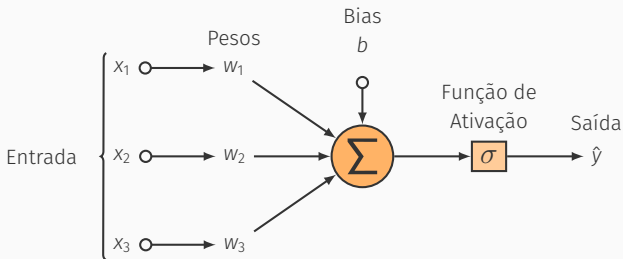


Figura 6: Representação gráfica da regressão logística  $\hat{y} = \sigma(Wx + b)$ .

## Função de custo<sup>22</sup>

- Para classificação pode-se usar *cross entropy cost* [Goodfellow et al., 2016]:

$$E(y, \hat{y}) = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad (11)$$

- **binary-cross-entropy**: problema de classificação binária,
- **categorical-cross-entropy**: problema de classificação multi-classe.

---

<sup>22</sup>Everything you need to know about Neural Networks — <https://hackernoon.com/>

# Regressão logística: Código GitHub

## Resumo parcial

**Tabela 1:** Regressão linear e regressão logística.

Método	Alvo	Hipótese	Custo
Regressão linear	Contínua	$\hat{y} = Wx + b$	MSE
Regressão logística	Categórica	$\hat{y} = \sigma(Wx + b)$	Cross Entropy

A maioria dos problemas do mundo real **não são linearmente separáveis**.

Para calcular **hipóteses não-lineares**, uma das formas computacionalmente mais eficientes é **conectar pequenas unidades** que fazem “regressão logística”<sup>23</sup>.

---

<sup>23</sup>Playing with machine learning — <https://medium.com/rocknnull>



## Redes Neurais Artificiais — RNA

São chamadas de redes porque elas são representadas pela composição de **várias funções** diferentes.

A informação **flui** através de funções conectadas em cadeia [Goodfellow et al., 2016]:

$$\hat{y} = f(x) \tag{12}$$

$$= f^{(2)}(f^{(1)}(x)) \tag{13}$$

- $x$ : Camada de entrada
- $f^{(1)}$ : Camada oculta
- $f^{(2)}$ : Camada de saída

## Camadas

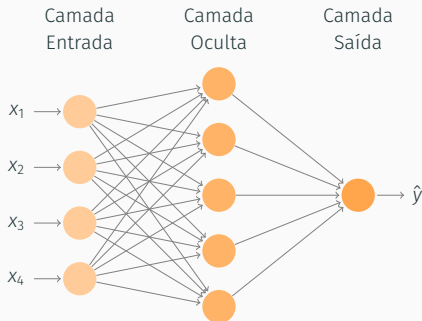
$$f^{(i)}(h) = \sigma_i(W_i h + b_i) \quad (14)$$

- $\sigma_i$ : função de ativação
- $W_i$ : Matriz de pesos
- $b_i$ : Vetor de bias

ou seja,

$$\hat{y} = \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) \quad (15)$$

## Representação gráfica



**Figura 7:** Representação gráfica de uma rede neural artificial com uma camada oculta  $\hat{y} = \sigma_2(W_2\sigma_1(W_1x + b_1) + b_2)$ .

# Redes neurais artificiais

## Funções de ativação<sup>24</sup>

Introduzem **não linearidade** nas redes neurais artificiais.

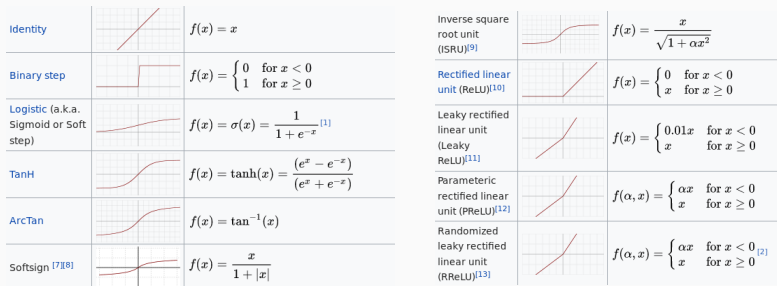


Figura 8: Algumas funções de ativação.

<sup>24</sup>Activation Functions: Neural Networks — <https://towardsdatascience.com/>

## Funções de ativação<sup>25</sup>

- Para problemas de **regressão**  $\sigma_2$ : função **identidade**.

$$\sigma(z) = z \quad (16)$$

- Para problemas de **classificação**  $\sigma_2$ : função **sigmóide**.
  - logistic,
  - tanh,
  - softmax,
  - etc.

---

<sup>25</sup>Activation Functions: Neural Networks — <https://towardsdatascience.com/>

## Função de custo

- Para problemas de **regressão**: minimização do **MSE**.

$$MSE = \frac{1}{n} \sum_i (y_i - \hat{y}_i)^2 \quad (17)$$

- Para problemas de **classificação**: minimização do **cross entropy**.

$$E(y, \hat{y}) = - \sum_{i=1}^K y_i \log(\hat{y}_i) \quad (18)$$

# Redes neurais artificiais: Código GitHub

## Resumo parcial

**Tabela 2:** Regressão linear, regressão logística e RNA.

Método	Alvo	Hipótese	Custo
R. Linear	Contínua	$\hat{y} = Wx + b$	MSE
R. Logística	Categórica	$\hat{y} = \sigma(Wx + b)$	Cross Entropy
RNA	Contínua, Categórica	$\hat{y} = \sigma_2(W_2\sigma_1(W_1x + b_1) + b_2)$	MSE, Cross Entropy



# Aprendizado Profundo

---

Rede neural  
artificial

+

Volume  
dados

+

Poder  
computacional



Aprendizado Profundo

## Aprendizado profundo

Abordagem do AM, que creceu em sua popularidade e utilidade, como resultado de **computadores poderosos, conjuntos de dados e técnicas maiores para treinar** redes mais profundas [Goodfellow et al., 2016].

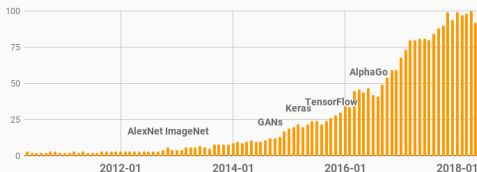


Figura 9: Interesse do aprendizado profundo nos ultimos 6 anos<sup>26</sup>.

---

<sup>26</sup><https://trends.google.com>

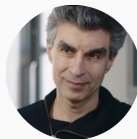
## Pesquisadores destacados



(a) Geoff Hinton  
Backpropag.



(b) Yann LeCun  
CNN.



(c) Yoshua Bengio  
RNA, GANs.



(d) Ian Goodfellow  
GANs.



(e) Andrew Ng  
Google Brain

**Figura 10:** Pesquisadores destacados de aprendizado profundo.

## Redes neurais artificiais profundas

Representadas pela composição de **várias camadas ocultas** [Goodfellow et al., 2016].

Pesquisadores concordam que o aprendizado profundo envolve mais que duas transformações da entrada para a saída (**CAP > 2**) [Schmidhuber, 2014].

$$\hat{y} = f(x) \tag{19}$$

$$= f^{(n)}(\dots f^{(2)}(f^{(1)}(x))) \tag{20}$$

$$= \sigma_n(W_n \dots \sigma_2(W_2 \sigma_1(W_1 x + b_1) + b_2) \dots + b_n) \tag{21}$$

- $f^{(1)}$ : Camada de entrada
- $f^{(i)} \forall i = [2, n - 1]$ : C. ocultas
- $f^{(n)}$ : Camada de saída
- $\sigma_i$ : função de ativação
- $W_i$ : Matriz de pesos
- $b_i$ : Vetor de bias

## Dificuldades no treinamento

Quanto mais profunda a RNA:

- Pode ser mais robusta,
- Porém, pode ser **mais difícil** de treiná-la.

## Problema do gradiente descendente<sup>27</sup>

Dada uma RNA Profunda com quatro camadas:

$$\hat{y} = f_4(f_3(f_2(f_1(x, z_1), z_2), z_3), z_4), \text{ sendo } z_n : \text{parâmetros camada } n \quad (22)$$

usamos **back-propagation** para atualizar,  $z_1 = z_1 - \epsilon \frac{\partial L}{\partial z_1}$ :

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial z_4} \frac{\partial z_4}{\partial z_3} \frac{\partial z_3}{\partial z_2} \frac{\partial z_2}{\partial z_1} \quad (23)$$

- Se os gradientes forem **maiores do que zero**, o produto **explode** para um número muito grande,
- Se os gradientes forem **menores do que zero**, o produto **desvanece** para um número muito próximo de zero.

---

<sup>27</sup><https://matheusfacure.github.io/2017/07/10/problemas-treinamento/>

## Função de ativação — ReLU<sup>28</sup>

Para evitar tais complicações e melhorar a precisão, pode ser usada a função de ativação baseadas em **ReLU** (*Rectified Linear Unit*) no treinamento da RNA Profunda.

$$\sigma(z) = \begin{cases} 0 & \text{para } z < 0 \\ z & \text{para } z \geq 0 \end{cases} \rightarrow \frac{d\sigma}{dz} = \begin{cases} 0 & \text{para } z < 0 \\ 1 & \text{para } z \geq 0 \end{cases} \quad (24)$$

As derivadas não vão mais desaparecer ou explodir, porque a derivada da função de ativação é limitada a os valores 0 e 1.<sup>29</sup>

---

<sup>28</sup>The vanishing gradient problem — <https://medium.com/@anishsingh20/>

<sup>29</sup>The vanishing gradient problem — <https://ayearofai.com/>



# Redes neurais artificiais profundas: Código GitHub

## Resumo parcial

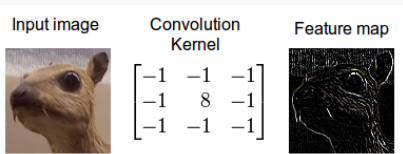
**Tabela 3:** Regressão linear, regressão logística, RNA e RNA profundas.

Método	Alvo	Hipótese	Custo
R. Linear	Continua	$\hat{y} = Wx + b$	MSE
R. Logística	Categórica	$\hat{y} = \sigma(Wx + b)$	Cross Entropy
RNA	Continua, Categórica	$\hat{y} = f^{(2)}(f^{(1)}(x))$	MSE, Cross Entropy
RNA Profundas	Continua, Categórica	$\hat{y} = f^{(n)}(\dots f^{(2)}(f^{(1)}(x)))$	MSE, Cross Entropy

## Redes neurais convolucionais — Redes convolucionais

Redes neurais que usam a **convolução** no lugar da multiplicação geral da matriz, em **pelo menos uma de suas camadas** [Goodfellow et al., 2016].

- Convolução é uma **operação matemática** que descreve uma regra de como misturar duas funções ou partes da informação<sup>30</sup>:
- Mapa de características  $I$ ,
- Kernel de convolução  $K$ , e
- Mapa de características transformadas  $S(i, j)$ .



$$S(i, j) = (I * K)(i, j) \quad (25)$$

<sup>30</sup><https://devblogs.nvidia.com/deep-learning-nutshell-core-concepts/>

## Parâmetros

As **camadas convolucionais** têm **parâmetros** que são aprendidos para que esses filtros sejam ajustados automaticamente para **extrair as informações** mais úteis para a tarefa em questão<sup>31</sup>.

- Entrada é uma matriz multidimensional de dados,
- Kernel é uma matriz multidimensional de parâmetros,
- Essas matrizes multidimensionais são **tensores** [Goodfellow et al., 2016].
  - Séries temporais: grade 1D intervalos de tempo regulares,
  - Dados de imagem: grade 2D de pixels.

---

<sup>31</sup><https://devblogs.nvidia.com/deep-learning-nutshell-core-concepts/>

# Redes neuronais convolucionais

## Camadas<sup>32</sup>

- **Convolution:** extrair características da imagem,
- **Pooling:** reduzir a dimensão da entrada, e
- **Dense/Fully connected:** conectar as camadas.

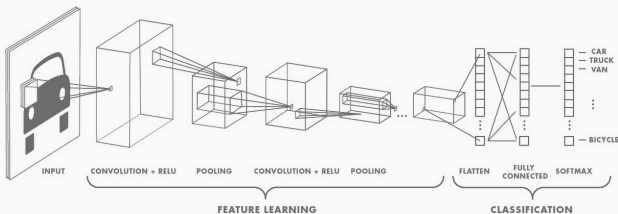


Figura 11: Exemplo de uma rede com varias camadas convolucionais<sup>33</sup>.

<sup>32</sup>A deeper understanding of CNNs — <https://towardsdatascience.com/>

<sup>33</sup><https://ch.mathworks.com/fr/discovery/convolutional-neural-network.html>

## Importância<sup>34</sup>

Ganharam popularidade através de **competições** como o **ImageNet** e, mais recentemente, eles são usados para **NLP** (Natural Language Processing) e **reconhecimento de fala** também.

---

<sup>34</sup>A deeper understanding of CNNs — <https://towardsdatascience.com/>

# Redes neuronais convolucionais: Código GitHub

## Resumo parcial

**Tabela 4:** Regressão linear, regressão logística, RNA e RNA profundas.

Método	Alvo	Hipótese	Custo
R. Linear	Continua	$\hat{y} = Wx + b$	MSE
R. Logística	Categórica	$\hat{y} = \sigma(Wx + b)$	Cross Entropy
RNA	Continua, Categórica	$\hat{y} = f^{(2)}(f^{(1)}(x))$	MSE, Cross Entropy
RNA Profundas	Continua, Categórica	$\hat{y} = f^{(n)}(\dots f^{(2)}(f^{(1)}(x)))$	MSE, Cross Entropy
RNA Convolucional	Continua, Categórica	$\hat{y} = f^{(n)}(\dots f_{(*)}^{(2)}(f_{(*)}^{(1)}(x)))$	MSE, Cross Entropy



# Revisão Final

---

## Resultados MNIST

- Épocas: 50
- Batch size: 256
- Otimizador: RMSProp
- Camada de saída: 10 unidades softmax

Modelo	Arquitetura	Ativação	Parâmetros	Precisão [%]
LOGREG	[]	[]	7.850	92.79
NN	[32]	[sigmoide]	25.450	96.27
DNN	[128, 64]	[relu, relu]	109.386	97.90
CDNN	[32*, 64*, 128]	[relu, relu, relu]	4.738.826	<b>98.84</b>

\*: tamanho do *kernel* convolucional ( $3 \times 3$ )

## Regularização<sup>35</sup>

Usado para superar o problema de *underfitting* e *overfitting*. Na regularização penalizamos a perda adicionando uma norma **L1** (LASSO) ou **L2** (Ridge) no vetor de peso  $W$ . Essas penalidades são incorporadas na **função de perda** que a rede otimiza<sup>36</sup>.

- L1: soma do valor absoluto dos coeficientes.
- L2: soma do valor ao quadrado dos coeficientes.
- **Dropout**: configura aleatoriamente uma fração de unidades da entrada para 0 a cada atualização durante o tempo de treinamento.

---

<sup>35</sup>Everything you need to know about Neural Networks — <https://hackernoon.com/>

<sup>36</sup><https://keras.io/regularizers/>

# Tópicos de interesse

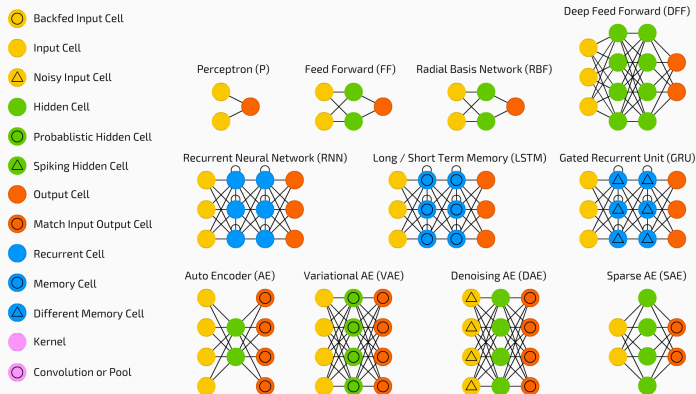


Figura 12: The neural network zoo<sup>37</sup>

<sup>37</sup><http://www.asimovinstitute.org/neural-network-zoo/>

# Tópicos de interesse

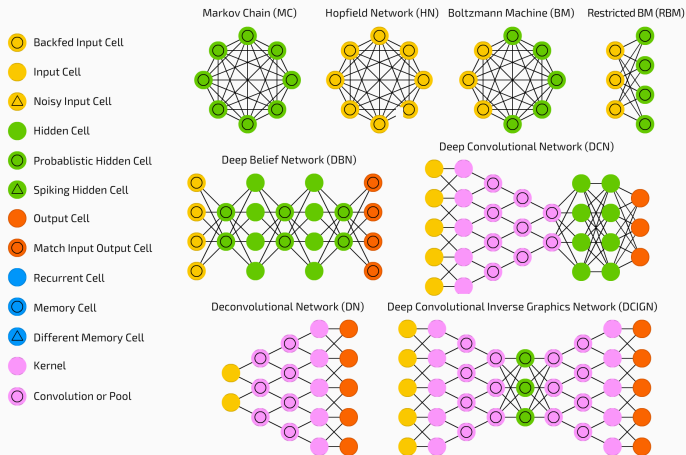


Figura 13: The neural network zoo<sup>38</sup>

<sup>38</sup><http://www.asimovinstitute.org/neural-network-zoo/>

# Tópicos de interesse

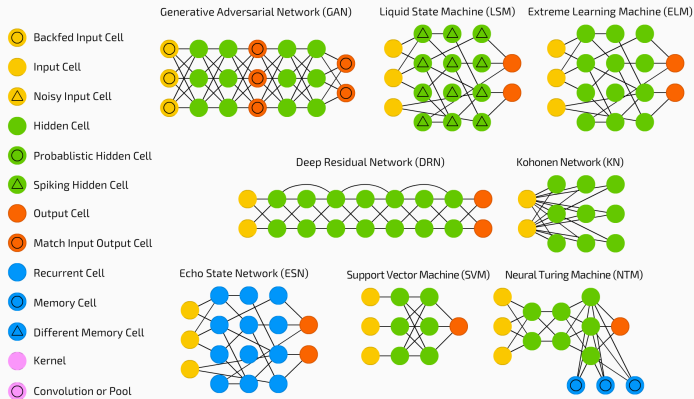


Figura 14: The neural network zoo<sup>39</sup>

<sup>39</sup><http://www.asimovinstitute.org/neural-network-zoo/>

Obrigada

## Códigos

[https://github.com/mafda/deep\\_learning\\_101](https://github.com/mafda/deep_learning_101)



# Referências i



Duchi, J., Hazan, E., and Singer, Y. (2011).

**Adaptive subgradient methods for online learning and stochastic optimization.**

*Journal of Machine Learning Research*, 12(Jul):2121–2159.



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

***Deep Learning.***

MIT Press.



Hinton, G. (2012).

**Neural networks for machine learning.**

Coursera, video lectures.



Kingma, D. P. and Ba, J. (2014).

**Adam: A method for stochastic optimization.**

*CoRR*, abs/1412.6980.



Schmidhuber, J. (2014).

**Deep learning in neural networks: An overview.**

*CoRR*, abs/1404.7828.