

Luiz Guilherme Moraes da Costa Faria

APRENDIZADO DE MÁQUINA

Brasília, DF
24 de setembro de 2025

Luiz Guilherme Moraes da Costa Faria

APRENDIZADO DE MÁQUINA

Universidade de Brasília

Orientador: Nome do Orientador/Revisor (se aplicável)

Brasília, DF
24 de setembro de 2025

Sumário

Sumário	3
I	HISTÓRIA DA IA E DO COMPUTADOR 11
1	UMA BREVE HISTÓRIA DO COMPUTADOR 13
1.1	A Necessidade de Contar ao Longo das Eras 13
1.1.1	Ábaco 13
1.1.2	Régua de Cálculo 13
1.1.3	Bastões de Napier 13
1.1.4	Pascalina 13
2	UMA BREVE HISTÓRIA DA INTELIGÊNCIA ARTIFICIAL . . . 15
2.1	Os Anos 1900 15
2.2	Os Anos 1910 15
2.3	Os Anos 1920 15
2.4	Os Anos 1930 15
2.5	Os Anos 1940 15
2.6	Os Anos 1950 15
2.7	Os Anos 1960 15
2.8	Os Anos 1970 15
2.9	Os Anos 1980 15
2.10	Os Anos 1990 15
2.11	Os Anos 2000 15
2.12	Atualidade 15
II	CONCEITOS MATEMÁTICOS 17
3	CÁLCULO PARA APRENDIZADO DE MÁQUINA 19
3.1	Funções: A Base do Cálculo 19
3.2	Derivadas Ordinárias 19
3.3	Integrais Simples 19
3.4	Derivadas Parciais 19
4	ÁLGEBRA LINEAR PARA APRENDIZADO DE MÁQUINA . . . 21
4.1	A Unidade Fundamental: Vetores e Espaços Vetoriais 21

4.2	Organizando Dados: Matrizes e Suas Operações	21
4.3	Tensores: A Estrutura de Dados do Deep Learning	21
4.4	Resolvendo Sistemas e Encontrando Propriedades: Autovalores e Autovetores	21
4.5	Decomposição de Matrizes (SVD e PCA)	21
5	PROBABILIDADE E ESTATÍSTICA PARA APRENDIZADO DE MÁQUINA	23
5.1	Medindo a Incerteza: Probabilidade Básica e Condicional	23
5.2	O Teorema de Bayes: Aprendendo com Evidências	23
5.3	Descrevendo os Dados: Estatística Descritiva: Média, mediana, variância, desvio padrão	23
5.4	Variáveis Aleatórias e Distribuições de Probabilidade	23
5.5	A Função de Máxima Verossimilhança (Maximum Likelihood Estimation - MLE)	23
III	PILARES DAS REDES NEURAIS	25
6	O ALGORITMO DA REPROPROPAGAÇÃO E OS OTIMIZADORES BASEADOS EM GRADIENTE	27
6.1	O Método do Gradiente Descendente	27
6.1.1	Exemplo Ilustrativo: Cadeia de Montanhas	27
6.1.2	O Método em Si	28
6.1.3	Implementação em Python	30
6.2	A Retropropagação: Aprendendo com os Erros	31
6.3	Otimizadores Baseados em Gradiente	32
6.3.1	Método do Gradiente Estocástico	32
6.3.1.1	Implementação em Python	32
6.3.2	Método do Gradiente com Momentum	32
6.3.2.1	Implementação em Python	32
6.3.3	Nesterov	32
6.3.3.1	Implementação em Python	32
6.3.4	AdaGrad	32
6.3.4.1	Implementação em Python	32
6.3.5	RMSPprop	32
6.3.5.1	Implementação em Python	32
6.3.6	Adam	32
6.3.6.1	Implementação em Python	32
6.3.7	Nadam	32

6.3.7.1	Implementação em Python	32
6.4	O Método de Newton: Indo Além do Gradiente	32
6.4.1	Implementação em Python	32
7	FUNÇÕES DE ATIVAÇÃO SIGMOIDAIS	33
7.1	Teoremas da Aproximação Universal	33
7.2	Exemplos Ilustrativo	33
7.3	A Sigmoide Logística	33
7.3.1	Implementação em Python	33
7.4	Tangente Hiperbólica	34
7.4.1	Implementação em Python	34
7.5	Softsign: Uma Sigmoidal Mais Barata	34
7.5.1	Implementação em Python	35
7.6	Hard Sigmoid e Hard Tanh: O Sacrifício da Suavidade em Prol do Desempenho	35
7.6.1	Implementação em Python	37
7.7	O Desaparecimento de Gradientes	38
7.8	Comparativo de Desempenho das Sigmoidais	38
8	FUNÇÕES DE ATIVAÇÃO RETIFICADORAS	39
8.1	Exemplo Ilustrativo	39
8.2	Rectified Linear Unit e Revolução Retificadora	39
8.2.1	Implementação em Python	39
8.3	Dying ReLUs Problem	40
8.4	Corrigindo o Dying ReLUs Problem: As Variantes com Vazamento	40
8.4.1	Leaky ReLU	40
8.4.1.1	Implementação em Python	40
8.4.2	Parametric ReLU	41
8.4.3	Randomized Leaky ReLU	41
8.5	Em Busca da Suavidade	41
8.5.1	Exponential Linear Unit	41
8.5.2	Scaled Exponential Linear Unit	42
8.5.3	Noisy ReLU	42
8.6	O Problema dos Gradientes Explosivos	42
8.7	Comparativo de Desempenho das Funções Retificadoras	42
9	FUNÇÕES DE ATIVAÇÃO MODERNAS E OUTRAS FUNÇÕES DE ATIVAÇÃO	43
10	FUNÇÕES DE PERDA PARA CLASSIFICAÇÃO BINÁRIA	45

10.1	A Intuição da Perda: Medindo o Erro do Modelo	45
10.2	Entropia Cruzada Binária (Binary Cross-Entropy): A função de perda padrão	45
10.3	Perda Hinge (Hinge Loss)	45
10.4	Comparativo Visual e Prático	45
11	FUNÇÕES DE PERDA PARA CLASSIFICAÇÃO MULTILABEL .	47
11.1	Softmax e a Distribuição de Probabilidades	47
11.2	Entropia Cruzada Categórica (Categorical Cross-Entropy)	47
11.3	Entropia Cruzada Categórica Esparsa (Sparse Categorical Cross- Entropy)	47
12	METAHEURÍSTICAS: OTIMIZANDO REDES NEURAI SEM O GRADIENTE	49
12.1	Algoritmos Evolutivos	49
12.2	Inteligência de Enxame	49
IV	APRENDIZADO DE MÁQUINA CLÁSSICO	51
13	TÉCNICAS DE REGRESSÃO	53
13.1	Exemplo Ilustrativo	53
13.2	Regressão Linear	53
13.2.1	Função de Custo MSE	53
13.2.2	Equação Normal	53
13.2.3	Implementação em Python	53
13.3	Regressão Polinomial	53
13.3.1	Implementação em Python	53
13.4	Regressão de Ridge	53
13.4.1	Implementação em Python	53
13.5	Regressão de Lasso	53
13.5.1	Implementação em Python	53
13.6	Elastic Net	53
13.6.1	Implementação em Python	53
13.7	Regressão Logística	53
13.7.1	Implementação em Python	53
13.8	Regressão Softmax	53
13.8.1	Implementação em Python	53
13.9	Outras Técnicas de Regressão	53
14	ÁRVORES DE DECISÃO E FLORESTAS ALEATÓRIAS	55

14.1	Exemplo Ilustrativo	55
14.2	Entendendo o Conceito de Árvores	55
14.2.1	Árvores Binárias	55
14.3	Árvores de Decisão	55
14.3.1	Implementação em Python	55
14.4	Florestas Aleatórias	55
14.4.1	Implementação em Python	55
15	MÁQUINAS DE VETORES DE SUPORTE	57
15.1	Exemplo Ilustrativo	57
16	ENSAMBLE	59
16.1	Exemplo Ilustrativo	59
17	DIMENSIONALIDADE	61
17.1	Exemplo Ilustrativo	61
17.2	A Maldição da Dimensionalidade	61
17.3	Seleção de Características (Feature Selection)	61
17.4	Extração de Características (Feature Extraction)	61
17.4.1	Análise de Componentes Principais (PCA)	61
17.4.2	t-SNE (t-Distributed Stochastic Neighbor Embedding) e UMAP	61
18	CLUSTERIZAÇÃO	63
18.1	Exemplo Ilustrativo	63
18.2	Aprendizado Não Supervisionado: Encontrando Grupos nos Dados	63
18.3	Clusterização Particional: K-Means	63
18.4	Clusterização Hierárquica	63
18.5	Clusterização Baseada em Densidade: DBSCAN	63
V	REDES NEURAIIS PROFUNDAS (DNNS)	65
19	PERCEPTRONS MLP - REDES NEURAIIS ARTIFICIAIS	67
20	REDES FEEDFORWARD (FFNS)	69
21	REDES DE CRENÇA PROFUNDA (DBNS) E MÁQUINAS DE BOLTZMANN RESTRITAS	71
22	REDES NEURAIIS CONVOLUCIONAIS (CNN)	74
22.1	Exemplo Ilustrativo	74
22.2	Camadas Convolucionais: O Bloco Fundamental para as CNNs . .	74

22.2.1	Implementação em Python	74
22.3	Camadas de Pooling: Reduzindo a Dimensionalidade	76
22.3.1	Max Pooling	76
22.3.1.1	Implementação em Python	76
22.3.2	Average Pooling	78
22.3.2.1	Implementação em Python	78
22.3.3	Global Average Pooling	79
22.3.3.1	Implementação em Python	79
22.4	Camada Flatten: Achatando os Dados	80
22.4.1	Implementação em Python	80
22.5	Criando uma CNN	81
22.6	Detecção de Objetos	81
22.7	Redes Totalmente Convolucionais (FCNs)	81
22.8	You Only Look Once (YOLO)	81
22.9	Algumas Arquiteturas de CNNs	81
22.9.1	LeNet-5	81
22.9.2	AlexNet	81
22.9.3	GoogLeNet	81
22.9.4	VGGNet	81
22.9.5	ResNet	81
22.9.6	Xception	81
22.9.7	SENet	81
23	REDES RESIDUAIS (RESNETS)	83
24	REDES NEURAIS RECORRENTES (RNN)	85
24.1	Exemplo Ilustrativo	85
24.2	Neurônios e Células Recorrentes	85
24.2.1	Implementação em Python	85
24.3	Células de Memória	85
24.3.1	Implementação em Python	85
24.4	Criando uma RNN	85
24.5	O Problema da Memória de Curto Prazo	85
24.5.1	Células LSTM	85
24.5.2	Conexões Peephole	85
24.5.3	Células GRU	85
25	TÉCNICAS PARA MELHORAR O DESEMPENHO DE REDES NEURAIS	87
25.1	Técnicas de Inicialização	87

25.2	Regulização L1 e L2	87
25.3	Normalização	87
25.3.1	Normalização de Camadas	87
25.3.2	Normalização de Batch	87
25.4	Clipping do Gradiente	87
25.5	Dropout: Menos Neurônios Mais Aprendizado	87
25.6	Data Augmentation	87
26	TRANSFORMERS	89
26.1	As Limitações das RNNs: O Gargalo Sequencial	89
26.2	A Ideia Central: Self-Attention (Query, Key, Value)	89
26.3	Escalando a Atenção: Multi-Head Attention	89
26.4	A Arquitetura Completa: O Bloco Transformer	89
26.5	Entendendo a Posição: Codificação Posicional	89
26.6	As Três Grandes Arquiteturas	89
26.6.1	Encoder-Only (Ex: BERT): Para tarefas de entendimento	89
26.6.2	Decoder-Only (Ex: GPT): Para tarefas de geração	89
26.6.3	Encoder-Decoder (Ex: T5): Para tarefas de tradução/sumarização	89
26.7	Além do Texto: Vision Transformers (ViT)	89
27	REDES ADVERSÁRIAS GENERATIVAS (GANS)	91
28	MIXTURE OF EXPERTS (MOE)	93
29	MODELOS DE DIFUSÃO	95
30	REDES NEURAIS DE GRAFOS (GNNS)	97
VI	APÊNDICES	99
	Referências	101

Parte I

História da IA e do Computador

1 Uma Breve História do Computador

1.1 A Necessidade de Contar ao Longo das Eras

1.1.1 Ábaco

1.1.2 Régua de Cálculo

1.1.3 Bastões de Napier

1.1.4 Pascalina

2 Uma Breve História da Inteligência Artificial

2.1 Os Anos 1900

2.2 Os Anos 1910

2.3 Os Anos 1920

2.4 Os Anos 1930

2.5 Os Anos 1940

2.6 Os Anos 1950

2.7 Os Anos 1960

2.8 Os Anos 1970

2.9 Os Anos 1980

2.10 Os Anos 1990

2.11 Os Anos 2000

2.12 Atualidade

Parte II

Conceitos Matemáticos

3 Cálculo para Aprendizado de Máquina

3.1 Funções: A Base do Cálculo

3.2 Derivadas Ordinárias

3.3 Integrais Simples

3.4 Derivadas Parciais

4 Álgebra Linear para Aprendizado de Máquina

4.1 A Unidade Fundamental: Vetores e Espaços Vetoriais

4.2 Organizando Dados: Matrizes e Suas Operações

4.3 Tensores: A Estrutura de Dados do Deep Learning

4.4 Resolvendo Sistemas e Encontrando Propriedades: Autovalores e Autovetores

4.5 Decomposição de Matrizes (SVD e PCA)

5 Probabilidade e Estatística para Aprendizado de Máquina

5.1 Medindo a Incerteza: Probabilidade Básica e Condicional

5.2 O Teorema de Bayes: Aprendendo com Evidências

5.3 Descrevendo os Dados: Estatística Descritiva: Média, mediana, variância, desvio padrão

5.4 Variáveis Aleatórias e Distribuições de Probabilidade

5.5 A Função de Máxima Verossimilhança (Maximum Likelihood Estimation - MLE)

Parte III

Pilares das Redes Neurais

6 O Algoritmo da Repropagação e Os Otimizadores Baseados em Gradiente

6.1 O Método do Gradiente Descendente

O Método do Gradiente faz parte de uma série de métodos numéricos que possuem como função otimizar diferentes funções. Métodos dessa forma veem sendo estudados a séculos, um exemplo disso é o trabalho *Méthode générale pour la résolution des systèmes d'équations simultanées* (Método geral para resolução de sistemas de equações simultâneas em português) do matemático francês do século XVIII Cauchy (1847), em que o autor apresenta um método que pode ser considerado um precursor para o método do gradiente atual.

Nesse texto, o autor apresenta uma forma de minimizar uma função de múltiplas variáveis ($u = f(x, y, z)$) que não assume valores negativos, para fazer isso, ele faz uso do cálculo de derivadas parciais dessa função de cada um dos seus componentes ($D_x u, D_y u, D_z u$), em seguida, ele realiza um passo de atualização, de forma que os valores de cada uma das variáveis sejam ligeiramente incrementados por valores (α, β, γ) (Cauchy, 1847). Um ponto importante destacado por Cauchy (1847) é de que esses incrementos devem ser proporcionais ao negativo das suas respectivas derivadas parciais, ele descreve que esse processo de calcular as derivadas e fazer pequenos incrementos deve ser feito de forma iterativa, assim, calculá-se as derivadas, faz-se os incrementos, e o passo é repetido até convergir para o valor mínimo de u .

Esse trabalho explica bem como aplicar o método do gradiente para se calcular mínimos de funções, mas para facilitar o entendimento do leitor, em seguida está um exemplo ilustrativo explicando o funcionamento dessa ferramenta.

6.1.1 Exemplo Ilustrativo: Cadeia de Montanhas

Imagine que você adora aventuras, e por isso, decidiu fazer uma trilha em uma floresta que fica em uma cadeia de montanhas que podem ser escaladas. Então, você teve a incrível ideia de ir para o menor ponto dessa cadeia de montanhas, pois, no guia que você estava seguindo, falava que lá havia um lago com uma água cristalina, perfeito para tirar fotos.

Para chegar até esse lago, você conta com uma bússula um tanto quanto diferente, ao invés dela apontar para o norte como uma bússola comum, ela aponta para a direção do lugar com menor altitude de uma região. Isso é perfeito para o que você precisa, pois ela irá apontar justamente para o lago de você quer ir.

Com isso em mente, você criou um plano de como irá chegar a esse lago, ele é método que segue dois passos diferentes, sendo eles:

1. Olha na bússola qual a direção ela está apontando;
2. Anda um metro na direção apontada pela bússola.

Você chegou na conclusão que se seguir essa estratégia repetidas vezes, em algum momento, você inevitavelmente irá chegar no lago que está querendo tirar as suas fotos.

Na matemática, existe um método semelhante a este, que busca com base em uma bússola (chamada de vetor gradiente), encontrar um ponto de mínimo de um determinado lugar (neste caso, uma função composta por múltiplas variáveis). Esse é o método do gradiente, ele é ponto central desse capítulo, pois, ele (e suas variações) junto com o algoritmo da retropropagação são algumas das principais ferramentas que colaboram para que os modelos de aprendizado de máquina possam aprender com os seus erros e com isso se tornarem melhores a cada iteração.

6.1.2 O Método em Si

A vantagem do método do gradiente é que ele é uma ferramenta matemática, e por isso pode ser representado utilizando notações mais formais e de forma enxuta. As notações utilizadas por Cauchy são diferentes das que são utilizadas hoje em dia, mas o seu significado não muda. Em *Deep Learning*, Goodfellow, Bengio e Courville (2016) explicam essa ferramenta através da equação 6.1 que deve ser repetida por múltiplos passos até o modelo convergir, ou seja, encontrar o ponto de mínimo da função estudada.

Método do Gradiente Descendente

$$x' = x - \epsilon \nabla f(x) \quad (6.1)$$

Em que:

- x' : representa as coordenadas do próximo ponto;
- x : representa as coordenadas do ponto atual;
- ϵ : representa o tamanho do passo, também chamado de taxa de aprendizado;
- $\nabla f(x)$: representa o vetor gradiente calculado na posição do ponto atual (x) para função que se deseja otimizar.

Note que assim como no método proposto por Cauchy, é pego como base o inverso do vetor gradiente. Isso ocorre pois o vetor gradiente é um vetor especial que tem como principal propriedade apontar para a direção de maior crescimento de uma função no ponto que está sendo calculada. Mas no método, o objetivo não é encontrar o ponto que irá gerar os maiores valores da função, e sim o contrário. Por isso, é tomado inverso do vetor gradiente, que, dessa forma, estará então apontando para a direção de menor crescimento de uma função.

Um ponto a ser destacado nesse método é na hora de escolher uma taxa de aprendizado para ser utilizada no método. Uma taxa de aprendizado muito pequena significa que o passo que o modelo irá dar de um ponto para outro será menor, e com isso implica que ele levará mais passos para encontrar um ponto de mínimo. É como se você fosse comparar a quantidade de passos que você gasta para andar do seu quarto até a sua cozinha com a quantidade de passos dados por uma formiga até lá, ambos vão chegar no local, mas a formiga certamente irá demorar bem mais. Considerando isso, surge então a hipótese de que quanto maior for o passo, mais rápido será a convergência, mas isso também não funciona muito bem, pois um passo muito largo pode ultrapassar o ponto de mínimo indo parar em outro canto da função, e ficará tentando chegar até o mínimo mas não irá conseguir pois caminha uma distância muito grande de uma só vez. Essas duas situações, em que o passo é pequeno demais e que o passo é grande demais, são ilustradas na figura 1.

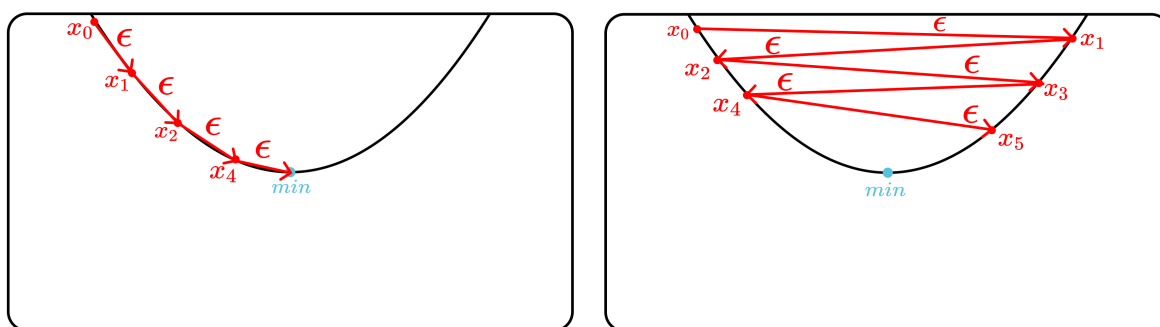


Figura 1 – Comparativo do tamanho de passos em uma função polinomial.

Na prática, escolher o valor da taxa de aprendizado é uma tarefa que irá depender de modelo em modelo, também irá variar com os diferentes métodos de otimização além da topologia da rede neural que está sendo construída. É sempre recomendado então experimentar diferentes tamanhos de passo, de forma que seja encontrado um que melhor se ajusta ao cenário que está sendo trabalhado.

Outro ponto que deve-se atentar é com relação as funções que estão sendo analisadas ao utilizar o método do gradiente mas também qualquer otimizador que seja baseado nele. Se tivermos uma função convexa, em que seu formato lembra um funil, será bem mais fácil para o modelo encontrar o ponto de mínimo global daquela função.

Mas se tivermos uma função não convexa, cheia de ondas e com muitos pontos de mínimos locais e pontos de sela, a convergência do modelo será pior, pois existe a chance de que ele fique preso em um ponto de mínimo local ou em um ponto de sela. Isso afeta diretamente o desempenho da rede neural que estará sendo criada, fazendo com que ela tenha métricas piores. O problema é que muitas das vezes a função $f(x)$ que estaremos interessados para calcular o desempenho do modelo será não convexa, dificultando o seu aprendizado.

Na figura 2 é possível ver o gráfico de duas funções diferentes, a primeira sendo uma função convexa e a segunda uma função não convexa.

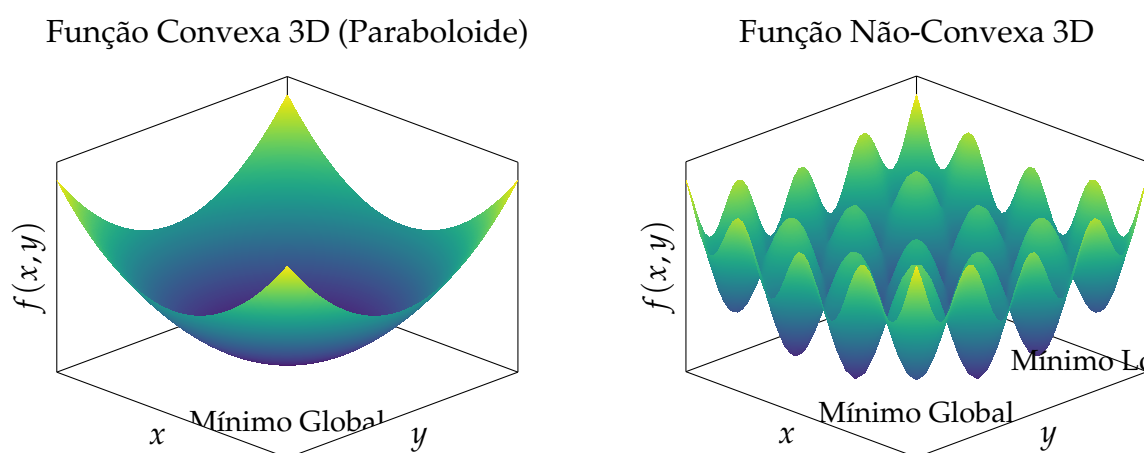


Figura 2 – Comparação entre funções 3D convexas e não-convexas.

6.1.3 Implementação em Python

Para implementar o método do gradiente utilizando Python e a biblioteca de cálculos Numpy, deve-se seguir como base a equação 6.1, criando uma classe implenta essa ferramenta, recebendo como parâmetros de entrada a taxa de aprendizado (`learning_rate`), a função que se quer encontrar o ponto de mínimo (`function`), e um ponto inicial (`initial_point`), que pode ser um conjunto de coordenadas aleatórias ou escolhidas pelo programador.

Outro ponto que deve ser destacado nas informações dessa função é de que ela também irá receber a derivada da função (`function_prime`) que se quer descobrir o ponto de mínimo, pois não será feito cálculo simbólico para calcular o vetor gradiente. Dessa forma os cálculos serão mais rápidos. Além disso, também é preciso definir outros parâmetros auxiliares, como a quantidade máxima de iterações que o modelo irá seguir (`tolerance`), que são chamadas de épocas `epochs`, também será definida um grau de tolerância para o modelo, pois, podem existir casos em que a norma do vetor gradiente são será exatamente zero, mas um valor muito próximo de zero, assim a

tolerância será responsável por definir qual valor dessa diferença será aceitável para o problema.

Por fim, um ponto interessante que é possível adicionar nessa função é uma lista, que irá armazenar todos os pontos que o modelo passou em cada uma de suas iterações, indicando o seu caminho pela função que está sendo estudada.

Bloco de Código : Classe completa do otimizador GradientDescent

```
1 class GradientDescent:
2
3     def __init__(self, function, function_prime, initial_point
4         , learning_rate=0.001, epochs=100, tolerance=1e-6):
5         self.f = function
6         self.fp = function_prime
7         self.ip = initial_point
8         self.lr = learning_rate
9         self.ep = epochs
10        self.tol = tolerance
11        self.path = []
12
13    def update_step(self):
14        for i in range(self.ep):
15            self.path.append(self.ip)
16            grad = self.fp(self.ip)
17            if abs(grad) < self.tol: break
18            self.ip = self.ip + self.lr * (-grad)
19        return self.ip, self.path
```

Note que a classe apresenta apenas dois métodos, o primeiro sendo o `__init__`, que inicializa os parâmetros da classe, e a `update_step` a qual é responsável por implementar de fato o método do gradiente, ela irá retornar o ponto de mínimo e também a lista com os pontos pelos quais o modelo passou ao longo das iterações.

6.2 A Retropropagação: Aprendendo com os Erros

Ainda no contexto de utilizar com o vetor gradiente para otimizar um modelo de rede neural, existe uma ferramenta que trabalha justamente com esse processo, ela é a retropropagação ou *backpropagation* em inglês.

Definição: A **retropropagação** é uma ferramenta que veio para permitir que redes que fazem o uso de unidades de neurônios possam aprender, para isso, o procedimento ajusta repetidamente os pesos das conexões da rede para minimizar a diferença entre o valor atual da saída do vetor da rede neural com o valor real desejado (Rumelhart; Hinton; Williams, 1986).

Essa ferramenta foi introduzida para a comunidade científica pelos pesquisadores Rumelhart, Hinton e Williams (1986) no texto *Learning Representations by Back-Propagating Errors*,

6.3 Otimizadores Baseados em Gradiente

6.3.1 Método do Gradiente Estocástico

6.3.1.1 Implementação em Python

6.3.2 Método do Gradiente com Momentum

6.3.2.1 Implementação em Python

6.3.3 Nesterov

6.3.3.1 Implementação em Python

6.3.4 AdaGrad

6.3.4.1 Implementação em Python

6.3.5 RMSProp

6.3.5.1 Implementação em Python

6.3.6 Adam

6.3.6.1 Implementação em Python

6.3.7 Nadam

6.3.7.1 Implementação em Python

6.4 O Método de Newton: Indo Além do Gradiente

6.4.1 Implementação em Python

7 Funções de Ativação Sigmoidais

7.1 Teoremas da Aproximação Universal

7.2 Exemplos Ilustrativo

7.3 A Sigmoide Logística

Sigmoide Logística

$$\sigma(z_i) = \frac{1}{1 + e^{-z_i}} \quad (7.1)$$

Derivada da sigmoide logística

$$\frac{d}{dz_i} \sigma(z_i) = \frac{e^{-z_i}}{(1 + e^{-z_i})^2} \quad (7.2)$$

7.3.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Sigmoid

```

1 import numpy as np
2
3 class Sigmoid(Layer):
4     def __init__(self):
5         super().__init__()
6         self.input = None
7         self.sigmoid = None
8
9     def forward(self, input_data):
10         self.input = input_data
11         self.sigmoid_output = 1 / (1 + np.exp(-input_data))
12         return self.sigmoid_output
13
14     def backward(self, grad_output):
15         sigmoid_grad = self.sigmoid_output * (1 - self.
sigmoid_output)
16         return grad_output * sigmoid_grad, None

```

7.4 Tangente Hiperbólica

Tangente Hiperbólica

$$\tanh(z_i) = \frac{\sinh(z_i)}{\cosh(z_i)} = \frac{e_i^z - e^{-z_i}}{e_i^z + e^{-z_i}} \quad (7.3)$$

Derivada da tangente hiperbólica

$$\frac{d}{dz_i} \tanh(z_i) = \text{sech}^2(z_i) \quad (7.4)$$

7.4.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Tangente Hiperbólica

```

1 import numpy as np
2 from layers.base import Layer # Assuming your base class is
   here
3
4 class Tanh(Layer):
5     def __init__(self):
6         super().__init__()
7         self.input = None
8         self.tanh_output = None
9
10    def forward(self, input_data):
11        self.input = input_data
12        self.tanh_output = np.tanh(self.input)
13        return self.tanh_output
14
15    def backward(self, grad_output):
16        tanh_grad = 1 - self.tanh_output**2
17        return grad_output * tanh_grad, None

```

7.5 Softsign: Uma Sigmoidal Mais Barata

Softsign

$$\text{softsign}(z_i) = \frac{z_i}{1 + |z_i|} \quad (7.5)$$

Derivada da softsign

$$\frac{d}{dz_i} \text{softsign}(z_i) = \frac{1}{(1 + |z_i|)^2} \quad (7.6)$$

7.5.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Softsign

```

1 from layers.base import Layer
2 import numpy as np
3
4 class Softsign(Layer):
5     def __init__(self):
6         super().__init__()
7         self.input = None
8
9     def forward(self, input_data):
10        self.input = input_data
11        return self.input / (1 + np.abs(self.input))
12
13    def backward(self, grad_output):
14        grad = (1 / (1 + np.abs(self.input)))**2
15        return grad_output * softsign_grad, None

```

7.6 Hard Sigmoid e Hard Tanh: O Sacrifício da Suavidade em Prol do Desempenho

Hard Sigmoid

$$\text{hard sigmoid}(z_i) = \begin{cases} 0 & \text{se } z_i < -3 \\ z_i/6 + 0.5 & \text{se } -3 \leq z_i \leq 3 \\ 1 & \text{se } z_i > 3 \end{cases} \quad (7.7)$$

Derivada da Hard Sigmoid

$$\frac{d}{dz_i} \text{hard sigmoid}(z_i) = \begin{cases} 0 & \text{se } z_i < -3 \\ 1/6 & \text{se } -3 < z_i < 3 \\ 0 & \text{se } z_i > 3 \end{cases} \quad (7.8)$$

Hard Tanh

$$\text{hard tanh}(z_i) = \begin{cases} -1 & \text{se } z_i < -1 \\ z_i & \text{se } -1 \leq z_i \leq 1 \\ 1 & \text{se } z_i > 1 \end{cases} \quad (7.9)$$

Derivada da Hard Tanh

$$\frac{d}{dz_i} \text{hard tanh}(z_i) = \begin{cases} 0 & \text{se } z_i < -1 \\ 1 & \text{se } -1 < z_i < 1 \\ 0 & \text{se } z_i > 1 \end{cases} \quad (7.10)$$

7.6.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Hard Sigmoid

```
1 from layers.base import Layer
2 import numpy as np
3
4 class HardSigmoid(Layer):
5
6     def __init__(self):
7         super().__init__()
8         self.input = None
9
10    def forward(self, input_data):
11        self.input = input_data
12
13        output = self.input / 6 + 0.5
14        output = np.clip(output, 0, 1)  # A more concise way
15        to handle the bounds
16
17        return output
18
19    def backward(self, grad_output):
20        hard_sigmoid_grad = np.full_like(self.input, 1 / 6)
21
22        hard_sigmoid_grad[self.input < -3] = 0
23        hard_sigmoid_grad[self.input > 3] = 0
24
25        return grad_output * hard_sigmoid_grad, None
```

Bloco de Código : Classe completa do função de ativação Hard Sigmoid

```
1 from layers.base import Layer
2 import numpy as np
3
4
5 class HardTanh(Layer):
6     def __init__(self):
7         super().__init__()
8         self.input = None
9
10    def forward(self, input_data):
11        self.input = input_data
12        return np.clip(self.input, -1, 1)
13
14    def backward(self, grad_output):
15
16        hard_tanh_grad = np.where((self.input > -1) & (self.
input < 1), 1, 0)
17
18        return grad_output * hard_tanh_grad, None
```

7.7 O Desaparecimento de Gradientes

7.8 Comparativo de Desempenho das Sigmoidais

8 Funções de Ativação Retificadoras

8.1 Exemplo Ilustrativo

8.2 Rectified Linear Unit e Revolução Retificadora

Rectified Linear Unit (ReLU)

$$\text{ReLU}(z_i) = \begin{cases} z_i, & \text{se } z_i > 0 \\ 0, & \text{se } z_i \leq 0 \end{cases} \quad (8.1)$$

Derivada Rectified Linear Unit (ReLU)

$$\frac{d}{dz_i}[\text{ReLU}](z_i) = \begin{cases} 1, & \text{se } z_i > 0 \\ 0, & \text{se } z_i \leq 0 \end{cases} \quad (8.2)$$

8.2.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Rectified Linear Unit

```

1 import numpy as np
2 from layers.base import Layer
3
4 class ReLU(Layer):
5     def __init__(self):
6         super().__init__()
7         self.input = None
8
9     def forward(self, input_data):
10        self.input = input_data
11        return np.maximum(0, self.input)
12
13    def backward(self, grad_output):
14        relu_grad = (self.input > 0)
15
16        # Apply the chain rule
17        return grad_output * relu_grad, None

```

8.3 Dying ReLUs Problem

8.4 Corrigindo o Dying ReLUs Problem: As Variantes com Vazamento

8.4.1 Leaky ReLU

Leaky ReLU (LReLU)

$$\text{LReLU}(z_i) = \begin{cases} z_i, & \text{se } z_i \geq 0 \\ \alpha \cdot z_i, & \text{se } z_i < 0 \end{cases} \quad (8.3)$$

Derivada Leaky ReLU (LReLU)

$$\frac{d}{dz_i}[\text{LReLU}](z_i) = \begin{cases} 1, & \text{se } z_i > 0 \\ \alpha, & \text{se } z_i \leq 0 \end{cases} \quad (8.4)$$

8.4.1.1 Implementação em Python

Bloco de Código : Classe completa do função de ativação Leaky ReLU

```

1 import numpy as np
2 from layers.base import Layer
3
4
5 class LeakyReLU(Layer):
6     def __init__(self, alpha=0.01):
7         super().__init__()
8         self.input = None
9         self.alpha = alpha
10
11     def forward(self, input_data):
12         self.input = input_data
13         return np.maximum(self.input * self.alpha, self.input)
14
15     def backward(self, grad_output):
16         leaky_relu_grad = np.where(self.input > 0, 1, self.alpha)
17         return grad_output * leaky_relu_grad, None

```

8.4.2 Parametric ReLU

Parametric ReLU (PReLU)

$$\text{PReLU}(z_i) = \begin{cases} z_i, & \text{se } z_i \geq 0 \\ \alpha_i \cdot z_i, & \text{se } z_i < 0 \end{cases} \quad (8.5)$$

Derivada Parametric ReLU (PReLU)

$$\frac{d}{dz_i}[\text{PReLU}](z_i) = \begin{cases} 1, & \text{se } z_i > 0 \\ \alpha_i, & \text{se } z_i < 0 \\ \nexists, & \text{se } z_i = 0 \end{cases} \quad (8.6)$$

8.4.3 Randomized Leaky ReLU

Randomized Leaky ReLU (RRReLU)

$$\text{RRReLU}(z_i) = \begin{cases} z_i, & \text{se } z_i > 0 \\ \alpha_i z_i, & \text{se } z_i \leq 0 \end{cases} \quad (8.7)$$

Derivada Randomized Leaky ReLU (RRReLU)

$$\frac{d}{dz_i}[\text{RRReLU}](z_i) = \begin{cases} 1, & \text{se } z_i > 0 \\ \alpha_i, & \text{se } z_i \leq 0 \end{cases} \quad (8.8)$$

8.5 Em Busca da Suavidade

8.5.1 Exponential Linear Unit

Exponential Linear Unit (ELU)

$$\text{ELU}(z_i) = \begin{cases} z_i, & \text{se } z_i \geq 0 \\ \alpha \cdot (e^{z_i} - 1), & \text{se } z_i < 0 \end{cases} \quad (8.9)$$

Derivada Exponential Linear Unit (ELU)

$$\frac{d}{dz_i}[\text{ELU}](z_i) = \begin{cases} 1, & \text{se } z_i > 0 \\ \alpha \cdot e^{z_i}, & \text{se } z_i \leq 0 \end{cases} \quad (8.10)$$

8.5.2 Scaled Exponential Linear Unit

Scaled Exponential Linear Unit (ELU)

$$\text{SELU}(z_i) = \lambda \begin{cases} z_i, & \text{se } z_i > 0 \\ \alpha \cdot (e^{z_i} - 1), & \text{se } z_i \leq 0 \end{cases} \quad (8.11)$$

Derivada Scaled Exponential Linear Unit (ELU)

$$\frac{d}{dz_i}[\text{SELU}](z_i) = \lambda \begin{cases} 1, & \text{se } z_i > 0 \\ \alpha \cdot e^{z_i}, & \text{se } z_i \leq 0 \end{cases} \quad (8.12)$$

8.5.3 Noisy ReLU

Noisy ReLU (NReLU)

$$\text{NReLU}(z_i) = \begin{cases} 0 & \text{se } z_i \leq 0 \\ z_i + \mathcal{N}(0, \sigma(z_i)) & \text{se } z_i > 0 \end{cases} \quad (8.13)$$

Derivada Noisy ReLU (NReLU)

$$\frac{d}{dz_i}[\text{NReLU}](z_i) = \begin{cases} 0 & \text{se } z_i \leq 0 \\ 1 & \text{se } z_i > 0 \end{cases} \quad (8.14)$$

8.6 O Problema dos Gradientes Explosivos

8.7 Comparativo de Desempenho das Funções Retificadoras

9 Funções de Ativação Modernas e Outras Funções de Ativação

O texto do seu capítulo começa aqui...

10 Funções de Perda para Classificação Binária

10.1 A Intuição da Perda: Medindo o Erro do Modelo

10.2 Entropia Cruzada Binária (Binary Cross-Entropy): A função de perda padrão

10.3 Perda Hinge (Hinge Loss)

10.4 Comparativo Visual e Prático

11 Funções de Perda para Classificação Multilabel

11.1 Softmax e a Distribuição de Probabilidades

11.2 Entropia Cruzada Categórica (Categorical Cross-Entropy)

11.3 Entropia Cruzada Categórica Esparsa (Sparse Categorical Cross-Entropy)

12 Metaheurísticas: Otimizando Redes Neurais Sem o Gradiente

O texto do seu capítulo começa aqui...

12.1 Algoritmos Evolutivos

12.2 Inteligência de Enxame

Parte IV

Aprendizado de Máquina Clássico

13 Técnicas de Regressão

13.1 Exemplo Ilustrativo

13.2 Regressão Linear

13.2.1 Função de Custo MSE

13.2.2 Equação Normal

13.2.3 Implementação em Python

13.3 Regressão Polinomial

13.3.1 Implementação em Python

13.4 Regressão de Ridge

13.4.1 Implementação em Python

13.5 Regressão de Lasso

13.5.1 Implementação em Python

13.6 Elastic Net

13.6.1 Implementação em Python

13.7 Regressão Logística

13.7.1 Implementação em Python

13.8 Regressão Softmax

13.8.1 Implementação em Python

13.9 Outras Técnicas de Regressão

14 Árvores de Decisão e Florestas Aleatórias

14.1 Exemplo Ilustrativo

14.2 Entendendo o Conceito de Árvores

14.2.1 Árvores Binárias

14.3 Árvores de Decisão

14.3.1 Implementação em Python

14.4 Florestas Aleatórias

14.4.1 Implementação em Python

15 Máquinas de Vetores de Suporte

15.1 Exemplo Ilustrativo

16 Ensamble

16.1 Exemplo Ilustrativo

17 Dimensionalidade

17.1 Exemplo Ilustrativo

17.2 A Maldição da Dimensionalidade

17.3 Seleção de Características (Feature Selection)

17.4 Extração de Características (Feature Extraction)

17.4.1 Análise de Componentes Principais (PCA)

17.4.2 t-SNE (t-Distributed Stochastic Neighbor Embedding) e UMAP

18 Clusterização

18.1 Exemplo Ilustrativo

18.2 Aprendizado Não Supervisionado: Encontrando Grupos nos Dados

18.3 Clusterização Particional: K-Means

18.4 Clusterização Hierárquica

18.5 Clusterização Baseada em Densidade: DBSCAN

Parte V

Redes Neurais Profundas (DNNs)

19 Perceptrons MLP - Redes Neurais Artificiais

O texto do seu capítulo começa aqui...

20 Redes FeedForward (FFNs)

O texto do seu capítulo começa aqui...

21 Redes de Crença Profunda (DBNs) e Máquinas de Boltzmann Restritas

O texto do seu capítulo começa aqui...

22 Redes Neurais Convolucionais (CNN)

22.1 Exemplo Ilustrativo

22.2 Camadas Convolucionais: O Bloco Fundamental para as CNNs

22.2.1 Implementação em Python

Bloco de Código : Classe completa de Convolution2D

```
1 import numpy as np
2 from layers.base import Layer
3
4 class Convolution2D(Layer):
5     def __init__(self, input_channels, num_filters,
6         kernel_size, stride=1, padding=0):
7         super().__init__()
8         self.input_channels = input_channels
9         self.num_filters = num_filters
10        self.kernel_size = kernel_size
11        self.stride = (stride, stride) if isinstance(stride,
12            int) else stride
13        self.padding = padding
14
15        kernel_height, kernel_width = self.kernel_size
16        self.kernels = np.random.randn(num_filters,
17            input_channels, kernel_height, kernel_width) * 0.01
18        self.biases = np.zeros((num_filters, 1))
19        self.params = [self.kernels, self.biases]
20        self.cache = None
21
22    def forward(self, input_data):
23        (batch_size, input_height, input_width, input_channels
24        ) = input_data.shape
25        filters, _, kernel_height, kernel_width = self.kernels
26        .shape
27        stride_height, stride_width = self.stride
28
29        pad_config = ((0, 0), (self.padding, self.padding), (
30            self.padding, self.padding), (0, 0))
31        input_padded = np.pad(input_data, pad_config, mode='
32            constant')
33        self.cache = input_padded
```


22.3 Camadas de Pooling: Reduzindo a Dimensionalidade

22.3.1 Max Pooling

22.3.1.1 Implementação em Python

Bloco de Código : Classe completa de MaxPooling2D

```

1 import numpy as np
2 from layers.base import Layer
3
4 class MaxPooling2D(Layer):
5     def __init__(self, pool_size=(2,2), stride=None):
6         super().__init__()
7         self.pool_size = pool_size
8         self.stride = stride if stride is not None else
pool_size
9         self.cache = None
10
11     def forward(self, input_data):
12         (batches, input_height, input_width, channels) =
input_data.shape
13
14         pool_height, pool_width = self.pool_size
15         stride_height, stride_width = self.stride
16
17         output_height = int((input_height - pool_height) /
stride_height) + 1
18         output_width = int((input_width - pool_width) /
stride_width) + 1
19
20         output_matrix = np.zeros((batches, output_height,
output_width, channels))
21         self.cache = np.zeros_like(input_data)
22
23         for b in range(batches):
24             for c in range(channels):
25                 for h in range(output_height):
26                     for w in range(output_width):
27
28                         start_height = h * stride_height
29                         start_width = w * stride_width
30                         end_height = start_height +
pool_height
31                         end_width = start_width + pool_width
32

```


22.3.2 Average Pooling

22.3.2.1 Implementação em Python

Bloco de Código : Classe completa de AveragePooling2D

```
1 import numpy as np
2 from layers.base import Layer
3
4 class AveragePooling2D(Layer):
5     def __init__(self, pool_size=(2, 2), stride=None):
6         super().__init__()
7         self.pool_size = pool_size
8         self.stride = stride if stride is not None else
pool_size
9
10    def forward(self, input_data):
11        (batches, input_height, input_width, channels) =
input_data.shape
12        pool_h, pool_w = self.pool_size
13        stride_h, stride_w = self.stride
14
15        output_height = int((input_height - pool_h) / stride_h
+ 1
16        output_width = int((input_width - pool_w) / stride_w)
+ 1
17
18        output_matrix = np.zeros((batches, output_height,
output_width, channels))
19
20        for b in range(batches):
21            for c in range(channels):
22                for h in range(output_height):
23                    for w in range(output_width):
24                        start_h = h * stride_h
25                        start_w = w * stride_w
26                        end_h = start_h + pool_h
27                        end_w = start_w + pool_w
28
29                        pooling_window = input_data[b, start_h
:end_h, start_w:end_w, c]
30                        output_matrix[b, h, w, c] = np.mean(
pooling_window) # Changed to mean
31
32        return output_matrix
33
```


22.3.3 Global Average Pooling

22.3.3.1 Implementação em Python

Bloco de Código : Classe completa de GlobalAveragePooling2D

```
1 import numpy as np
2 from layers.base import Layer
3
4 class GlobalAveragePooling2D(Layer):
5     def __init__(self):
6         super().__init__()
7         self.input_shape = None
8
9     def forward(self, input_data):
10         self.input_shape = input_data.shape
11
12         output = np.mean(input_data, axis=(1, 2), keepdims=
13             True)
14
15         return output
16
17     def backward(self, output_gradient):
18         _, input_h, input_w, _ = self.input_shape
19
20         distributed_grad = output_gradient / (input_h *
21             input_w)
22
23         upsampled_grad = np.ones(self.input_shape) *
24             distributed_grad
25
26         return upsampled_grad, None
```

22.4 Camada Flatten: Achatando os Dados

22.4.1 Implementação em Python

Bloco de Código : Classe completa de Flatten

```
1 import numpy as np
2 from layers.base import Layer
3
4 class Flatten(Layer):
5     def __init__(self):
6         super().__init__()
7         self.input_shape = None
8
9     def forward(self, input_data):
10         self.input_shape = input_data.shape
11
12         flatten_output = input_data.reshape(input_data.shape
13 [0], -1)
14
15         return flatten_output
16
17     def backward(self, output_gradient):
18         input_gradient = output_gradient.reshape(self.
19 input_shape)
20         return input_gradient, None
```

22.5 Criando uma CNN

22.6 Detecção de Objetos

22.7 Redes Totalmente Convolucionais (FCNs)

22.8 You Only Look Once (YOLO)

22.9 Algumas Arquiteturas de CNNs

22.9.1 LeNet-5

22.9.2 AlexNet

22.9.3 GoogLeNet

22.9.4 VGGNet

22.9.5 ResNet

22.9.6 Xception

22.9.7 SENet

23 Redes Residuais (ResNets)

O texto do seu capítulo começa aqui...

24 Redes Neurais Recorrentes (RNN)

O texto do seu capítulo começa aqui...

24.1 Exemplo Ilustrativo

24.2 Neurônios e Células Recorrentes

24.2.1 Implementação em Python

24.3 Células de Memória

24.3.1 Implementação em Python

24.4 Criando uma RNN

24.5 O Problema da Memória de Curto Prazo

24.5.1 Células LSTM

24.5.2 Conexões Peephole

24.5.3 Células GRU

25 Técnicas para Melhorar o Desempenho de Redes Neurais

25.1 Técnicas de Inicialização

25.2 Regularização L1 e L2

25.3 Normalização

25.3.1 Normalização de Camadas

25.3.2 Normalização de Batch

25.4 Clipping do Gradiente

25.5 Dropout: Menos Neurônios Mais Aprendizado

25.6 Data Augmentation

26 Transformers

26.1 As Limitações das RNNs: O Gargalo Sequencial

26.2 A Ideia Central: Self-Attention (Query, Key, Value)

26.3 Escalando a Atenção: Multi-Head Attention

26.4 A Arquitetura Completa: O Bloco Transformer

26.5 Entendendo a Posição: Codificação Posicional

26.6 As Três Grandes Arquiteturas

26.6.1 Encoder-Only (Ex: BERT): Para tarefas de entendimento

26.6.2 Decoder-Only (Ex: GPT): Para tarefas de geração

26.6.3 Encoder-Decoder (Ex: T5): Para tarefas de tradução/sumarização

26.7 Além do Texto: Vision Transformers (ViT)

27 Redes Adversárias Generativas (GANs)

O texto do seu capítulo começa aqui...

28 Mixture of Experts (MoE)

O texto do seu capítulo começa aqui...

29 Modelos de Difusão

O texto do seu capítulo começa aqui...

30 Redes Neurais de Grafos (GNNs)

O texto do seu capítulo começa aqui...

Parte VI

Apêndices

Referências

CAUCHY, Augustin-Louis. Méthode générale pour la résolution des systèmes d'équations simultanées. *Comptes Rendus Hebdomadaires des Séances de l'Académie des Sciences*, v. 25, p. 536–538, 1847. Citado na p. 27.

GOODFELLOW, Ian; BENGIO, Yoshua; COURVILLE, Aaron. *Deep Learning*. [S. l.]: MIT Press, 2016. Citado na p. 28.

RUMELHART, David E.; HINTON, Geoffrey E.; WILLIAMS, Ronald J. Learning Representations by Back-Propagating Errors. *Nature*, v. 323, p. 533–536, 1986. Citado na p. 32.