



# UNIVERSIDADE FEDERAL DO RIO GRANDE DO NORTE DISCIPLINA REDES NEURAIS - DEEP LEARNING Iª LISTA DE EXERCÍCIO - 2024.1

Aluno: Vinícius Lamas von Sohsten  
Matrícula: 20231011310  
Professor: Adrião Duarte Doria Neto

## 1. Métricas de distância

A representação do conhecimento em uma rede neural é definida pelos valores dos pesos sinápticos da rede, sendo a representação do conhecimento a chave para seu desempenho. Apesar deste tema ser complicado, existem quatro regras para representação do conhecimento. A primeira regra diz respeito à similaridade entre entrada e representação da rede, para que assim uma informação possa ser classificada. Existem diversas medidas para determinar a similaridade das informações, a seguir, apresenta-se um estudo sobre algumas dessas métricas.

**I. Distância Euclidiana:** A distância Euclidiana entre dois vetores  $\mathbf{x}_i$  e  $\mathbf{x}_j$  em um espaço Euclidiano é a distância do segmento de reta que os une, definida pela seguinte fórmula:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \|\mathbf{x}_i - \mathbf{x}_j\| = \left[ \sum_{k=1}^m (x_{ik} - x_{jk})^2 \right]^{\frac{1}{2}} \quad (1)$$

onde  $m$  é a dimensão do espaço euclidiano  $\mathbb{R}^m$ , e  $x_{ik}$  e  $x_{jk}$  são os  $k$ -ésimos elementos dos vetores de entrada.

Essa métrica é adequada para problemas de classificação de padrão, onde a forma ou magnitude das diferenças entre padrões são importantes, por exemplo, reconhecimento de dígitos. Também pode ser utilizada em algoritmos de clustering como K-means.

**II. Distância de Minkowski:** É uma generalização da distância Euclidiana que permite ajustar o parâmetro de ordem da fórmula, como mostra-se a seguir:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \left[ \sum_{k=1}^m |x_{ik} - x_{jk}|^p \right]^{\frac{1}{p}} \quad (2)$$

onde  $p$ , o parâmetro a ser ajustado, é um inteiro e  $p \geq 1$ , pois quando  $p < 1$  a métrica viola a desigualdade triangular tornando a métrica inválida. Quando  $p$  assume o valor de 1 ou 2, a métrica corresponde à distância de Manhattan e a distância Euclidiana, respectivamente.

Essa métrica é bem semelhante ao caso da distância Euclidiana, sendo útil para classificar padrões, mas permitindo ajustar a sensibilidade às diferenças de magnitude das entradas. Também utiliza-se em clustering quando se deseja dar mais peso a certas características a outras.

**III. Distância City Block:** Também conhecida como distância de Manhattan, é a soma das diferenças absolutas das coordenadas entre dois pontos, um caso particular da distância de Minkowski, apresenta-se da seguinte maneira:

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sum_{k=1}^m |x_{ik} - x_{jk}| \quad (3)$$

Aplicada em casos de classificação de padrões quando as característica têm diferentes escalas, pois a distância não é afetada por transformações lineares. Também em caso de clustering para agrupamento de dados esparsos em alta dimensão.

IV. **Distância de Mahalanobis:** Essa medida leva em consideração a covariância e os valores médios entre duas populações de dados. Descrita pela seguinte fórmula:

$$d_{ij}^2 = (\mathbf{x}_i - \boldsymbol{\mu}_i) \boldsymbol{\Sigma}^{-1} (\mathbf{x}_j - \boldsymbol{\mu}_j) \quad (4)$$

onde  $\boldsymbol{\Sigma}^{-1}$  é a inversa da matriz de covariância  $\boldsymbol{\Sigma}$ , e  $\boldsymbol{\mu}_i$  e  $\boldsymbol{\mu}_j$  representam os valores médios dos vetores  $\mathbf{x}_i$  e  $\mathbf{x}_j$ , respectivamente, de tal maneira:

$$\boldsymbol{\mu}_i = E[\mathbf{x}_i] \quad (5)$$

onde  $E$  é o operador estatístico esperado.

Essa métrica se faz útil em problemas de classificação de padrões onde as variâncias das diferentes dimensões dos dados são importantes e não são iguais. Também pode ser aplicada em algoritmos de clustering como DBSCAN.

V. **Coefficiente de Correlação de Pearson:** Este coeficiente mede o grau de correlação linear entre duas variáveis, sendo denotado como  $\rho$  quando utilizado como parâmetro populacional e representado por  $r$  para uma estatística amostral. Se faz útil quando ambas as variáveis em estudo têm distribuição normal, pois caso não possuam tal distribuição o coeficiente pode ser afetado por valores extremos tornando a análise pelo coeficiente inadequada. O valor deste coeficiente assume apenas valores entre -1 e 1, para valores positivos tem-se a correlação positiva e o mesmo vale para valores negativos. Porém, no caso em que o coeficiente assume valor zero significa que o resultado é inconclusivo e deve-se procurar outros meios de investigar a correlação, pois pode existir uma dependência não linear. O coeficiente pode assumir as seguintes fórmulas:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{[\sum_{i=1}^n (x_i - \bar{x})^2][\sum_{i=1}^n (y_i - \bar{y})^2]}} \quad (6)$$

$$\rho_{ab} = \frac{E[(Y_a - E[Y_a])(Y_b - E[Y_b])]}{\sigma_a \sigma_b} \quad (7)$$

onde  $x_i$  e  $y_i$  são os valores que estão contidos nos vetores  $\mathbf{x}$  e  $\mathbf{y}$ , respectivamente, que se deseja medir a correlação, os valores médios dos vetores são representados por  $\bar{x}$  e  $\bar{y}$ , respectivamente. Para a outra equação, vale salientar que o numerador representa a covariância dos dados e no denominador se encontram os desvios padrões dos respectivos dados representados pelo símbolo  $\sigma$ , cujo o valor é a raiz quadrada da variância.

Esta métrica se mostra eficiente em casos de reconhecimento de padrões para encontrar a similaridade entre dois conjuntos de dados multidimensionais. Pelo mesmo motivo, para análise de dados ajuda a identificar padrões de correlação entre diferentes variáveis.

VI. **Similaridade Cosseno:** Esse é um caso simples utilizado para medir o cosseno do ângulo entre dois vetores no espaço, fazendo o produto interno dos vetores e dividindo pela norma dos vetores, da seguinte maneira:

$$\cos(\mathbf{x}_i, \mathbf{x}_j) = \frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|} \quad (8)$$

Útil para calcular a similaridade entre vetores de características em sistemas de recomendação e em sistemas de busca para medir a similaridades entre documentos e consultas de pesquisa.

## 2. Demonstrações com base na expansão em série de Taylor

### a) Método do Gradiente da Descida mais Íngreme

Partindo da expressão em série de Taylor da função custo:

$$\mathcal{E}(\mathbf{w}(n) + \Delta \mathbf{w}(n)) = \mathcal{E}(\mathbf{w}(n)) + \mathbf{g}^t(\mathbf{w}(n)) \Delta \mathbf{w}(n) + \frac{1}{2} \Delta \mathbf{w}^t(n) \mathbf{H}(\mathbf{w}(n)) \Delta \mathbf{w}(n) + O(\|\Delta \mathbf{w}\|^3) \quad (9)$$

Nesse método é realizado o truncamento da expansão de Taylor na primeira ordem em torno de  $\mathbf{w}(n)$  para aproximar o valor da função custo. Dessa forma os termos de segunda ordem e maiores são desprezados, obtém-se:

$$\mathbf{w}(n) + \Delta \mathbf{w}(n) = \mathbf{w}(n + 1) \quad (10)$$

$$\mathcal{E}(\mathbf{w}(n + 1)) \simeq \mathcal{E}(\mathbf{w}(n)) + \mathbf{g}^t(\mathbf{w}(n)) \Delta \mathbf{w}(n) \quad (11)$$

O algoritmo da descida mais íngreme é descrito formalmente por:

$$\Delta \mathbf{w}(n) = -\eta \mathbf{g}(\mathbf{w}(n)) \Rightarrow \mathbf{w}(n+1) = \mathbf{w}(n) - \eta \mathbf{g}(\mathbf{w}(n)) \quad (12)$$

onde  $\eta$  é uma constante positiva chamada taxa de aprendizagem. Pode-se então substituir Eq. 12 na Eq. 11 e utilizar a Eq. 10 para obter:

$$\mathcal{E}(\mathbf{w}(n+1)) \leq \mathcal{E}(\mathbf{w}(n)) \Rightarrow \mathcal{E}(\mathbf{w}(n+1)) \simeq \mathcal{E}(\mathbf{w}(n)) - \eta \|\mathbf{g}(n)\|^2 \quad (13)$$

#### b) Método de Newton

Diferentemente do método do gradiente, o método de Newton faz uso da expansão de Taylor de segunda ordem da função de custo em torno do ponto  $\mathbf{w}(n)$ . Desta forma, obtém-se:

$$\mathcal{E}(\mathbf{w}(n+1)) - \mathcal{E}(\mathbf{w}(n)) = \Delta \mathcal{E}(\mathbf{w}(n)) \simeq \mathbf{g}^t(\mathbf{w}(n)) \Delta \mathbf{w}(n) + \frac{1}{2} \Delta \mathbf{w}^t(n) \mathbf{H}(\mathbf{w}(n)) \Delta \mathbf{w}(n) \quad (14)$$

Para minimizar a variação da função custo  $\Delta \mathcal{E}(\mathbf{w})$  deve-se diferenciar a Eq. 14 em relação a  $\Delta \mathbf{w}$ . Obtém-se, então:

$$\mathbf{g}(\mathbf{w}(n)) + \mathbf{H}(\mathbf{w}(n)) \Delta \mathbf{w}(n) = \mathbf{0} \quad (15)$$

Agora, basta isolar  $\Delta \mathbf{w}(n)$  e utilizar a relação da Eq. 10 para obter:

$$\mathbf{w}(n+1) = \mathbf{w}(n) + \mathbf{H}^{-1}(\mathbf{w}(n)) \mathbf{g}(\mathbf{w}(n)) \quad (16)$$

#### c) Passo ótimo

Pode-se utilizar a definição da Eq. 12 e substituir na Eq. 14, obtendo-se:

$$\mathcal{E}(\mathbf{w} - \eta \mathbf{g}(\mathbf{w})) \simeq \mathcal{E}(\mathbf{w}) - \eta \mathbf{g}^t(\mathbf{w}) \mathbf{g}(\mathbf{w}) + \frac{1}{2} \eta^2 \mathbf{g}^t(\mathbf{w}) \mathbf{H}(\mathbf{w}) \mathbf{g}(\mathbf{w}) \quad (17)$$

O passo ótimo ocorre quando:

$$\frac{d}{d\eta} (\mathcal{E}(\mathbf{w} - \eta \mathbf{g}(\mathbf{w}))) = 0 \quad (18)$$

Portanto, ao derivar Eq. 17 em relação a  $\eta$ , obtém-se a relação do passo ótimo ou taxa de aprendizagem ótima, para uma  $\mathbf{H}$  positiva, da seguinte maneira:

$$\eta^* = \frac{\mathbf{g}(\mathbf{w}) \mathbf{g}^t(\mathbf{w})}{\mathbf{g}^t(\mathbf{w}) \mathbf{H}(\mathbf{w}) \mathbf{g}(\mathbf{w})} \quad (19)$$

#### d) Valor do passo se o vetor gradiente estiver alinhado com o autovetor correspondente ao autovalor máximo

Se o vetor gradiente  $\mathbf{g}(\mathbf{w})$  estiver alinhado com o autovetor de  $\mathbf{H}(\mathbf{w})$  correspondente ao autovalor máximo  $\lambda_{\max}$ , tem-se as seguintes condições:

$$\mathbf{g}(\mathbf{w}) = \alpha \mathbf{v}_{\max} \quad (20)$$

$$\mathbf{H}(\mathbf{w}) \mathbf{v}_{\max} = \lambda_{\max} \mathbf{v}_{\max} \quad (21)$$

onde  $\alpha$  é um escalar constante e  $\mathbf{v}_{\max}$  é o autovetor máximo da relação. Com isso, é possível obter o passo sob tais condições substituindo as Eq. 20 e Eq. 21 na Eq. 19, em que obtém-se:

$$\eta^* = \frac{\alpha^2 \mathbf{v}_{\max}^2}{\alpha^2 \mathbf{v}_{\max} \mathbf{H} \mathbf{v}_{\max}} \quad (22)$$

$$\eta^* = \frac{\alpha^2 \mathbf{v}_{\max}^2}{\alpha^2 \lambda_{\max} \mathbf{v}_{\max}^2} \quad (23)$$

Portanto, para esse cenário, o passo ótimo é:

$$\eta^* = \frac{1}{\lambda_{\max}} \quad (24)$$

### 3. Matriz Hessiana

#### a) Demonstração

Essa demonstração pode ser realizada como mostrado nos seguintes passos:

$$\mathbf{H}(n-1) = \sum_{k=1}^{n-1} \mathbf{g}(k)\mathbf{g}^t(k) \quad (25)$$

$$\mathbf{H}(n) = \mathbf{g}(n)\mathbf{g}^t(n) + \sum_{k=1}^{n-1} \mathbf{g}(k)\mathbf{g}^t(k) \quad (26)$$

Substituindo Eq. 25 na Eq. 26, obtém-se o resultado desejado:

$$\mathbf{H}(n) = \mathbf{H}(n-1) + \mathbf{g}(n)\mathbf{g}^t(n) \quad (27)$$

#### b) Lema da inversão matricial

Partindo do resultado do lema da inversão matricial:

$$\mathbf{A}^{-1} = \mathbf{B} - \mathbf{B}\mathbf{C}(\mathbf{D} + \mathbf{C}^t\mathbf{B}\mathbf{C})^{-1}\mathbf{C}^t\mathbf{B} \quad (28)$$

Para obter-se a forma recursiva da inversa da hessiana, pode-se utilizar o resultado do lema da inversão matricial utilizando as seguintes relações:

$$\mathbf{A} = \mathbf{H}(n) \quad (29)$$

$$\mathbf{B}^{-1} = \mathbf{H}(n-1) \quad (30)$$

$$\mathbf{C} = \mathbf{g}(n) \quad (31)$$

$$\mathbf{D} = 1 \quad (32)$$

Portando, obtém-se nos seguintes passos, a demonstração desejada:

$$\mathbf{H}^{-1}(n) = \mathbf{H}^{-1}(n-1) - \mathbf{H}^{-1}(n-1)\mathbf{g}(n)(1 + \mathbf{g}(n)\mathbf{H}^{-1}(n-1)\mathbf{g}(n))^{-1}\mathbf{g}^t(n)\mathbf{H}^{-1}(n-1) \quad (33)$$

$$\mathbf{H}^{-1}(n) = \mathbf{H}^{-1}(n-1) - \frac{\mathbf{H}^{-1}(n-1)\mathbf{g}(n)\mathbf{g}^t(n)\mathbf{H}^{-1}(n-1)}{1 + \mathbf{g}(n)\mathbf{H}^{-1}(n-1)\mathbf{g}(n)} \quad (34)$$

### 4. Função de múltiplas variáveis

#### a) Representação na forma matricial

A função dada é:

$$E(w_1, w_2, w_3) = w_1^2 - w_2^2 + 2w_3^2 - 0.1w_1w_2 + 0.5w_1w_3 + 0.08w_2w_3 - 0.2w_1 + 0.3w_2 - 0.7w_3 + 0.4 \quad (35)$$

A forma matricial é dada por:

$$E(\mathbf{w}) = \mathbf{w}^t\mathbf{A}\mathbf{w} + \mathbf{b}^t\mathbf{w} + \alpha \quad (36)$$

onde  $\mathbf{w} = [w_1, w_2, w_3]^t$ , onde  $\alpha \in \mathbf{R}$ ,  $\mathbf{b} \in \mathbf{R}^3$  e  $\mathbf{A} \in \mathbf{R}^{3 \times 3}$ .

Além disso, para resolver esse problema, pode-se representar a matriz  $\mathbf{A}$  e o vetor  $\mathbf{b}$  das seguintes maneiras:

$$\mathbf{A} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix} \quad (37)$$

$$\mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix} \quad (38)$$

Realizando a substituição das Eq. 37 e Eq. 38 na Eq. 36 e realizando as devidas manipulações algébricas, obtém-se a seguinte igualdade polinomial:

$$E = A_{11}w_1^2 + A_{22}w_2^2 + A_{33}w_3^2 + (A_{12} + A_{21})w_1w_2 + (A_{13} + A_{31})w_1w_3 + (A_{23} + A_{32})w_2w_3 + b_1w_1 + b_2w_2 + b_3w_3 + \alpha$$

(39)

Dessa maneira, pode-se resolver e obter:

$$\mathbf{A} = \begin{bmatrix} 1 & -0.05 & 0.25 \\ -0.05 & -1 & 0.04 \\ 0.25 & 0.04 & 2 \end{bmatrix} \quad (40)$$

$$\mathbf{b} = \begin{bmatrix} -0.2 \\ 0.3 \\ -0.7 \end{bmatrix} \quad (41)$$

$$\alpha = 0.4 \quad (42)$$

Com isso, foi obtido a forma matricial da função.:

$$E(\mathbf{w}) = [w_1 \ w_2 \ w_3] \begin{bmatrix} 1 & -0.05 & 0.25 \\ -0.05 & -1 & 0.04 \\ 0.25 & 0.04 & 2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + [-0.2 \ 0.3 \ -0.7] \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + \alpha \quad (43)$$

#### b) Determinar o gradiente

Para determinar o gradiente da função dada, deve-se derivá-la parcialmente em relação a  $\mathbf{w}$ :

$$\mathbf{g}(\mathbf{w}) = \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = 2\mathbf{A}\mathbf{w} + \mathbf{b}^t \quad (44)$$

$$\mathbf{g}(\mathbf{w}) = 2 \begin{bmatrix} 1 & -0.05 & 0.25 \\ -0.05 & -1 & 0.04 \\ 0.25 & 0.04 & 2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} + [-0.2 \ 0.3 \ -0.7] \quad (45)$$

$$\mathbf{g}(\mathbf{w}) = [2.3w_1 - 0.2 \quad 1.29w_2 + 0.3 \quad 2.29w_3 - 0.7] \quad (46)$$

#### c) Determinar matriz Hessiana

Para determinar a matriz Hessiana da função dada é necessário derivar o seu gradiente parcialmente em relação a  $\mathbf{w}$ , ou seja, derivar duas vezes parcialmente a função custo em relação a  $\mathbf{w}$ . Obtém-se:

$$\mathbf{H}(\mathbf{w}) = \frac{\partial^2 E(\mathbf{w})}{\partial \mathbf{w}^2} = 2\mathbf{A} = \begin{bmatrix} 2 & -0.1 & 0.5 \\ -0.1 & -2 & 0.08 \\ 0.5 & 0.08 & 4 \end{bmatrix} \quad (47)$$

#### d) Verificar se a função é estritamente convexa

Uma função estritamente convexa é aquela em que sua matriz Hessiana é definida positiva, isto é, todos os autovalores da matriz são positivos. Para descobrir os autovalores de forma rápida e eficiente, foi utilizado o módulo "linalg" da biblioteca Numpy em um código Python, obtém-se:

$$\lambda = [4.12 \quad 1.88 \quad -2.00] \quad (48)$$

Portanto, de acordo com a definição, a função não é estritamente convexa.

#### e) Teste da segunda derivada para múltiplas dimensões

Em um ponto crítico, isto é, em que  $\mathbf{g}(\mathbf{w}) = \mathbf{0}$ , pode-se examinar os autovalores da matriz Hessiana para determinar a condição desse ponto, da seguinte maneira:

- I. Para o ponto ser um mínimo local, a matriz Hessiana deve ser definida positiva, todos seus autovalores devem ser positivos;
- II. para o ponto ser um máximo local, a matriz Hessiana deve ser definida negativa, todos seus autovalores devem ser negativos;
- III. para o caso de ser um ponto de sela, é necessário que ao menos um autovalor seja positivo e ao menos um autovalor seja negativo.

## 5. Estudo dos algoritmos de otimização

### I. Gradiente Estocástico Descendente

O Gradiente Estocástico Descendente (Stochastic Gradient Descent - SGD) é um algoritmo de otimização que se baseia no método da descida mais íngreme, porém calcula o gradiente usando apenas um pequeno conjunto de exemplos escolhidos aleatoriamente chamado *minibatch*. É um algoritmo muito eficiente em termos computacionais, principalmente para modelos de grande conjuntos de dados, no entanto o método de estimar o gradiente introduz fontes de ruídos que implica em menores taxas de convergência, ou seja, um aprendizado mais lento.

### II. AdaGrad

Algoritmo do Gradiente Adaptativo (Adaptative Gradient Algorithm - AdaGrad) é um algoritmo de otimização que individualmente adapta a taxa de aprendizagem para todos os parâmetros do modelo dividindo a taxa pela raiz da soma acumulada dos quadrados dos gradientes anteriores para cada parâmetro. Isso faz com que parâmetros com maiores derivadas parciais do custo tenham sua taxa de aprendizagem aumentada, enquanto parâmetros com menores derivadas parciais tenham sua taxa de aprendizagem diminuída relativamente. Isso torna o algoritmo eficaz em problemas em que gradientes possuem escalas muito diferentes e suaviza a convergência. Entretanto, na prática, acumular o quadrado dos gradientes pode resultar em diminuição excessiva da taxa de aprendizagem, fazendo o algoritmo performar bem alguns, mas não todos modelos de aprendizagem de máquina.

### III. RMSProp

Propagação da Raiz Quadrática Média (Root Mean Square Propagation - RMSProp) é um algoritmo que modifica o AdaGrad para performar melhor em situações não convexas, resolvendo o problema da diminuição muito rápida da taxa de aprendizagem ao longo do tempo. RMSProp faz isso utilizando uma média móvel exponencialmente ponderada dos gradientes anteriores, basicamente utiliza uma média decaindo exponencialmente para tirar a importância de valores passados extremos para que possa convergir mais rapidamente. Isto é possível, pois o algoritmo faz uso de um novo hiperparâmetro que controla o tamanho da escala da média móvel. Esse algoritmo mostra-se eficaz e prático para otimização de redes neurais profundas.

### IV. Adam

Estimativa de Momento Adaptativa (Adaptive Moment Estimation - Adam) é um algoritmo de otimização estocástico adaptativo que combina as ideias do SGD com momento (variação do SGD que inclui um parâmetro para acelerar a convergência) e RMSProp. Adam faz uso de dois momentos, um linear e um quadrático, para atualizar a taxa de aprendizagem com base numa relação desses momentos com suas respectivas taxas de decaimento exponencial, portanto esse algoritmo possui dois hiperparâmetros.

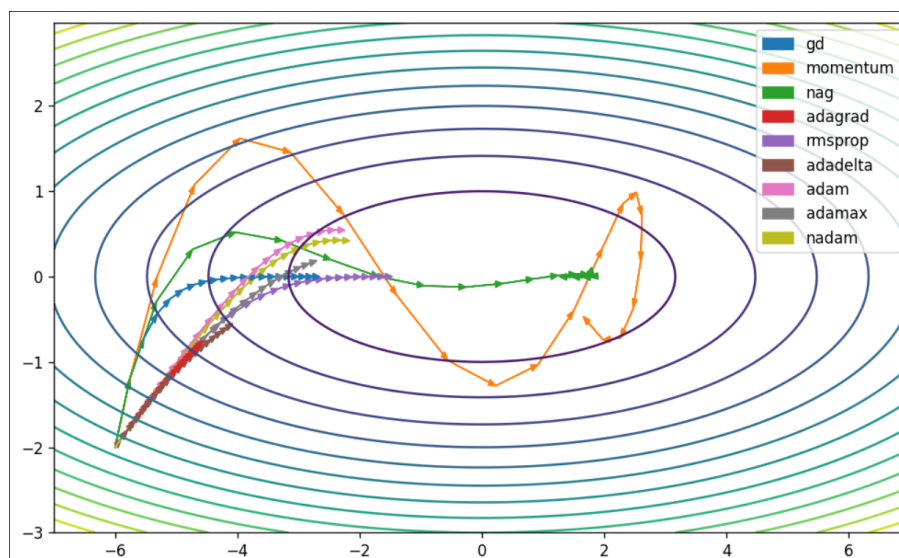


Figure 1. Comparativo dos algoritmos em um problema 2D.

## 6. Funções de ativação

### I. Análise comparativa das funções de ativação e suas respectivas derivadas

Cada uma dessas funções tem suas próprias características e vantagens, que podem influenciar significativamente o desempenho e o comportamento da rede neural. A seguir é realizada uma análise comparativa dos gráficos dessas funções, bem como de suas respectivas derivadas.

- **Função Sigmoid**

A função sigmoide possui uma forma de "S" suave e contínua, como mostrado na Fig. 2. Ela mapeia qualquer valor de  $v$  para o intervalo  $(0, 1)$ , o que a torna adequada para problemas de classificação binária, onde a saída desejada está entre 0 e 1. A derivada da função sigmoide atinge valores máximos em torno de  $v = 0$  e diminui para zero em direção aos extremos do intervalo. Isso significa que a taxa de mudança é alta em torno do ponto médio do intervalo e diminui à medida que nos aproximamos de 0 ou 1. Isso pode causar o problema de desaparecimento do gradiente, especialmente em redes profundas.

- **Função Tangsigmoide**

A função tangente hiperbólica é semelhante à sigmoide, mas sua faixa de saída está entre -1 e 1, como mostrado na Fig. 3. Ela é simétrica em torno de  $v = 0$  e possui uma forma de "S" suave. Ela é útil em problemas de classificação e regressão. A derivada da tangente hiperbólica também atinge seus valores máximos em torno de  $v = 0$ , mas é mais pronunciada do que a derivada da sigmoide. Isso pode ajudar a mitigar o problema do desaparecimento do gradiente em comparação com a função sigmoide.

- **Função ReLU**

A função ReLU é linear e não linear apenas em um lado, onde retorna 0 para valores de  $v$  negativos e  $v$  para valores positivos, como mostrado na Fig. 4. Ela é simples e computacionalmente eficiente, e tem sido amplamente utilizada em redes neurais profundas. A derivada da ReLU é simples: ela retorna 0 para valores negativos de  $v$  e 1 para valores positivos de  $v$ . Isso significa que não há desaparecimento do gradiente para valores positivos de  $v$ , tornando a ReLU uma escolha popular para redes profundas.

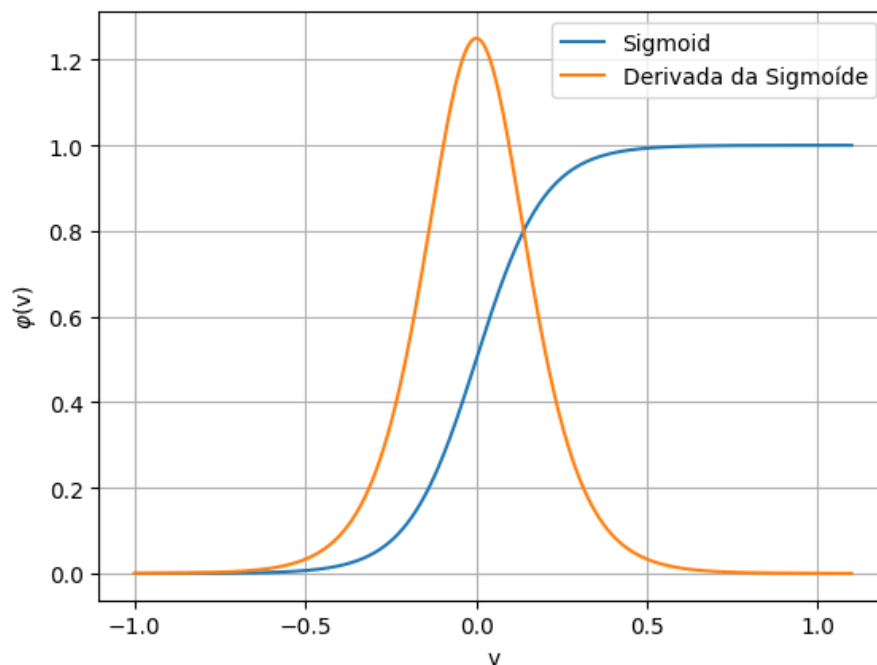


Figure 2. Função de Ativação Sigmoid e sua Derivada

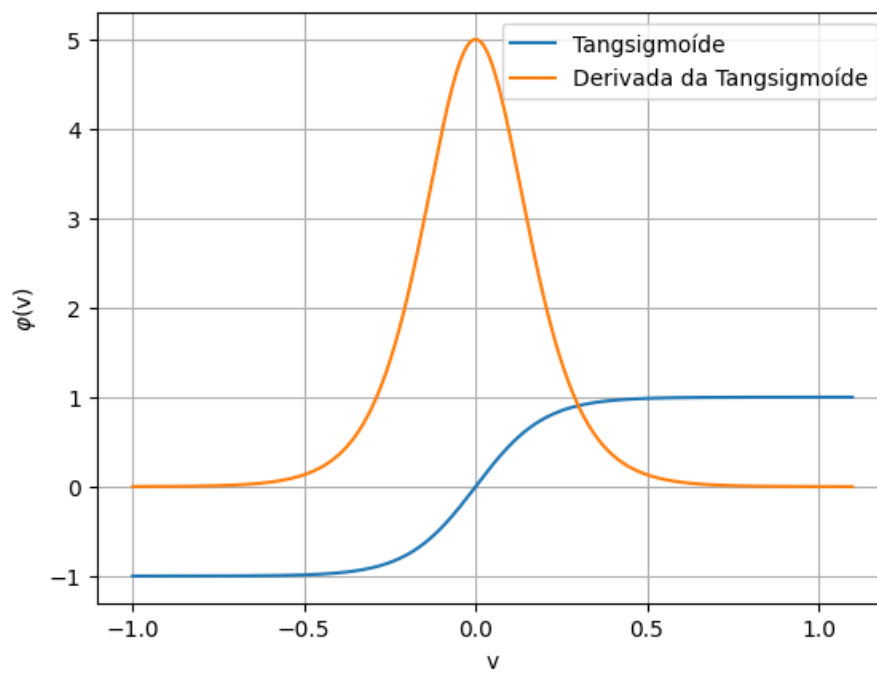


Figure 3. Função de Ativação Tangsigmoide e sua Derivada

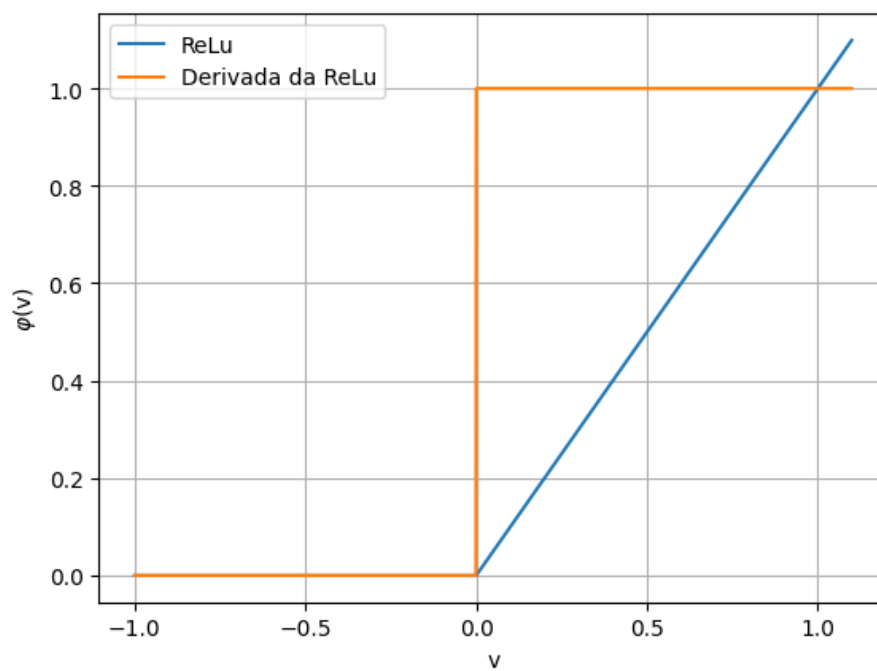


Figure 4. Função de Ativação ReLu e sua Derivada



## II. Derivada função sigmoíde

$$\varphi(v) = \frac{1}{1 + e^{-av}} \Rightarrow \frac{d\varphi(v)}{dv} = \frac{\frac{d1}{dv} \cdot (1 + e^{-av}) - \frac{d(1 + e^{-av})}{dv} \cdot 1}{(1 + e^{-av})^2} = \frac{ae^{-av}}{(1 + e^{-av})^2} \quad (49)$$

$$\Rightarrow \frac{d\varphi(v)}{dv} = \frac{a}{1 + e^{-av}} \left( \frac{1 + e^{-av} - 1}{1 + e^{-av}} \right) = \frac{a}{1 + e^{-av}} \left( 1 - \frac{1}{1 + e^{-av}} \right) = a\varphi(v)[1 - \varphi(v)] \quad (50)$$

## III. Derivada função tangsigmoíde

$$\varphi(v) = \frac{1 - e^{-av}}{1 + e^{-av}} \quad (51)$$

$$\Rightarrow \frac{d\varphi(v)}{dv} = \frac{\frac{d(1 - e^{-av})}{dv} \cdot (1 + e^{-av}) - \frac{d(1 + e^{-av})}{dv} \cdot (1 - e^{-av})}{(1 + e^{-av})^2} = \frac{ae^{-av}(1 + e^{-av}) + ae^{-av}(1 - e^{-av})}{(1 + e^{-av})^2} \quad (52)$$

$$\Rightarrow \frac{a(e^{-av} + e^{-2av} + e^{-av} - e^{-2av})}{(1 + e^{-av})^2} = \frac{a(e^{-av} + e^{-2av} + e^{-av} - e^{-2av})}{(1 + e^{-av})^2} \quad (53)$$

$$\Rightarrow \frac{a[(1 + 2e^{-av} + e^{-2av}) - (1 - 2e^{-av} + e^{-2av})]}{2(1 + e^{-av})^2} = \frac{a[(1 + e^{-av})^2 - (1 - e^{-av})^2]}{2(1 + e^{-av})^2} \quad (54)$$

$$\frac{d\varphi(v)}{dv} = \frac{a}{2} \left[ 1 - \frac{(1 - e^{-av})^2}{(1 + e^{-av})^2} \right] = \frac{a}{2} [1 - \varphi^2(v)] \quad (55)$$

## 7. Funções de saída

### I. Classificação de padrões com duas classes

Para problemas de classificação binária a função de ativação mais comumente utilizada na camada de saída é a **função sigmoíde** (ou logística). Esta função tem a propriedade de mapear os valores de entrada para o intervalo [0,1], o que é útil para representar a probabilidade de pertencer a uma das duas classes. função é dada pela seguinte fórmula:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (56)$$

### II. Classificação de padrões com múltiplas classes

Em casos de classificação com mais de duas classes, o mais comum é utilizar a **função softmax**. A função softmax é uma generalização da função sigmoíde para múltiplas classes. Ela transforma os valores de entrada em uma distribuição de probabilidade sobre as classes, de modo que a soma das saídas seja igual a 1. A função é dada pela seguinte fórmula:

$$\text{softmax}(x)_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}} \quad (57)$$

### III. Problema de regressão (aproximação de funções)

Para problemas de regressão, onde o objetivo é prever um valor contínuo, a função de ativação mais comumente usada na camada de saída é a **função linear**. Esta função simplesmente passa os valores de entrada sem alterá-los. Ela é dada pela fórmula:

$$f(x) = x \quad (58)$$

## 8. Método Max-Likelihood em problema de distribuição gaussiana multivariada

Para modelar um conjunto de dados com várias variáveis, pode-se utilizar a distribuição gaussiana multivariada. Esta distribuição é uma generalização da distribuição gaussiana univariada para múltiplas dimensões. Ela é caracterizada por dois parâmetros principais: o vetor de médias  $\mu$  e a matriz de covariância  $\Sigma$ . O vetor de médias  $\mu$  especifica o centro da distribuição, enquanto a matriz de covariância  $\Sigma$  controla a forma e a orientação da distribuição em torno de  $\mu$ . Para determinar os melhores valores desse dois parâmetros para representar o conjunto de dados se usa o método da máxima verossimilhança (Max-Likelihood). Para isso, utiliza-se a função de densidade de probabilidade da distribuição gaussian multivariada, definida da seguinte maneira:

$$p(\mathbf{x}|\mu, \Sigma) = (2\pi)^{-\frac{D}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right) \quad (59)$$

onde  $\mathbf{x} \in \mathbb{R}^D$ . Para simplificar futuras equações, escreve-se  $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu, \Sigma)$ .

Para formular o aprendizado de parâmetro via máxima verossimilhança, assume-se um dataset  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , onde  $\mathbf{x}_n$ ,  $n = 1, \dots, N$  possuem uma distribuição  $p(\mathbf{x})$ . O objetivo é encontrar uma boa aproximação da distribuição a partir de uma quantidade  $K$  de componentes da mistura das distribuições. Sendo esses parâmetros as médias  $\mu_k$ , a covariância  $\Sigma_k$  e o peso de cada distribuição na mistura  $\pi_k$ . Pode-se compactar todos esses parâmetros em  $\theta = \{\pi_k, \mu_k, \Sigma_k : k = 1, \dots, K\}$ . A função de verossimilhança é o produto das densidades de probabilidade de cada exemplo nos dados de treinamento, dada a distribuição gaussiana multivariada:

$$p(\mathcal{X}|\theta) = \prod_{n=1}^N p(\mathbf{x}_n|\theta), \quad p(\mathbf{x}_n|\theta) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \quad (60)$$

Para encontrar os parâmetros  $\theta_{ML}^*$  que maximizam a função de verossimilhança, pode-se maximizar o logaritmo da função de verossimilhança (pois o logaritmo é uma função monótona crescente). Isso simplifica o cálculo e não afeta a localização dos máximos. Então, calcula-se:

$$\log p(\mathcal{X}|\theta) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k) \quad (61)$$

Está sendo considerado o caso genérico de densidade desejada, porém se considerar uma única distribuição gaussiana como densidade desejada, a soma de  $k$  desaparece, e o log pode ser aplicado diretamente, resultando:

$$\log \mathcal{N}(\mathbf{x}|\mu, \Sigma) = -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \quad (62)$$

Para maximizar a função log-verossimilhança, tomamos as derivadas parciais em relação a  $\mu$  e  $\Sigma$ , e igualamos a zero para encontrar os valores ótimos. Para  $\mu$ :

$$\frac{\partial}{\partial \mu} \log p(\mathcal{X}|\mu, \Sigma) = 0 \Rightarrow \sum_{n=1}^N \frac{\partial}{\partial \mu} \left[ -\frac{D}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right] = 0 \quad (63)$$

$$\Rightarrow \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}_n - \mu) = 0 \Rightarrow \sum_{n=1}^N (\mathbf{x}_n - \mu) = 0 \quad (64)$$

$$\Rightarrow \hat{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \bar{\mathbf{x}} \quad (65)$$

Para  $\Sigma$ :

$$\frac{\partial}{\partial \Sigma^{-1}} \log p(\mathcal{X}|\mu, \Sigma) = 0 \Rightarrow \sum_{n=1}^N \frac{\partial}{\partial \Sigma^{-1}} \left[ -\frac{D}{2} \log(2\pi) + \frac{1}{2} \log |\Sigma^{-1}| - \frac{1}{2} \text{tr}[(\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T \Sigma^{-1}] \right] = 0 \quad (66)$$

$$\Rightarrow \frac{N}{2} \Sigma - \frac{1}{2} \sum_{n=1}^N (\mathbf{x}_n - \mu)(\mathbf{x}_n - \mu)^T = 0 \quad (67)$$

$$\Rightarrow \hat{\Sigma} = \frac{1}{N} \sum_{n=1}^m (\mathbf{x}_n - \hat{\mu})(\mathbf{x}_n - \hat{\mu})^T \quad (68)$$

## 9. Implementações Computacionais de Redes Neurais

### 9.1 MLP para função multivariada

Para resolver esse problema, foi contruído um código em Python utilizando a biblioteca TensorFlow para realizar os cálculos da arquitetura definida de uma rede perceptron de múltiplas camadas. Este é um problema de regressão em que há duas variáveis de entrada e uma de saída, dada pela seguinte função:

$$f(\mathbf{x}) = 16x_1^2 + x_1x_2 + 8x_2^2 - x_1 - x_2 + \ln(1 + x_1^2 + x_2^2) \quad (69)$$

A arquitetura escolhida, após alguns testes, foi de 2:10:10:1, ou seja, duas camadas ocultas de 10 neurônios cada para resolver o problema. O conjunto de dados é separado em 80% para treinamento e 20% para validação. As camadas ocultas utilizam a função de ativação ReLu, por se tratar de um problema de regressão, e a camada de saída é linear. O algoritmo de otimização utilizado foi o Adam devido a sua boa adaptabilidade ao problema. A função custo escolhida foi a dos mínimos médios quadráticos, que é um bom parâmetro para casos de regressão em que dados muito distantes são alamente penalizados e dados mais próximos da função são menos penalizados. A figura a seguir mostra o diagrama da rede neural:

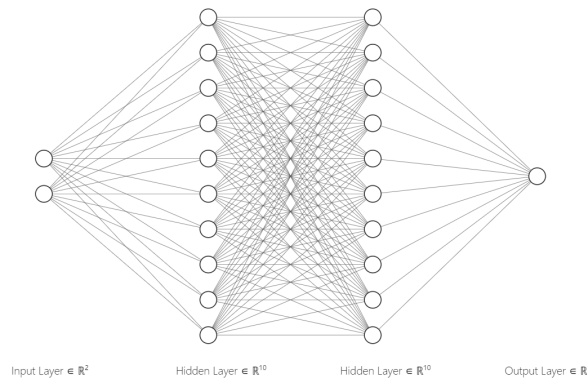


Figure 5. Arquitetura de rede perceptron de múltiplas camadas 2:10:10:1

Ao rodar o código que encontra-se no apêndice, com 500 epochs e uma arquitetura relativamente simples, obtém-se um resultado consideravelmente bom com um erro baixo e uma função bem aproximada da real, como mostrado nas Fig. 6 e Fig. 7 a seguir:

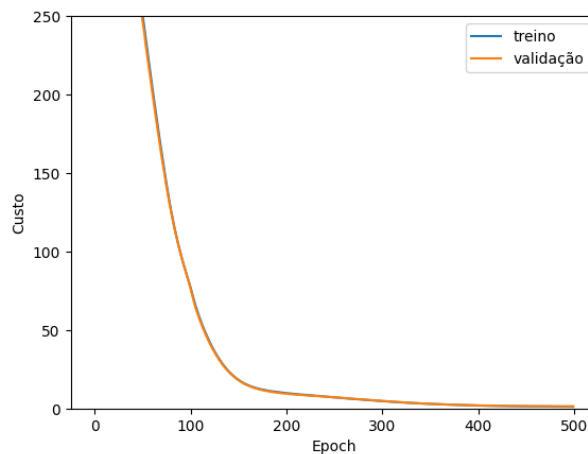


Figure 6. Custo para teste e validação para 2:10:10:1.

Observa-se que uma rede neural relativamente simples é capaz de solucionar esse problema com um valor final de erro baixo. Porém, foi realizado um teste utilizando uma rede neural profunda com a seguinte arquitetura 2:50:100:200:100:50:10, ou seja, bem mais complexa. Com esta mudança há resultados notórios, como convergência em um número menor de epochs e uma função prevista bem mais próxima da função real, como mostrado nas Fig. 8 e Fig. 9.

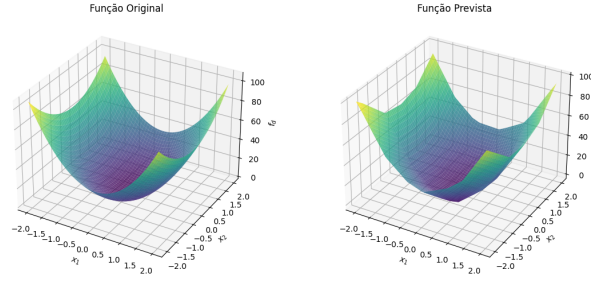


Figure 7. Representação 3D da função resultado treinada pelo modelo para 2:10:10:1.

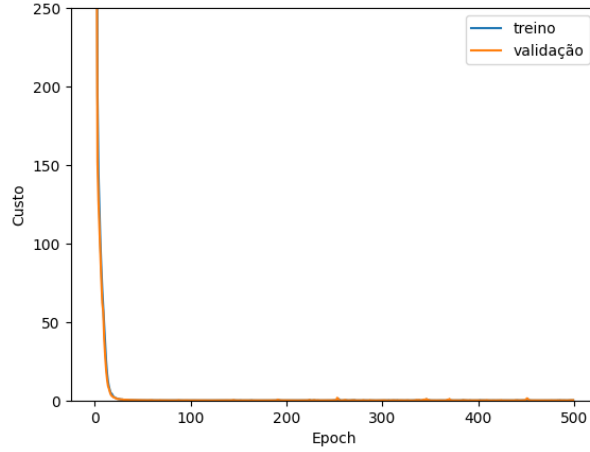


Figure 8. Custo para o caso da rede mais complexa.

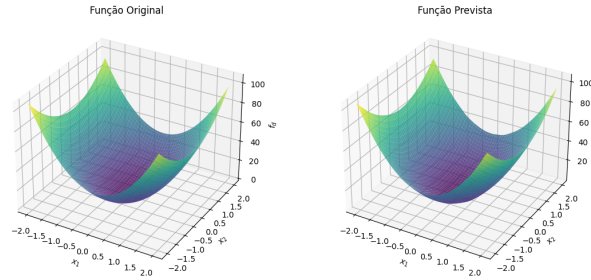


Figure 9. Representação 3D da função prevista pela rede mais complexa.

## 9.2 Problema das espirais

As curvas das espirais são dadas pelas seguintes equações:

$$x = \frac{\theta}{4} \cos \theta, \quad y = \frac{\theta}{4} \sin \theta, \quad \theta \geq 0 \quad (70)$$

$$x = \left(\frac{\theta}{4} + 0.8\right) \cos \theta, \quad y = \left(\frac{\theta}{4} + 0.8\right) \sin \theta, \quad \theta \geq 0 \quad (71)$$

As equações resultam na Fig. 10. Este é um problema de classificação relativamente complexo para uma rede neural, portanto se faz necessário múltiplas camadas ocultas. Para resolver esse problema, foi feito um código em Python utilizando TensorFlow que encontra-se no apêndice. A arquitetura utilizada foi de uma rede neural de múltiplas camadas em forma de funil com duas variáveis de entrada e apenas uma de saída com 6 camadas ocultas formadas da seguinte maneira 2:500:400:300:200:100:50:1. Foi escolhida essa arquitetura após diversos testes, tendo esse o melhor resultado dentro dos limites computacionais aceitáveis. As camadas ocultas utilizam a função de ativação ReLu para que a derivada da função seja maior que zero e o resultado não atinja divergência. Já a camada de saída utiliza a função de ativação sigmoide, devido o tipo de problema ser de classificação, essa função indica a probabilidade de cada classe. O algoritmo de otimização foi o Adam e a função custo a entropia cruzada binária, também devido ao problema de classificação. Após rodar o código, que obteve ótimos resultados, tem-se a seguinte matriz de confusão para os dados de teste e o seguinte gráfico do custo durante o treinamento, mostrado nas Fig. 11 e Fig. 12, respectivamente.

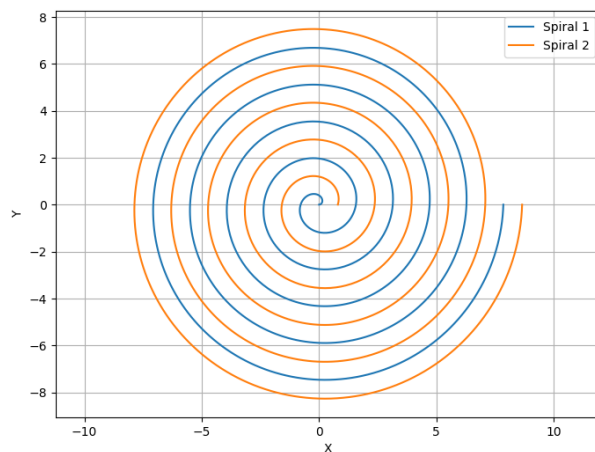


Figure 10. Espirais geradas pelas equações.

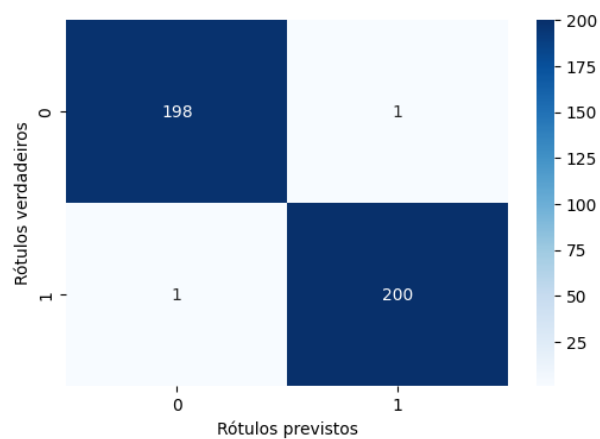


Figure 11. Matriz de confusão dos dados de teste.

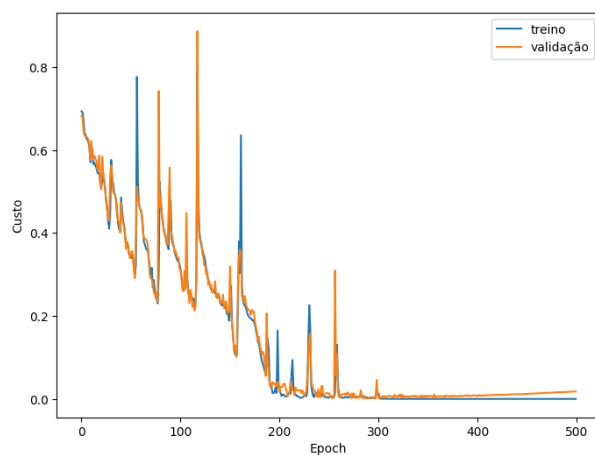


Figure 12. Custo vs epochs do treinamento.

Pode-se agora utilizar o modelo para prever as classes de uma grade de pontos, para testar efetivamente e visualmente o modelo, como mostrado na Fig. 13.

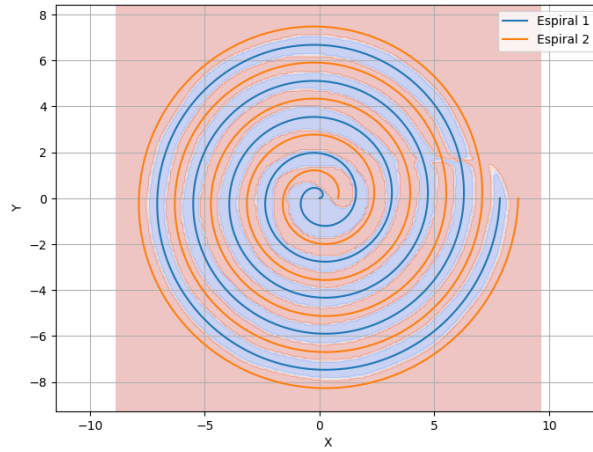


Figure 13. Previsão do modelo.

### 9.3 Problema de classificação de 5 padrões

Este problema propõe uma geometria de um quadrado com 4 semicírculos dispostos dentro dele de tal forma que 5 área delimitadas são geradas. Pode-se visualizar a geometria do problema na Fig. 14. As equações que delimitam a região azul, amarela, verde e vermelha são descritas pelas Eq. 72, Eq. 73, Eq. 74, Eq. 75, respectivamente, e o quadrado pela Eq. 76.

$$-\sqrt{1-x^2} + 1 \leq y \leq \sqrt{1-(x-1)^2}, \quad 0 < x < 1 \quad (72)$$

$$-\sqrt{1-(x-1)^2} \leq y \leq \sqrt{1-x^2} - 1, \quad 0 < x < 1 \quad (73)$$

$$-\sqrt{1-(x+1)^2} \leq y \leq \sqrt{1-x^2} - 1, \quad -1 < x < 0 \quad (74)$$

$$-\sqrt{1-x^2} + 1 \leq y \leq \sqrt{1-(x+1)^2}, \quad -1 < x < 0 \quad (75)$$

$$x = 1, \quad x = -1, \quad y = 1, \quad y = -1 \quad (76)$$

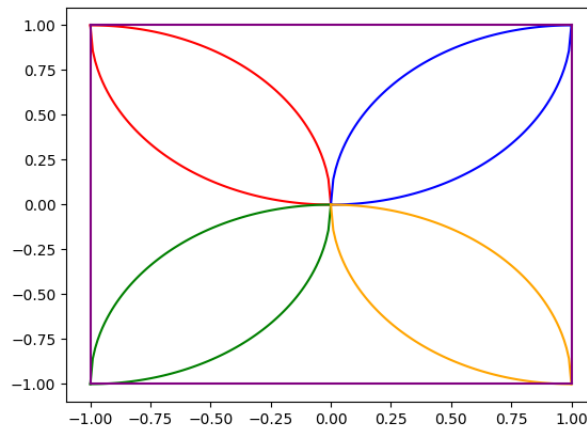


Figure 14. Geometria proposta pela problema.

Então, com a geometria definida, pode-se gerar os dados do problema de classificação, 10000 pontos foram gerados aleatoriamente no total, de tal forma que a distribuição dos pontos por classe pode ser visualizado no histograma da Fig. 15 e também pode-se visualizar os pontos gerados na própria geometria na Fig. 16. Naturalmente, há mais pontos gerados na classe 5, devido possuir maior área na geometria.

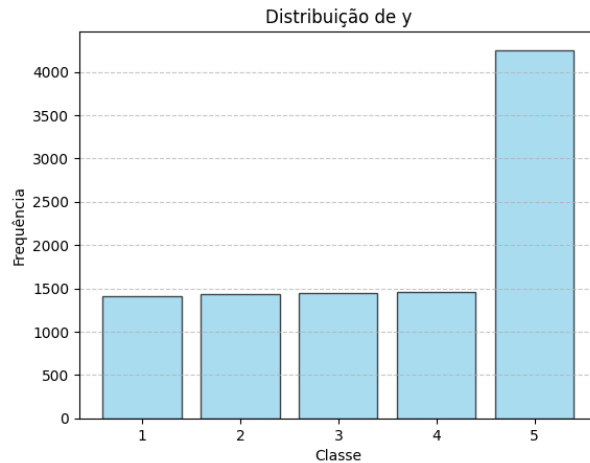


Figure 15. Histograma da distribuição dos pontos gerados por classe.

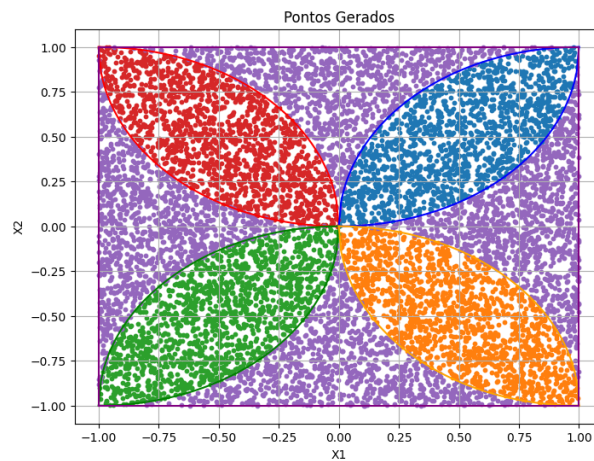


Figure 16. Pontos gerados na geometria.

Após gerar os pontos, os conjuntos de treinamento e teste foram divididos, de tal maneira que 20% foram destinados para o treinamento e 80% para o teste. Foi decidida esta divisão, pois há um grande número de pontos gerados, portanto é interessante que apenas uma quantidade suficiente de pontos seja destinada ao treinamento para que faça sentido o teste do modelo, com pontos excessivos de treinamento, basicamente qualquer modelo será capaz de resolver. Após diversos testes, encontrou-se uma arquitetura de números de neurônios e camadas equilibradas de 2:50:50:50:50:6. Durante os testes de diferentes modelos, observou-se que um número elevado de camadas causava ruído na função custo, atrapalhando a convergência, e um número elevado de neurônios impossibilitava de o modelo atingir valores mais baixos de erro. O algoritmo otimizador utilizado foi o Adam e a função custo foi a entropia cruzada categórica esparsa, esta que é mais indicada para problemas de classificação de multiclases, também é responsável por existir 6 neurônios na camada de saída, mesmo existindo apenas 5 padrões para classificar, devido a como o limite do seu cálculo é realizada. As funções de ativação foram ReLu nas camadas ocultas e na camada de saída foi a Softmax que proporciona a probabilidade de em qual classe está o conjunto de entrada de forma eficaz. Por fim, o modelo é treinado por 500 epochs e um tamanho de lote de 500. O erro do teste foi de 0,059 e a acurácia de 97,5%. O código realizado em python utilizando Tensorflow pode ser encontrado no apêndice. Os resultados são mostrados a seguir:

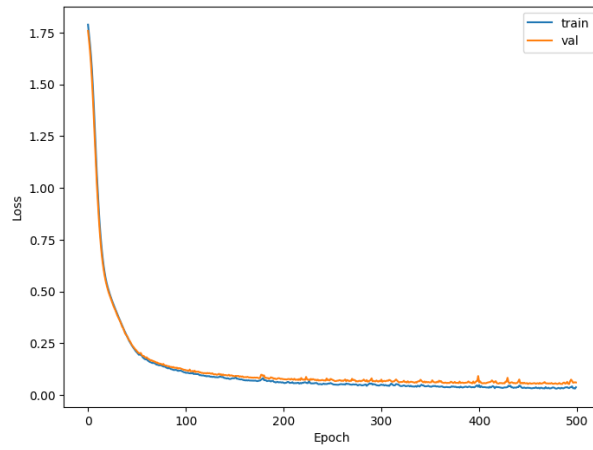


Figure 17. Função custo.

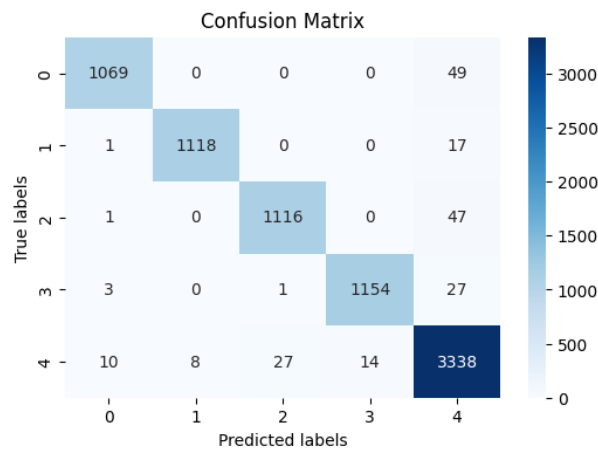


Figure 18. Matriz de confusão.

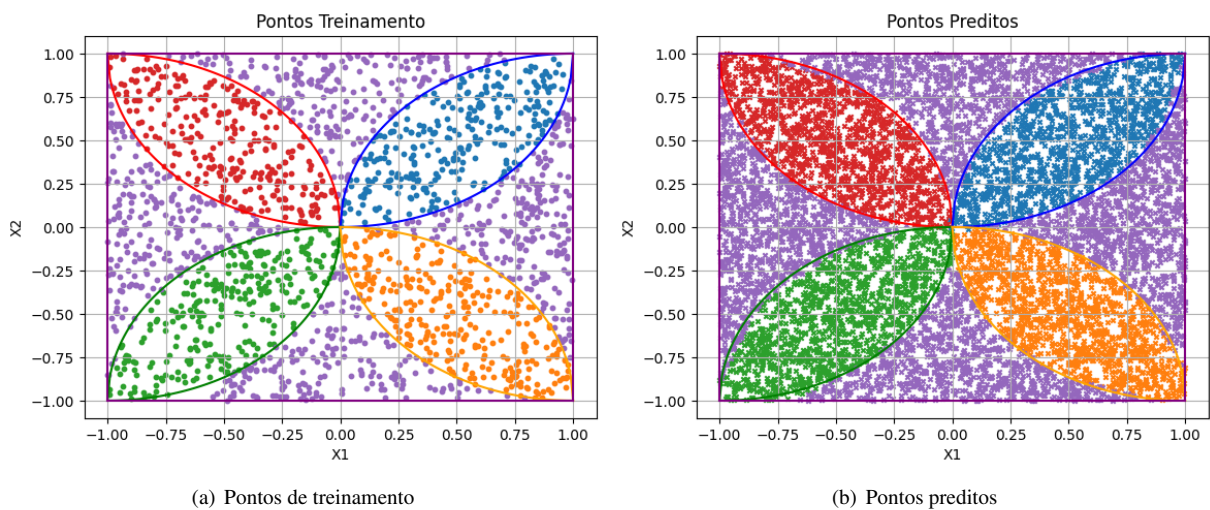


Figure 19. Comparação entre os pontos entregues para treinamento e os pontos de teste preditos pelo modelo.



## 10. Aprendizado Profundo para predições de Simulações Aerodinâmicas

A tese de Crichigno, *Deep Learning for prediction of Aerodynamic Simulations*, faz uma aplicação de Redes Neurais Artificiais (RNNs) na área de aerodinâmica, com o objetivo de desenvolver uma RNN de Perceptron de múltiplas camadas para prever os resultados de simulações de Dinâmica dos Fluidos Computacional (CFD) de um aerofólio bidimensional. O CFD é uma ferramenta extremamente valiosa para análises aerodinâmicas no geral, atualmente existem diversos *softwares*, tanto comerciais como código aberto, capazes de simular modelos complexos de fluido e geometria. No estudo são implementadas duas RNNs, uma para prever os coeficientes aerodinâmicos de sustentação ( $C_L$ ) e de arrasto ( $C_D$ ) do aerofólio. E uma segunda RNN para prever o campo de velocidade ao redor do aerofólio.

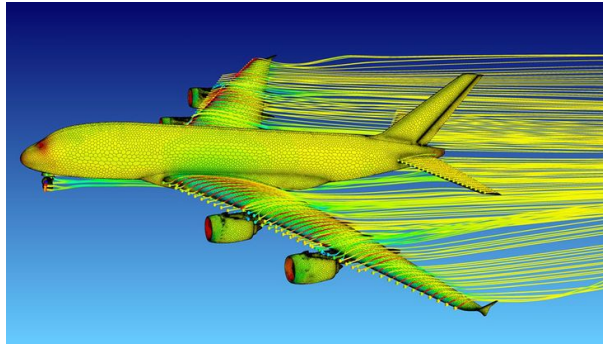


Figure 20. Exemplo CFD de aeronave.

### 10.1 Predição dos Coeficientes Aerodinâmicos

Primeiramente, para uma boa implementação de Perceptron de Múltiplas Camadas (MLP) é necessário um bom conjunto de dados para treinar, validar e testar, para isso o estudo dispõe de 7056 resultados de simulações CFD. Essa primeira MLP utiliza como parâmetros de entrada as condições do escoamento: ângulo de ataque  $\alpha$ , número de Mach  $M_\infty$  e número de Reynolds  $Re_\infty$ .

DADOS DE ENTRADA	TARGETS
$[\alpha, M_\infty, Re_\infty]_1$	$[C_L, C_D]_1$
$[\alpha, M_\infty, Re_\infty]_2$	$[C_L, C_D]_2$
$\vdots$	$\vdots$
$[\alpha, M_\infty, Re_\infty]_N$	$[C_L, C_D]_N$

Table 1. Conjunto de dados para os coeficientes aerodinâmicos,  $N = 7056$  amostras

#### 10.1.1 Preparação do conjunto de dados

A divisão para o conjunto de dados do estudo é: Treinamento 70%, Validação 20% e Teste 10%. Outra ação realizada pelo estudo consiste em normalizar os dados de entrada, uma vez que a magnitude dos parâmetros varia consideravelmente em ordem de magnitude. Isso pode levar a estimativas de importância discrepantes para cada parâmetro, prejudicando os resultados do treinamento. A normalização é realizada em relação ao valor médio e ao desvio padrão do conjunto de Treinamento apenas. Esta abordagem é adotada porque o conjunto de treinamento compreende as amostras nas quais a rede aprende o padrão que liga a entrada à saída.

#### 10.1.2 Função Custo

A função custo escolhida pelo estudo é a do Erro Quadrático Médio (MSE), visto que este é um problema de regressão, então essa é uma métrica comumente adotada.

#### 10.1.3 Algoritmo de Treinamento

O algoritmo adotado para essa primeira MLP é o algoritmo Adam com uma taxa de aprendizagem dinâmica.

#### 10.1.4 Hiperparâmetros: Tamanho de lote

O estudo realiza *grid search* para encontrar os melhores valores para alguns hiperparâmetros selecionados. Um deles é o tamanho do lote, com os valores de 25, 50, 100 e 150. Esse hiperparâmetro coloca na balança a acurácia e o tempo de treinamento.

#### 10.1.5 Hiperparâmetros: Arquitetura da Rede

A MLP do estudo possui 3 entradas e 2 saídas, e busca variar a quantidade de camadas ocultas e números de neurônios para encontrar uma arquitetura otimizada. Todas as camadas ocultas utilizam a função de ativação ReLu e a camada de saída é linear. As arquiteturas testadas foram:

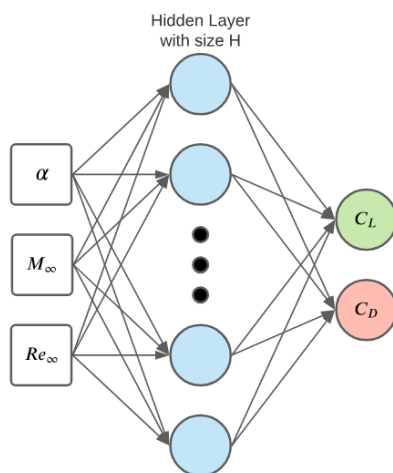


Figure 21. Arquitetura de Única Camada Oculta

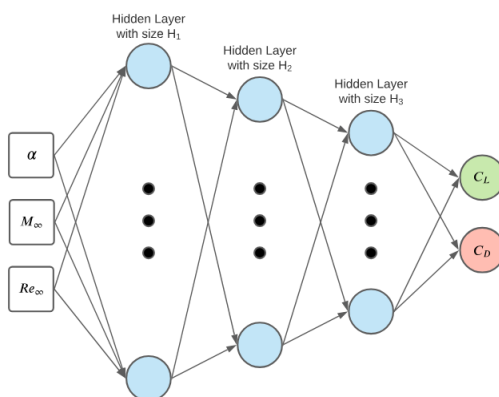


Figure 22. Arquitetura de Funil

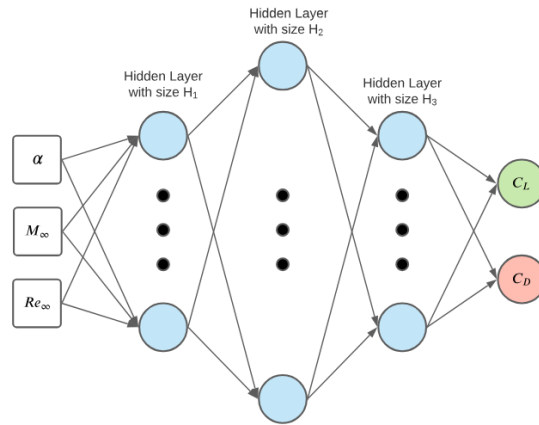


Figure 23. Arquitetura Rhombus

Camadas Ocultas						
H <sub>1</sub>	H <sub>2</sub>	H <sub>3</sub>	H <sub>4</sub>	H <sub>5</sub>	H <sub>6</sub>	H <sub>7</sub>
40						
100						
40	20					
100	10					
80	40	20				
160	80	40	20			
20	40	20				
20	40	80	40	20		
20	40	80	160	80	40	20

Table 2. Quantidade de neurônios por camada oculta.

## 10.2 Resultados

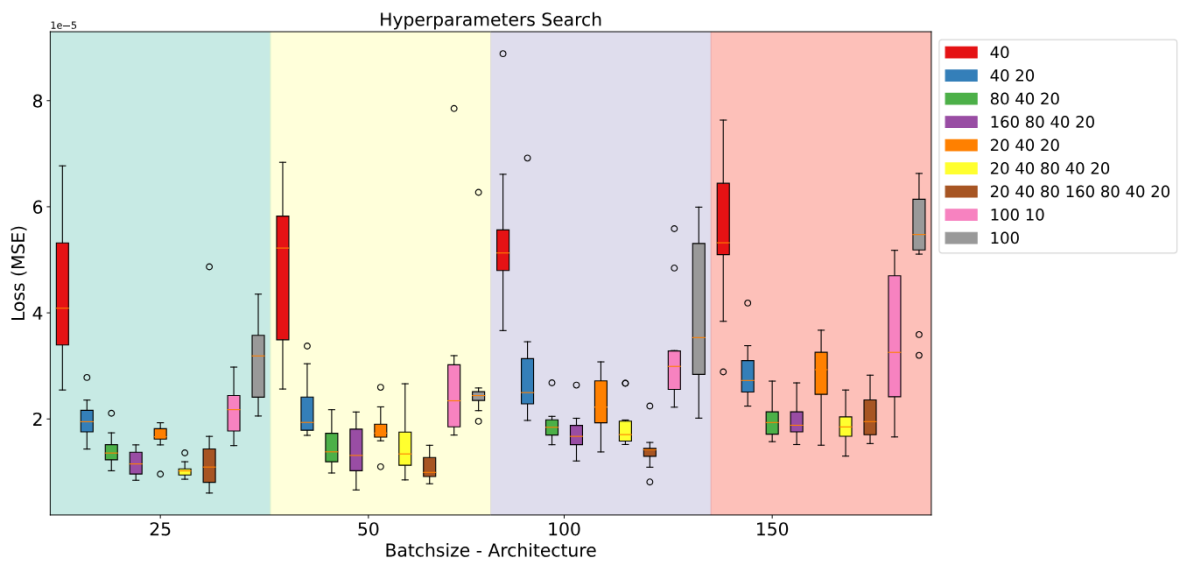


Figure 24. Resultados do gridsearch para o MLP dos coeficientes  $C_L$  e  $C_D$ .

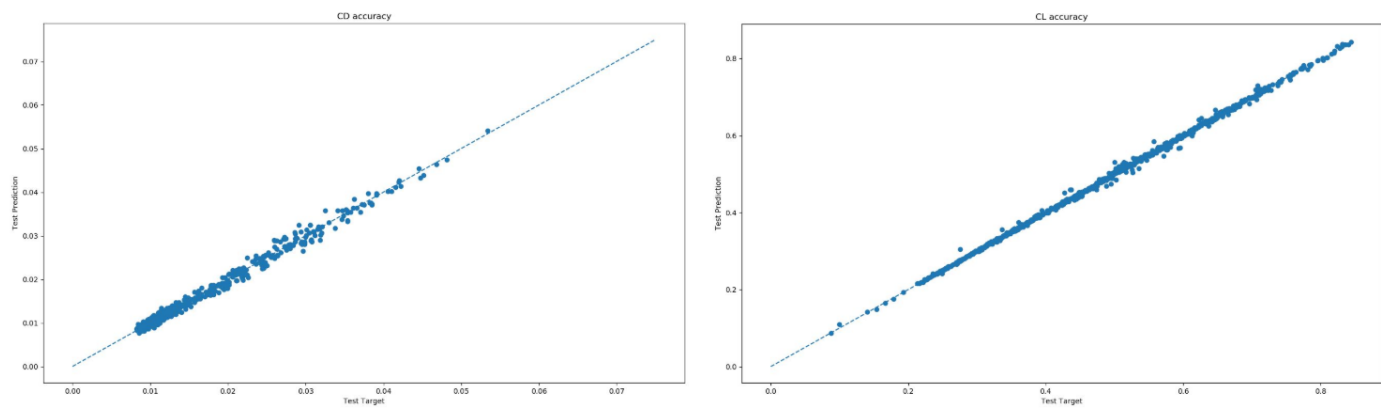


Figure 25. Acurácia para a predição de teste para  $C_L$  (direita) e  $C_D$  (esquerda).

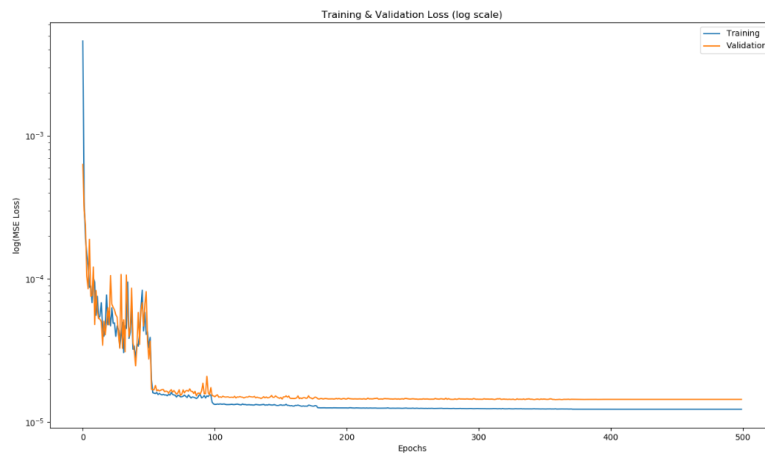


Figure 26. Custo vs Epochs da MLP.

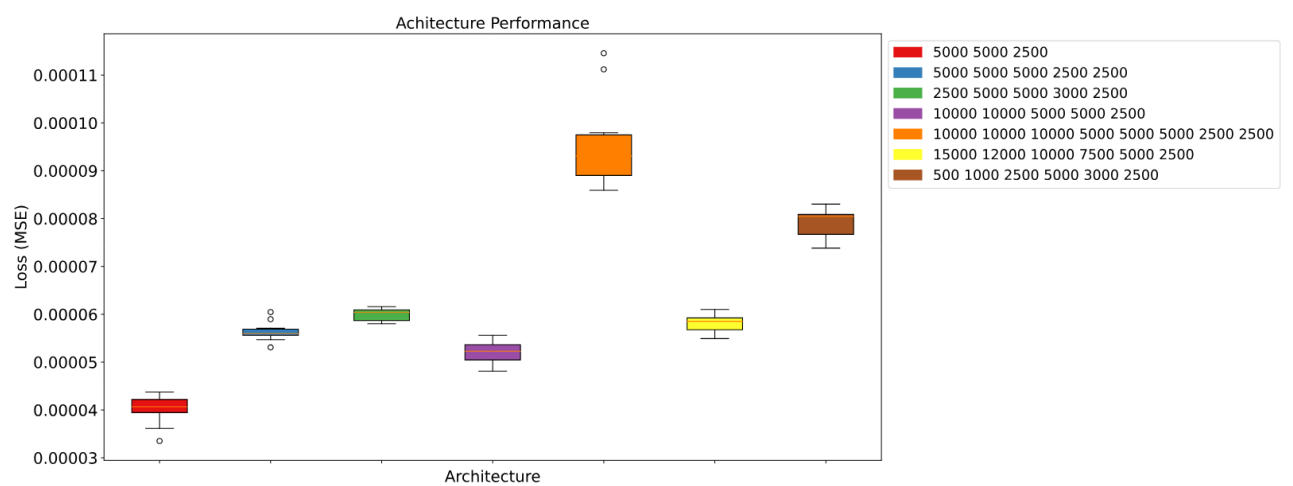


Figure 27. Resultados da gridsearch para a MLP da imagem cinza.

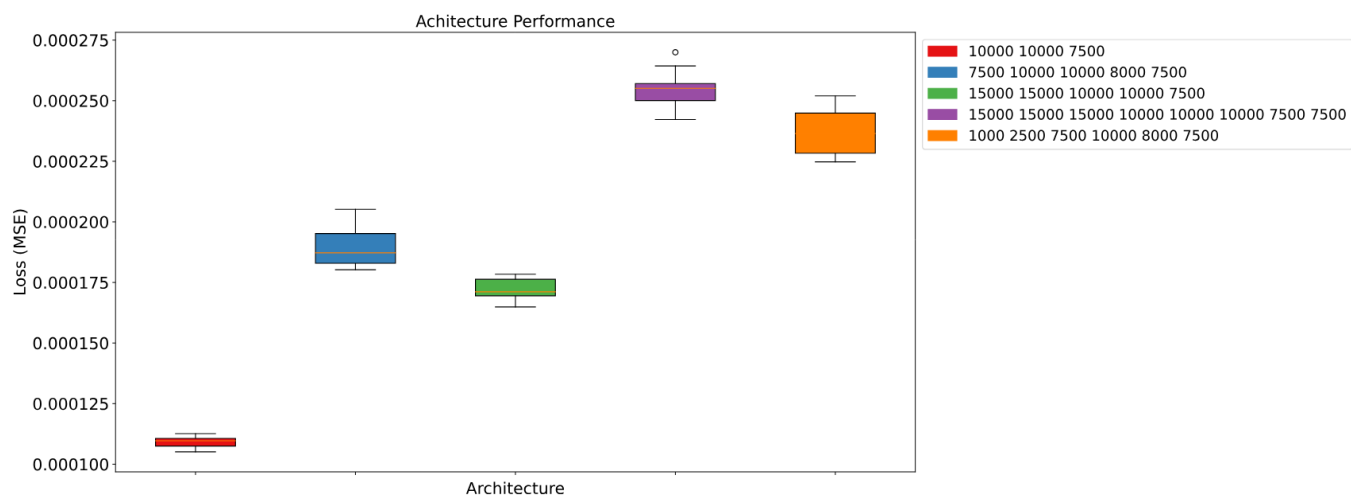


Figure 28. Resultados da gridsearch para a MLP da imagem colorida.

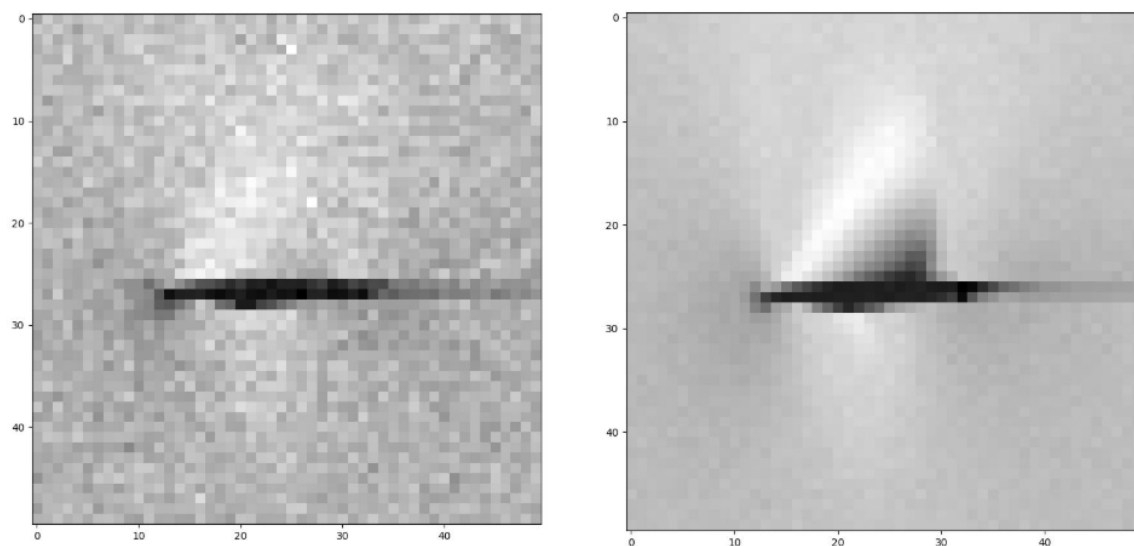


Figure 29. Predição do teste no começo do treinamento (esquerda) e após 500 epochs.

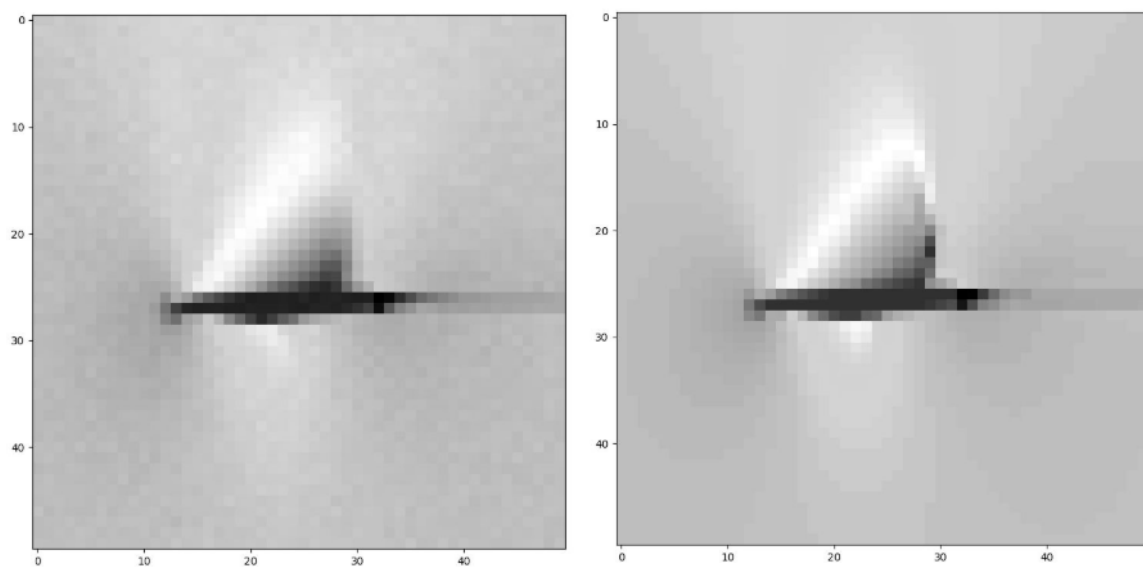


Figure 30. Comparação da imagem após o treinamento (esquerda) com a imagem desejada (direita).

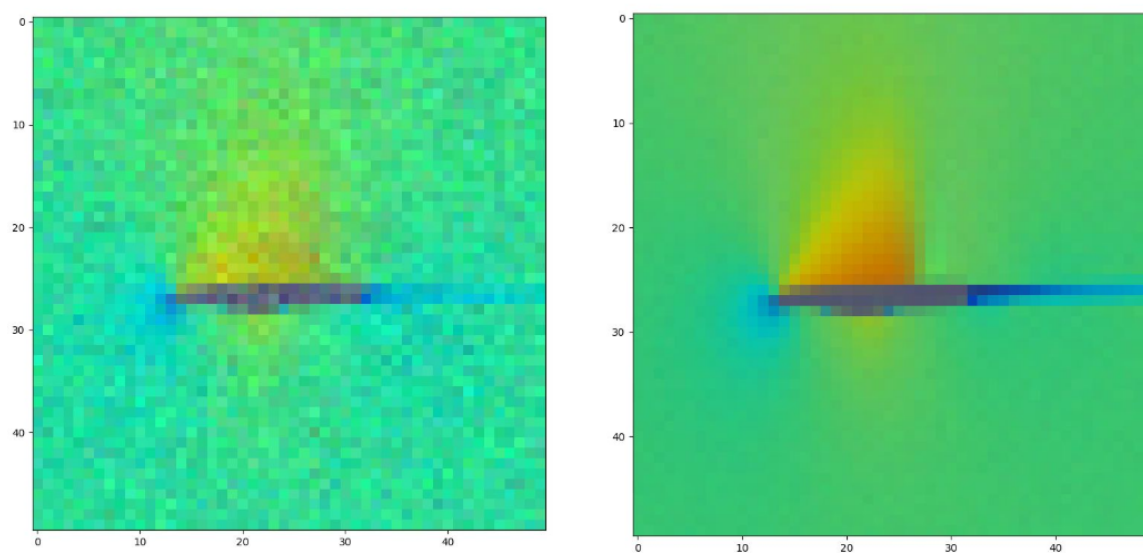


Figure 31. Predição do teste no começo do treinamento (esquerda) e após 500 epochs.

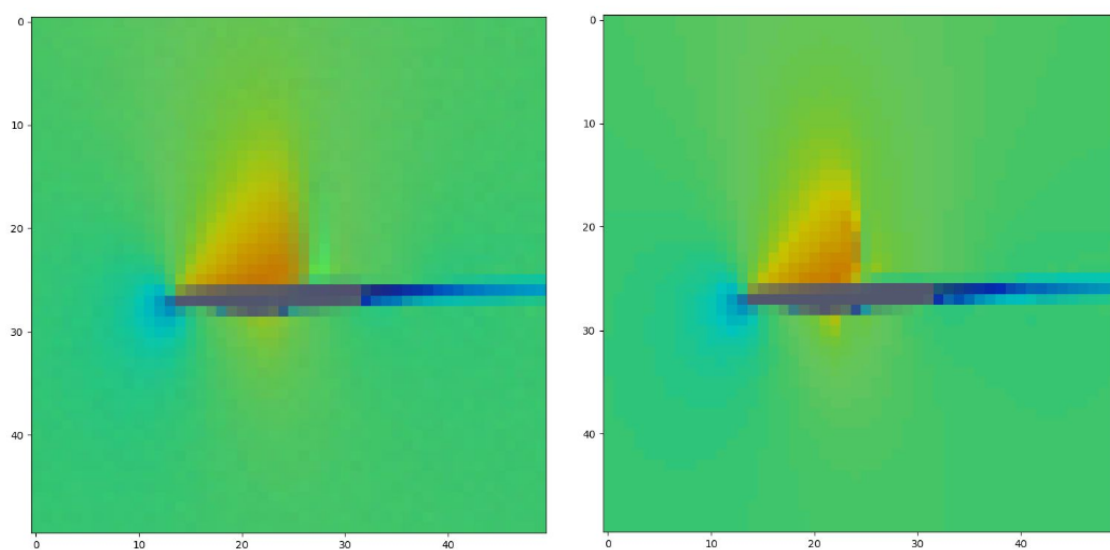


Figure 32. Comparação da imagem após o treinamento (esquerda) com a imagem desejada (direita).

# Appendices

## A Código Python questão 9-1

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense

# Definir a função alvo
def f_d(x1, x2):
    return 16*x1**2 + x1*x2 + 8*x2**2 - x1 - x2 + np.log(1 + x1**2 + x2**2)

# Gerar dados de entrada
np.random.seed(42)
n_samples = 1000
X = np.random.uniform(-2, 2, (n_samples, 2))

# Calcular os valores de saída correspondentes
y = np.array([f_d(x[0], x[1]) for x in X])

# Dividir os dados em conjuntos de treinamento e validação
X_train, X_val, y_train, y_val = train_test_split(X, y, test_size=0.2, random_state=42)

# Construir o modelo da rede neural
model = Sequential([
    Dense(10, activation='relu', input_shape=(2,)),
    Dense(10, activation='relu'),
    Dense(1) # Camada de saída
])

# Compilar o modelo
model.compile(optimizer='Adam', loss='mean_squared_error')

# Treinar o modelo
history = model.fit(X_train, y_train, epochs=500, batch_size=50, validation_data=(X_val, y_val))

# Visualizar a curva do erro de treinamento e validação
plt.plot(history.history['loss'], label='treino')
plt.plot(history.history['val_loss'], label='validação')
plt.ylim([0, 250])
plt.xlabel('Epoch')
plt.ylabel('Custo')
plt.legend()
plt.show()

# Gerar uma grade de pontos para plotagem
x1_values = np.linspace(-2, 2, 100)
x2_values = np.linspace(-2, 2, 100)
x1_grid, x2_grid = np.meshgrid(x1_values, x2_values)
X_grid = np.column_stack([x1_grid.ravel(), x2_grid.ravel()])

# Calcular as previsões da rede neural para a grade de pontos
y_pred = model.predict(X_grid).reshape(x1_grid.shape)

fig = plt.figure(figsize=(12, 5))

ax1 = fig.add_subplot(1, 2, 1, projection='3d')
ax1.plot_surface(x1_grid, x2_grid, f_d(x1_grid, x2_grid), cmap='viridis', alpha=0.8)
ax1.set_title('Função Original')
ax1.set_xlabel('$x_1$')
ax1.set_ylabel('$x_2$')
ax1.set_zlabel('$f_d$')

ax2 = fig.add_subplot(1, 2, 2, projection='3d')
ax2.plot_surface(x1_grid, x2_grid, y_pred, cmap='viridis', alpha=0.8)
ax2.set_title('Função Prevista')
ax2.set_xlabel('$x_1$')
ax2.set_ylabel('$x_2$')
```



```
ax2.set_zlabel('$f_d$')

plt.tight_layout()
plt.show()
```

## B Código Python questão 9-2

```
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
import seaborn as sns
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense

# Definir as equações das espirais
def spiral_1(theta):
    x = (theta / 4) * np.cos(theta)
    y = (theta / 4) * np.sin(theta)
    return x, y

def spiral_2(theta):
    x = ((theta / 4) + 0.8) * np.cos(theta)
    y = ((theta / 4) + 0.8) * np.sin(theta)
    return x, y

# Gerar pontos das espirais
theta_1 = np.linspace(0, 10*np.pi, 1000)
x1, y1 = spiral_1(theta_1)

theta_2 = np.linspace(0, 10*np.pi, 1000)
x2, y2 = spiral_2(theta_2)

# Plotar as espirais
plt.figure(figsize=(8, 6))
plt.plot(x1, y1, label='Spiral 1')
plt.plot(x2, y2, label='Spiral 2')
plt.title('Spirals')
plt.xlabel('X')
plt.ylabel('Y')
plt.legend()
plt.grid(True)
plt.axis('equal')

# Gerar dados para espiral 1
spiral_1_data = np.column_stack((x1, y1))
labels_1 = np.zeros(len(theta_1))

# Gerar dados para espiral 2
spiral_2_data = np.column_stack((x2, y2))
labels_2 = np.ones(len(theta_2))

# Combinar dados e rótulos
X = np.concatenate((spiral_1_data, spiral_2_data), axis=0)
y = np.concatenate((labels_1, labels_2))

# Dividir os dados em conjuntos de treinamento e teste
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Construir o modelo da rede neural
model = Sequential([
    Dense(500, activation='relu', input_shape=(2,)),
    Dense(400, activation='relu'),
    Dense(300, activation='relu'),
    Dense(200, activation='relu'),
    Dense(100, activation='relu'),
    Dense(50, activation='relu'),
    Dense(1, activation='sigmoid') # Camada de saída
])

# Compilar o modelo
```

```

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Treinar o modelo
history = model.fit(X_train, y_train, epochs=500, batch_size=50, validation_data=(X_test,
                                                                                   y_test), verbose=0)

# Avaliar o modelo
test_loss, test_accuracy = model.evaluate(X_test, y_test, verbose=0)
print(f'Test Loss: {test_loss}')
print(f'Test Accuracy: {test_accuracy}')

# Gerar previsões
y_pred = model.predict(X_test)
y_pred_binary = np.round(y_pred)

# Calcular matriz de confusão
conf_matrix = confusion_matrix(y_test, y_pred_binary)
print('Confusion Matrix:')
print(conf_matrix)

# Plotar matriz de confusão
plt.figure(figsize=(6, 4))
sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='g')
plt.title('Confusion Matrix')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.show()

# Visualizar a curva do erro de treinamento e validação
plt.figure(figsize=(8, 6))
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='val')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

# Gerar uma grade de pontos para fazer previsões
x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
xx, yy = np.meshgrid(np.arange(x_min, x_max, 0.1),
                     np.arange(y_min, y_max, 0.1))

# Prever as classes para cada ponto na grade
Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
Z = Z.reshape(xx.shape)

# Plotar as espirais reais
plt.figure(figsize=(8, 6))
plt.plot(x1, y1, label='Espiral 1')
plt.plot(x2, y2, label='Espiral 2')
plt.xlabel('X')
plt.ylabel('Y')
plt.legend()
plt.grid(True)
plt.axis('equal')

# Plotar as regiões previstas pela rede neural
plt.contourf(xx, yy, Z, alpha=0.3, cmap='coolwarm')
plt.show()

```

## C Código Python questão 9-3

```

import numpy as np
import matplotlib.pyplot as plt

def plot_problem():
    xp = np.linspace(0,1,100) #Espaço de pontos x>0
    xm = np.linspace(-1,0,100) #Espaço de pontos x<0

    ycm = -np.sqrt(1-xm**2)+1 # Curva superior direita
    ycp = -np.sqrt(1-xp**2)+1 # Curva superior esquerda

```

```

ybm = np.sqrt(1-xm**2)-1 # Curva inferior esquerda
ybp = np.sqrt(1-xp**2)-1 # Curva superior direita
ydc = np.sqrt(1-(xp-1)**2) # Curva direita superior
ydb = -np.sqrt(1-(xp-1)**2) # Curva direita inferior
yeb = -np.sqrt(1-(xm+1)**2) # Curva esquerda inferior
yec = np.sqrt(1-(xm+1)**2) # Curva esquerda superior

plt.plot(xm,ycm, 'r') # Curva superior direita
plt.plot(xm,yec, 'r') # Curva esquerda superior
plt.plot(xp,ycp, 'b') # Curva superior esquerda
plt.plot(xp,ydc, 'b') # Curva direita superior
plt.plot(xm,ybm, 'g') # Curva inferior esquerda
plt.plot(xm,yeb, 'g') # Curva esquerda inferior
plt.plot(xp,ybp, 'orange') # Curva superior direita
plt.plot(xp,ydb, 'orange') # Curva direita inferior

#plot do quadrado
plt.plot([-1, -1], [-1, 1], 'purple')
plt.plot([1, 1], [-1, 1], 'purple')
plt.plot([-1, 1], [-1, -1], 'purple')
plt.plot([-1, 1], [1, 1], 'purple')

plt.show()

plot_problem()
# Função para gerar os dados
def generate_data(num_samples):
    X = np.zeros((num_samples, 2))
    y = np.zeros(num_samples, dtype=int)

    for i in range(num_samples):
        x1 = np.random.uniform(-1, 1)
        x2 = np.random.uniform(-1, 1)

        if x1 >= 0 and x1 <=1:
            if x2 >= -np.sqrt(1-x1**2)+1 and x2 <= np.sqrt(1-(x1-1)**2):
                y[i] = 1
            elif x2 >= -np.sqrt(1-(x1-1)**2) and x2 <= np.sqrt(1-x1**2)-1:
                y[i] = 2
            else:
                y[i] = 5
        elif x1 >=-1 and x1 <0:
            if x2 >= -np.sqrt(1-(x1+1)**2) and x2 <= np.sqrt(1-x1**2)-1:
                y[i] = 3
            elif x2 >= -np.sqrt(1-x1**2)+1 and x2 <= np.sqrt(1-(x1+1)**2):
                y[i] = 4
            else:
                y[i] = 5
        else:
            y[i] = 5

        X[i] = [x1, x2]

    return X, y

# Gerar dados
num_samples = 10000
X, y = generate_data(num_samples)

# Plotar histograma de y
plt.hist(y, bins=np.arange(1, 7) - 0.5, rwidth=0.8, color='skyblue', edgecolor='black', alpha=0.7)

plt.xticks(range(1, 6))
plt.xlabel('Classe')
plt.ylabel('Frequência')
plt.title('Distribuição de y')
plt.grid(axis='y', linestyle='--', alpha=0.7)
plt.show()

# Plotar os pontos gerados
plt.figure(figsize=(8, 6))

```

```

for label in range(1, 6):
    plt.scatter(X[y == label][:, 0], X[y == label][:, 1], label=f'Classe {label}', s=10)

plt.xlabel('X1')
plt.ylabel('X2')
plt.title('Pontos Gerados')
plt.grid(True)
plot_problem()
plt.show()
import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
import tensorflow as tf
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import Dense
import seaborn as sns

# Dividir os dados em conjuntos de treinamento e teste
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.8, random_state=42)

# Construir o modelo da rede neural
model = Sequential([
    Dense(50, activation='relu', input_shape=(2,)),
    Dense(50, activation='relu'),
    Dense(50, activation='relu'),
    Dense(50, activation='relu'),
    Dense(6, activation='softmax') # Camada de saída
])

# Compilar o modelo
model.compile(optimizer='adam', loss='sparse_categorical_crossentropy', metrics=['
    sparse_categorical_accuracy'])

# Treinar o modelo
history = model.fit(X_train, y_train, epochs=500, batch_size=500, validation_data=(X_test,
    y_test), verbose=0)

# Avaliar o modelo
test_loss, test_accuracy = model.evaluate(X_test, y_test, verbose=0)
print(f'Test Loss: {test_loss}')
print(f'Test Accuracy: {test_accuracy}')

# Gerar previsões
y_pred = model.predict(X_test)
y_pred_labels = np.argmax(y_pred, axis=1)

# Calcular matriz de confusão
conf_matrix = confusion_matrix(y_test, y_pred_labels)
print('Confusion Matrix:')
print(conf_matrix)

# Plotar matriz de confusão
plt.figure(figsize=(6, 4))
sns.heatmap(conf_matrix, annot=True, cmap='Blues', fmt='g')
plt.title('Confusion Matrix')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.show()

# Visualizar a curva do erro de treinamento e validação
plt.figure(figsize=(8, 6))
plt.plot(history.history['loss'], label='train')
plt.plot(history.history['val_loss'], label='val')
plt.xlabel('Epoch')
plt.ylabel('Loss')
plt.legend()
plt.show()

# Plotar os pontos previstos
for label in range(1, 6):

```

```

plt.scatter(X_test[y_pred_labels == label][:, 0], X_test[y_pred_labels == label][:, 1],
            label=f'Predito {label}', marker='x', s=10)

plt.xlabel('X1')
plt.ylabel('X2')
plt.title('Pontos Preditos')
#plt.legend()
plt.grid(True)
plot_problem()
plt.show()

# Plotar os pontos de cada classe no conjunto de teste
for label in range(1, 6):
    plt.scatter(X_train[y_train == label][:, 0], X_train[y_train == label][:, 1], label=f'
                Classe {label}', s=10)

plt.xlabel('X1')
plt.ylabel('X2')
plt.title('Pontos Treinamento')
#plt.legend()
plt.grid(True)
plot_problem()
plt.show()

```

Texto apenas para inserir referências Goodfellow *et al.* (2016) Haykin (2009) Deisenroth *et al.* (2020) Crichigno (2021) Bourret (2018) Lenail (2019) class (2019)

## D REFERENCES

- Bourret, X., 2018. “Maximum likelihood estimators - multivariate gaussian”. <https://stats.stackexchange.com/questions/351549/maximum-likelihood-estimators-multivariate-gaussian>. Accessed 05 May 2024.
- class, S.C., 2019. “Cs231n: Convolutional neural networks for visual recognition”. <https://cs231n.github.io/neural-networks-1/actfun>. Accessed 05 May 2024.
- Crichigno, N., 2021. *Deep Learning for prediction of Aerodynamic Simulations*. Master’s thesis, POLITECNICO DI TORINO.
- Deisenroth, M.P., Faisal, A.A. and Ong, C.S., 2020. *MATHEMATICS FOR MACHINE LEARNING*. Cambridge University Press. <https://mml-book.com>.
- Goodfellow, I., Bengio, Y. and Courville, A., 2016. *Deep Learning*. MIT Press. <http://www.deeplearningbook.org>.
- Haykin, S., 2009. *Neural Networks and Learning Machines*. Pearson.
- Lenail, A., 2019. “Nn-svg: Publication-ready neural network architecture schematics”. <https://alexlenail.me/NN-SVG/index.html>. Accessed 05 May 2024.