

UNIVERSIDADE ESTADUAL PAULISTA "JÚLIO DE MESQUITA FILHO"

FACULDADE DE CIÊNCIAS - CAMPUS BAURU

DEPARTAMENTO DE COMPUTAÇÃO

BACHARELADO EM CIÊNCIA DA COMPUTAÇÃO

VINÍCIUS DE PAULA PILAN

**TÉCNICAS DE INTELIGÊNCIA ARTIFICIAL PARA DIAGNÓSTICO
DE ACIDENTE VASCULAR CEREBRAL ATRAVÉS DE IMAGENS E
DADOS TEXTUAIS SOBRE POSSÍVEIS VÍTIMAS**

BAURU

Janeiro/2023

VINÍCIUS DE PAULA PILAN

**TÉCNICAS DE INTELIGÊNCIA ARTIFICIAL PARA DIAGNÓSTICO
DE ACIDENTE VASCULAR CEREBRAL ATRAVÉS DE IMAGENS E
DADOS TEXTUAIS SOBRE POSSÍVEIS VÍTIMAS**

Trabalho de Conclusão de Curso do Curso
de Ciência da Computação da Universidade
Estadual Paulista “Júlio de Mesquita Filho”,
Faculdade de Ciências, Campus Bauru.
Orientador: Prof. Dr. Clayton Reginaldo Pereira

BAURU
Janeiro/2023

P637t Pilan, Vinícius de Paula
Técnicas de Inteligência Artificial para diagnóstico de Acidente Vascular Cerebral através de imagens e dados textuais sobre possíveis vítimas / Vinícius de Paula Pilan. -- , 2023
58 p. : il., tabs., fotos

Trabalho de conclusão de curso (-) - Universidade Estadual Paulista (Unesp), Faculdade de Ciências, Bauru,
Orientador: Prof. Dr. Clayton Reginaldo Pereira

1. Inteligência Artificial. 2. Acidente Vascular Cerebral. 3. Aprendizado de Máquina. 4. Aprendizado Profundo. 5. Rede Neural Convolutacional. I. Título.

Sistema de geração automática de fichas catalográficas da Unesp. Biblioteca da Faculdade de Ciências, Bauru. Dados fornecidos pelo autor(a).

Essa ficha não pode ser modificada.

Vinícius de Paula Pilan

TÉCNICAS DE INTELIGÊNCIA ARTIFICIAL PARA DIAGNÓSTICO DE ACIDENTE VASCULAR CEREBRAL ATRAVÉS DE IMAGENS E DADOS TEXTUAIS SOBRE POSSÍVEIS VÍTIMAS

Trabalho de Conclusão de Curso do Curso de Ciência da Computação da Universidade Estadual Paulista "Júlio de Mesquita Filho", Faculdade de Ciências, Campus Bauru.

Banca Examinadora

Prof. Dr. Clayton Reginaldo Pereira

Orientador

Universidade Estadual Paulista "Júlio de
Mesquita Filho"

Faculdade de Ciências

Departamento de Computação

Prof. Dra. Simone Prado

Universidade Estadual Paulista "Júlio de
Mesquita Filho"

Faculdade de Ciências

Departamento de Computação

**Prof. Dr. Kelton Augusto Pontara da
Costa**

Universidade Estadual Paulista "Júlio de
Mesquita Filho"

Faculdade de Ciências

Departamento de Computação

Bauru, 27 de Janeiro de 2023.

Agradecimentos

Agradeço a Deus por todas as oportunidades que tive acesso, pelas conquistas que obtive e pelas experiências que vivi até o presente momento em minha vida, que me fizeram ser quem sou hoje.

Agradeço imensamente aos meus pais, Adriana e Pedro, e a minha namorada, Isabella, que sempre estiveram presentes ao meu lado, celebrando comigo os momentos felizes e me dando apoio nos momentos difíceis.

Agradeço ao Prof. Dr. Clayton Reginaldo Pereira, que aceitou ser meu orientador neste trabalho, que me guiou e forneceu todo o suporte necessário para que meu trabalho de conclusão de curso pudesse ser realizado.

E, por fim, agradeço também às minhas amizades que fiz durante todo o período da graduação e aos demais professores, colegas e todos que, de alguma forma, contribuíram para meu desenvolvimento pessoal e profissional nessa ótima jornada de aprendizado.

Resumo

O Acidente Vascular Cerebral (AVC) é uma das doenças que mais matam e incapacitam no mundo todo e quanto mais tardio é o seu diagnóstico, maiores podem ser os prejuízos para a vítima. Portanto, formas de agilizar e auxiliar o processo de diagnóstico da doença podem ser bastante relevantes e benéficas. Este trabalho abordou essa questão, aplicando técnicas de Aprendizado de Máquina e Aprendizagem Profunda, sendo ambas subáreas da Inteligência Artificial, a qual vem sendo amplamente aplicada em diversos segmentos na busca de otimizar tarefas, principalmente no diagnóstico de doenças. Modelos classificadores de fatores de risco foram criados a partir da aplicação dos modelos de Regressão Logística e Floresta Aleatória, assim como uma Rede Neural Convolucional (CNN) para classificação de imagens de tomografia computadorizada da região cerebral.

Palavras-chave: Inteligência Artificial, Acidente Vascular Cerebral, Aprendizado de Máquina, Aprendizagem Profunda, Rede Neural Convolucional.

Abstract

Stroke (AVC) is one of the most deadly and disabling diseases worldwide and a later diagnostic can be more harmful to the victim. Therefore, ways to speed up and help the process of diagnosing the disease can be very relevant and beneficial, and this work addressed this issue, applying Machine Learning and Deep Learning techniques, both of which are subareas of Artificial Intelligence, which has been widely applied in several segments in the search to optimize tasks, mainly in the diagnosis of diseases. Risk factor classification models were created based on Logistic Regression and Random Forest models application, as well as a Convolutional Neural Network (CNN) for classification of computed tomography images of the brain region.

Keywords: Artificial intelligence, Stroke, Machine Learning, Deep Learning, Convolutional Neural Networks.

Lista de figuras

Figura 1 – Relação entre Inteligência Artificial (<i>Artificial Intelligence</i>), Aprendizado de Máquina (<i>Machine Learning</i>) e Aprendizado Profundo (<i>Deep Learning</i>). . .	14
Figura 2 – Gráfico de uma função logística $f(x)$ sendo utilizada para divisão de duas classes a partir do valor 0.5 (limiar entre as classes)	20
Figura 3 – Exemplo de uma Árvore de Decisão para um conjunto de dados que contenha as variáveis Idade, Problema Cardíaco, Alcoolismo, Gênero e Tabagismo, com as possíveis rotulações de classe não AVC e AVC.	21
Figura 4 – Estrutura de um modelo de Floresta Aleatória para m variáveis, contendo n Árvores de Decisão com duas possibilidades de classificação: Classe 01 e Classe 02.	22
Figura 5 – Ilustração de um neurônio biológico	23
Figura 6 – Ilustração de um neurônio artificial simples	23
Figura 7 – Ilustração simplificada de uma Rede Neural Artificial simples	25
Figura 8 – Representação de uma rede da arquitetura de Rede Neural Convolucional VGG-16	27
Figura 9 – Exemplo de matriz de confusão para um classificador binário em que as classes são rotuladas como 0 (negativo) e 1 (positivo). No eixo vertical tem-se a rotulação 0 ou 1 para os valores previstos pelo modelo e no eixo horizontal, a rotulação dos valores reais	29
Figura 10 – Exemplo Curva ROC (em laranja) constituída a partir da variação de TPR e FPR do modelo de Regressão Logística criado. A curva azul é uma curva de referência, indicando pontos de TPR e FPR de valores iguais.	30
Figura 11 – Exemplo de correção para a coluna com apenas duas opções de valores (variável binária).	39
Figura 12 – Exemplo de aplicação da técnica <i>One-Hot Encoding</i> para correção de uma coluna com quatro opções de valores diferentes.	39
Figura 13 – Gráfico de distribuição dos atributos contínuos presentes no conjunto de dados com a presença de anomalias, representadas por círculos pretos. . . .	40
Figura 14 – Gráfico de distribuição dos atributos contínuos presentes no conjunto de dados sem a presença de anomalias, após as alterações realizadas nos valores. . . .	41
Figura 15 – Ilustração do processo de validação cruzada, com cinco dobras, feito para treinamento e avaliação dos dois classificadores criados. A junção das cinco dobras é o conjunto de dados total utilizado.	42
Figura 16 – Matriz de confusão para o modelo de Regressão Logística, criado a partir da média dos valores obtidos nos testes da validação cruzada.	50

Figura 17 – Matriz de confusão para o modelo de Floresta Aleatória, criado a partir da média dos valores obtidos nos testes da validação cruzada. 50

Figura 18 – Matriz de confusão com os valores gerados na classificação feita pela rede a partir do conjunto de dados de treino, com valores conhecidos pelo classificador. 53

Figura 19 – Matriz de confusão com os valores gerados na classificação feita pela rede a partir do conjunto de dados de validação, com valores desconhecidos pelo classificador. 53

Lista de tabelas

Tabela 1 – Frequência de cada possível valor dos atributos categóricos presentes no conjunto de dados.	38
Tabela 2 – Distribuição das imagens do conjunto de dados antes da nova rotulação. .	44
Tabela 3 – Distribuição das imagens do conjunto de dados depois da nova rotulação. .	44
Tabela 4 – Alterações feitas para aumento dos dados do conjunto de imagens utilizado.	45
Tabela 5 – Estrutura final da Rede Neural Convolutacional criada para imagens RGB de dimensão 512x512 pixels.	47
Tabela 6 – Valor médio e desvio padrão das acurácias calculadas em cada iteração de teste da validação cruzada, para ambos os classificadores criados.	49
Tabela 7 – Média e desvio padrão da taxa de valores verdadeiros positivos gerados em cada iteração de teste da validação cruzada.	51
Tabela 8 – Média e desvio padrão da taxa de valores falsos positivos gerados em cada iteração de teste da validação cruzada.	51
Tabela 9 – Média e desvio padrão dos valores de AUC ROC calculados em cada iteração de teste da validação cruzada para ambos os modelos.	52
Tabela 10 – Acurácia alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.	52
Tabela 11 – Taxa de valores verdadeiros positivos alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.	54
Tabela 12 – Taxa de valores falsos positivos alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida	54
Tabela 13 – AUC ROC calculadas na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.	54

Lista de abreviaturas e siglas

AVC	Acidente Vascular Cerebral.
ACVh	Acidente Vascular Cerebral Hemorrágico.
ACVi	Acidente Vascular Cerebral Isquêmico.
TC	Tomografia computadorizada.
AUC	<i>Area Under Curve.</i>
KNN	K-Vizinhos Próximos.
SVM	Máquina de Vetores de Suporte.
ANN	<i>Artificial Neural Network.</i>
CNN	<i>Convolutional Neural Network.</i>
IA	Inteligência Artificial.
IQR	<i>Interquartile range.</i>
FN	<i>False negative.</i>
TN	<i>True negative.</i>
FP	<i>False positive.</i>
TP	<i>True positive.</i>
FPR	<i>False positive rate.</i>
TPR	<i>True positive rate.</i>
ROC	<i>Receiver Operating Characteristic.</i>

Sumário

1	INTRODUÇÃO	13
2	FUNDAMENTAÇÃO TEÓRICA	16
2.1	Acidente Vascular Cerebral – AVC	16
2.1.1	Fatores de risco da doença	16
2.1.2	Diagnóstico com tomografias computadorizadas (TC)	17
2.2	Aprendizado de Máquina	18
2.2.1	Aprendizado Supervisionado	18
2.2.1.1	Regressão Logística	19
2.2.1.2	Floresta Aleatória	20
2.3	Aprendizagem profunda	22
2.3.1	Neurônios artificiais	22
2.3.2	Redes Neurais Artificiais	24
2.3.3	Redes Neurais Convolucionais	25
2.4	Métricas para avaliação do desempenho de modelos classificadores	27
2.4.1	Acurácia	28
2.4.2	Matriz de confusão	28
2.4.3	Taxa de FP, taxa de TP, Curva ROC e AUC ROC	29
3	TRABALHOS CORRELATOS	31
4	METODOLOGIA	33
4.1	Base de dados para os fatores de risco	33
4.1.1	Variáveis categóricas	33
4.1.2	Variáveis quantitativas	34
4.2	Base de dados para as imagens de tomografia computadorizada	34
4.3	Ferramentas utilizadas	34
4.3.1	Python	34
4.3.2	Google Colab	35
4.3.3	Bibliotecas utilizadas da linguagem	35
5	DESENVOLVIMENTO	37
5.1	Classificador para fatores de risco	37
5.1.1	Pré-processamento: preparação dos dados	37
5.1.2	Modelagem: criação do classificador	41
5.2	Rede Neural classificadora de imagens de tomografia computadorizada	43

5.2.1	Correção da rotulação dos dados e criação dos conjuntos de treino e de validação	43
5.2.2	Pré-processamento: redimensionamento dos valores e aumento dos dados do conjunto de dados original	44
5.2.3	Estrutura da Rede Neural Convolucional desenvolvida	45
5.2.4	Treinamento da rede estruturada	47
6	RESULTADOS	49
6.1	Resultados dos modelos criados para classificação de fatores de risco	49
6.2	Resultados da rede neural criada para classificação de imagens de tomografia computadorizada	52
7	CONCLUSÃO	55
	REFERÊNCIAS	57

1 Introdução

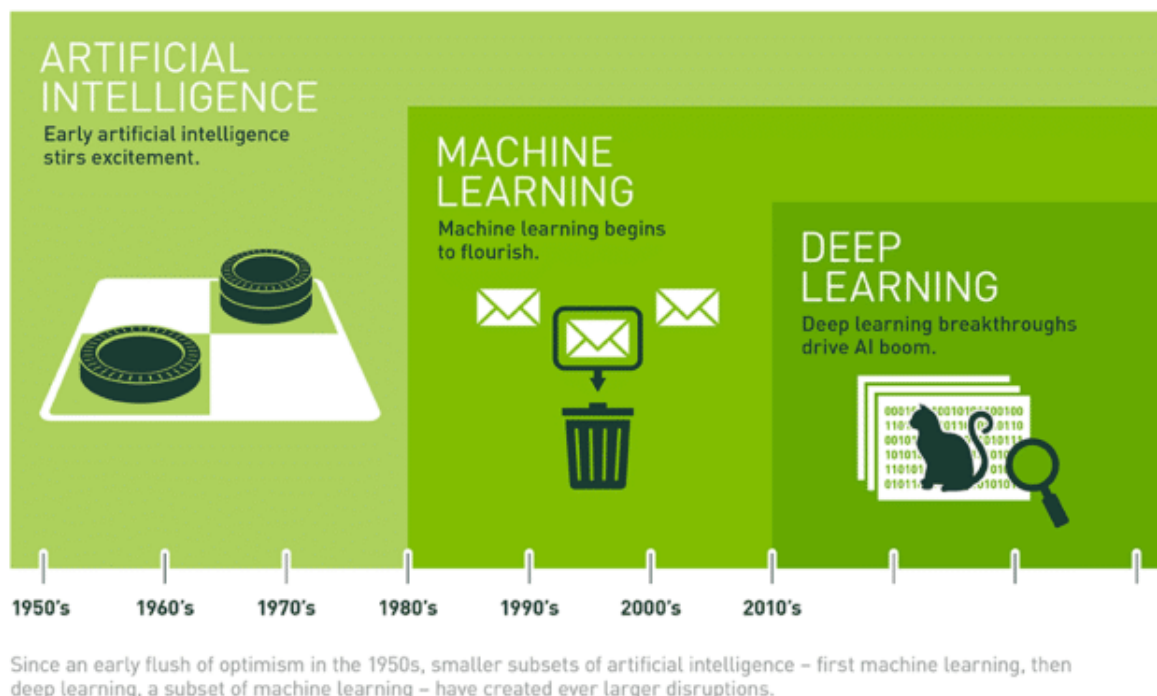
O Acidente Vascular Cerebral, também chamado de derrame cerebral ou simplesmente AVC, é uma das doenças que mais causa mortes, incapacitações e internações no mundo de acordo com [OLIVEIRA e ANDRADE \(2001\)](#). A doença ocorre em virtude da alteração do fluxo de sangue que chega ao cérebro. Quando acontece uma interrupção ou diminuição drástica do suprimento sanguíneo que chega na região cerebral, células nervosas dessa região são atingidas e ficam sem o abastecimento necessário de oxigênio e de nutrientes, o que as levam à morte. [Saúde \(2015\)](#) explica que tal alteração do fluxo sanguíneo pode ser causada por obstrução ou ruptura de vasos sanguíneos da região .

Quanto mais tardio é realizado o diagnóstico e tratamento, maiores são os prejuízos e sequelas, já que as células vão morrendo pela falta de abastecimento sanguíneo. Dessa forma, é de grande importância facilitar e agilizar o diagnóstico da doença, para que, assim, as consequências e sequelas sejam evitadas/diminuídas ao máximo, auxiliando na redução da quantidade de casos com maior gravidade e até mesmo casos de óbitos.

Desenvolver formas para otimizar o processo de diagnóstico do AVC no menor tempo possível pode reduzir consideravelmente as sequelas causadas no indivíduo, facilitando muito no processo de reabilitação. Utilizando a computação, especificamente as áreas relacionadas a Inteligência Artificial, ou IA, é possível sintetizar sistemas inteligentes que consigam auxiliar no diagnóstico de casos da doença.

Modelos de Aprendizado de Máquina e Aprendizagem Profunda, subárea da IA e subárea do próprio Aprendizado de Máquina, respectivamente, como mostrado na Figura 1, possuem boa capacidade de reconhecimento de padrões, de acordo com [Tiwari, Tiwari e Tiwari \(2018\)](#), e possibilitam a criação de modelos classificadores. Com a ajuda de um classificador de fatores de risco, por exemplo, um indivíduo conseguiria ter a noção básica de sua situação atual com relação ao AVC e procurar auxílio médico de um especialista o mais rápido possível.

Figura 1 – Relação entre Inteligência Artificial (*Artificial Intelligence*), Aprendizado de Máquina (*Machine Learning*) e Aprendizado Profundo (*Deep Learning*).



Fonte: [COPELAND \(2021\)](#)

O objetivo do presente trabalho de conclusão de curso desenvolvido está diretamente relacionado a utilização de técnicas de Aprendizado de máquina e aprendizado profundo para possibilitar a classificação da situação de um indivíduo com relação a um possível quadro de AVC. Baseado nas informações do indivíduo, seus fatores de risco e, também, imagens de tomografia computadorizada (TC), será possível alertá-lo, através do processamento de suas informações sobre fatores de risco e imagens de tomografia computadorizada. o quanto antes sobre uma possível suspeita. Isso auxilia o médico no processo de diagnóstico, que irá tomar as devidas providências, e agiliza ainda mais o processo.

Vale ressaltar que o intuito da utilização de métricas inteligentes junto ao diagnóstico de doenças não busca, em hipótese alguma, substituir qualquer opinião médica e muito menos algum profissional de saúde, mas apenas informar e conscientizar indivíduos para que busquem auxílio médico o mais rápido possível em caso de suspeita e agilizar o processo de diagnóstico da doença, auxiliando na tomada de decisões.

Baseado nessas informações, para esse trabalho de conclusão de curso (TCC), foram aplicados dois modelos classificadores de Aprendizado de Máquina para classificação de casos de AVC através da análise de fatores de risco. A aplicação desses modelos apresentou resultados eficientes e bastante similares. Além disso, foi construída uma rede neural, do tipo convolucional,

que consegue classificar imagens de tomografia computadorizada, rotulando-as como sendo imagens da região cerebral de um indivíduo de condição patológica vítima de AVC ou de condição normal, através do processamento dessas imagens.

2 Fundamentação Teórica

2.1 Acidente Vascular Cerebral – AVC

O AVC, causado pela alteração do fluxo de sangue do cérebro, pode ser originado da obstrução total ou parcial de alguma artéria da região, o que é conhecido como Acidente Vascular Cerebral Isquêmico (AVCi), ou simplesmente infarto cerebral. Esse tipo de AVC é o mais comum entre os casos da doença e pode ocorrer devido a um trombo (caso de trombose) ou êmbolo (caso de embolia) presente na vítima.

Outra origem da alteração do fluxo de sangue mencionada é o caso do rompimento de determinado vaso sanguíneo do cérebro, o que causa uma hemorragia na região. Esse caso caracteriza o Acidente Vascular Cerebral Hemorrágico (AVCh). Tal rompimento de determinado vaso altera o nível da pressão intracraniana e pode dificultar também a chegada de sangue em outras áreas não afetadas. Embora seja menos comum, esse tipo de AVC é o mais grave, tendo maiores índices de mortalidade quando comparado ao outro tipo abordado (AVC isquêmico), segundo [Saúde \(2015\)](#).

Com relação aos sintomas presentes nas vítimas, dor de cabeça forte sem causa aparente; fraqueza ou formigamento da face, braço ou perna (especialmente em apenas um dos lados do corpo); alteração da fala e/ou da compreensão; alteração da visão de um ou ambos os olhos; e dificuldade, ou incapacidade, de se movimentar são sintomas bastante comuns em ambos os tipos da doença.

Em específico para o caso de AVCi, tem-se de forma comum os sintomas de tontura e perda de equilíbrio ou coordenação. Já para o caso de AVCh, pode ocorrer também náusea, vômito, confusão mental e, em alguns casos, perda de consciência. Em tal caso também podem aparecer os sintomas de sonolência exagerada, alterações na frequência cardíaca e respiratória, e até mesmo convulsões. Vale ressaltar que é de grande importância que um indivíduo com os sintomas apresentados procure assistência médica o mais rápido possível. Dessa forma, a doença pode ser diagnosticada e tratada rapidamente, possivelmente resultando em maiores chances de sobrevivência da vítima e menores sequelas, de acordo com [Saúde \(2015\)](#).

2.1.1 Fatores de risco da doença

Existem alguns fatores que aumentam a probabilidade da ocorrência do acidente vascular cerebral, facilitando o desencadeamento da doença. Estes são denominados fatores de risco e podem ser inerentes à vida humana, maus hábitos, estilo de vida inadequado ou até mesmo questões genéticas. Os principais fatores são:

- Hipertensão
- Diabetes
- Obesidade
- Tabagismo
- Consumo excessivo e frequente de álcool e drogas
- Estresse
- Idade avançada (envelhecimento)
- Histórico familiar
- Sexo masculino
- Colesterol elevado
- Doenças cardiovasculares (principalmente as que produzem arritmia cardíaca)
- Sedentarismo
- Doenças do sangue

2.1.2 Diagnóstico com tomografias computadorizadas (TC)

Além da análise de fatores de risco e sintomas, outra forma de identificar o quadro de acidente vascular cerebral em vítimas da doença é através da análise de tomografias computadorizadas. Os quadros de AVC podem apresentar características aparentes quando analisados em tais tomografias, permitindo a identificação da doença e o nível de ocorrência da mesma na região cerebral de determinado indivíduo. Embora haja certa complexidade no diagnóstico do AVC via TC, é possível aplicar a computação nesse tipo de situação, pois modelos de Inteligência Artificial possuem alta capacidade de reconhecimento de padrões e isso pode auxiliar o profissional de saúde a realizar o diagnóstico da doença.

De forma simplificada, a hipodensidade, presença de áreas hipodensas (mais escuras), é a representação radiológica nas TCs de um edema citotóxico¹, o qual ocorre devido ao mal funcionamento da bomba de sódio/potássio, causando uma área acinzentada mais escura que o habitual no córtex. O corpo humano resolve esse edema citotóxico deixando uma cicatriz na região, denominada de malácia cerebral.

Todo esse processo gera manchas na tomografia e essa diferença de escala de cinza pode ser um dos fatores para classificar se o quadro analisado é um AVC ou não, pois, dependendo

¹Retenção de água e sódio no interior das células cerebrais em razão de uma anomalia no metabolismo celular. Fonte: <https://saude.ccm.net/faq/4569-edema-cerebral-tipos-e-causas>. Acesso em: 10 Dezembro 2022.

do caso, é possível identificar um quadro da doença em fase inicial, fase subaguda ou até mesmo algum quadro de AVC com ocorrência mais antiga na vítima, segundo [Rueda \(2022\)](#).

2.2 Aprendizado de Máquina

O Aprendizado de Máquina, ou Aprendizagem de Máquina, é uma área da inteligência artificial que lida com a criação de programas “inteligentes”, capazes de aprender determinados conceitos e comportamentos, sem serem explicitamente programados para isso, de adquirir conhecimento, de aprender e tomar decisões a partir desse aprendizado. Tal área explora também a capacidade de programas e sistemas conseguirem realizar previsões sobre determinadas situações, aprendendo através da análise de dados já existentes sobre tais situações.

Envolvendo conceitos estatísticos, reconhecimento de padrões e programação, a área de Aprendizado de Máquina possibilita que resultados sejam especulados de acordo com acontecimentos passados semelhantes, através de um processamento dos dados usando o reconhecimento de padrões. "Um sistema de aprendizado é um programa de computador que toma decisões baseado em experiências acumuladas através da solução bem sucedida de problemas anteriores." ([MONARD; BARANAUSKAS, 2003](#), p. 39).

Com relação à criação de uma solução utilizando conceitos de Aprendizado de Máquina, o processo geralmente envolve duas etapas, o treinamento e a avaliação. A etapa de treinamento é quando um determinado conjunto de dados é passado para o modelo processá-lo, aprendendo a reconhecer padrões e também a fazer previsões. A etapa de avaliação consiste em passar dados semelhantes ao determinado conjunto de dados utilizado no treino para que o desempenho do modelo seja medido, analisando os resultados esperados com os resultados obtidos.

No campo do Aprendizado de Máquina, existem basicamente três tipos diferentes de aprendizagem: aprendizado supervisionado, aprendizado não supervisionado e aprendizado por reforço. Os algoritmos de Aprendizado de Máquina utilizados neste trabalho pertencem à área de aprendizado supervisionado e, portanto, apenas este será abordada com mais detalhes.

2.2.1 Aprendizado Supervisionado

O Aprendizado supervisionado utiliza dados rotulados no treinamento do modelo, o qual prevê uma variável dependente a partir de uma ou mais variáveis independentes, conforme explica [Honda, Facure e Yaohao \(2017\)](#). Existem diversos algoritmos desse tipo de aprendizado, sendo os principais: regressão linear, regressão logística, máquina de suporte vetorial, árvores de decisão, k-vizinhos mais próximos, floresta aleatória, entre outros. Os utilizados neste trabalho para desenvolver o classificador de fatores de risco foram apenas a Regressão Logística e o Classificador de Floresta Aleatória.

O modelo de Regressão Logística foi escolhido pela sua popularidade na área de

Aprendizado de Máquina e por conseguir ter capacidade de calcular, de forma eficiente, a probabilidade de determinado evento pertencer a alguma das classes do problema, e não só simplesmente classificar a entrada processada. O algoritmo de Floresta Aleatória foi escolhido pela sua popularidade na área, mas também pelo fato de ser um modelo que tem seu funcionamento bastante diferente de como o modelo de Regressão Logística processa os dados que recebe.

2.2.1.1 Regressão Logística

De acordo com [IBM \(2022b\)](#), a Regressão Logística é um modelo estatístico capaz de modelar determinada probabilidade de um evento acontecer a partir da combinação linear de variáveis independentes entre si associadas a tal evento. Muito utilizado em situações de classificação binária, esse modelo combina variáveis independentes com respectivos pesos atribuídos a elas de modo a obter um valor a partir deste processo. Tal valor é passado para uma função logística que resulta na probabilidade de algum evento acontecer. No caso de uma classificação, tal resultado indica a probabilidade de os valores analisados pertencerem a determinada classe do problema em questão.

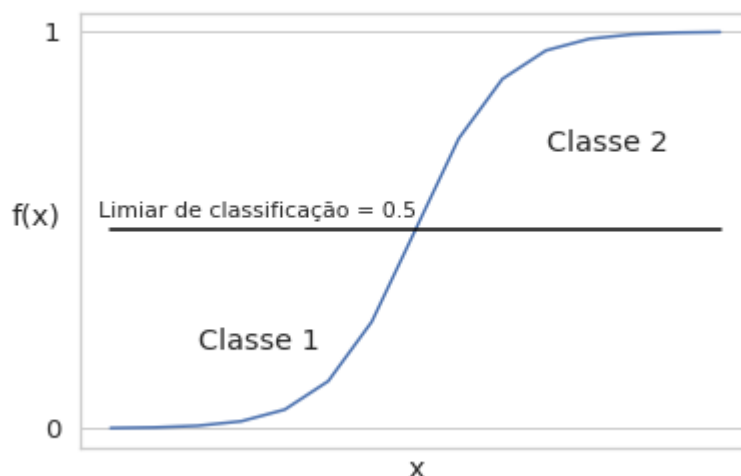
De maneira geral e simplificada, algoritmos de regressão ajustam uma curva fazendo com que variáveis independentes consigam gerar, através de uma função, determinado valor resultante, denominado de variável dependente. Segundo [Scikit-learn \(entre 2007 e 2022a\)](#), o algoritmo de Regressão Logística, embora seja utilizado para classificação de dados, é um algoritmo de regressão o qual ajusta uma curva através da função logística, também conhecida como função sigmoide, com sua estrutura sendo apresentada na Equação 2.1:

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (2.1)$$

Cada variável independente passada para o modelo de Regressão Logística é interpretada como x_i , com i variando de 1 ao número total de variáveis independentes do problema abordado, e é multiplicada por um valor denominado peso. Pesos de um modelo são valores comumente descritos pelas letras θ ou w e são gerados, inicialmente, de forma aleatória, sendo otimizados durante a etapa de treinamento do modelo para gerarem menores erros de classificação (modelo mais eficiente).

Através da combinação linear dos pontos $x * \theta$ (ou $x * w$), é obtida a variável dependente passada para a função logística (equação 1) e o resultado desse processo é um valor que varia de 0 a 1. Esse valor, após comparado com um limiar de classificação, termo de valor geralmente igual a 0.5 para classificações binárias que serve como comparativo para classificação de determinado resultado que varia de 0 a 1, resulta na classificação da entrada recebida como uma das classes pertencentes ao problema binário abordado, como mostrado na Figura 2:

Figura 2 – Gráfico de uma função logística $f(x)$ sendo utilizada para divisão de duas classes a partir do valor 0.5 (limiar entre as classes)



Fonte: Elaborado pelo autor.

2.2.1.2 Floresta Aleatória

Floresta Aleatória é um modelo de aprendizado de máquina composto por um conjunto aleatório de Árvores de Decisão que pode ser usado em situações de regressão ou classificação, por exemplo. Especificamente para problemas de classificação, como o abordado neste trabalho, o modelo classifica de acordo com a classificação feita pelas árvores que o compõem. Dessa forma, se há um resultado classificativo mais frequente, isto é, se a maioria das árvores gerou determinada classificação, então esta será a que o modelo de floresta aleatória irá resultar em sua conclusão, conforme explicado por [Navlani \(2018\)](#).

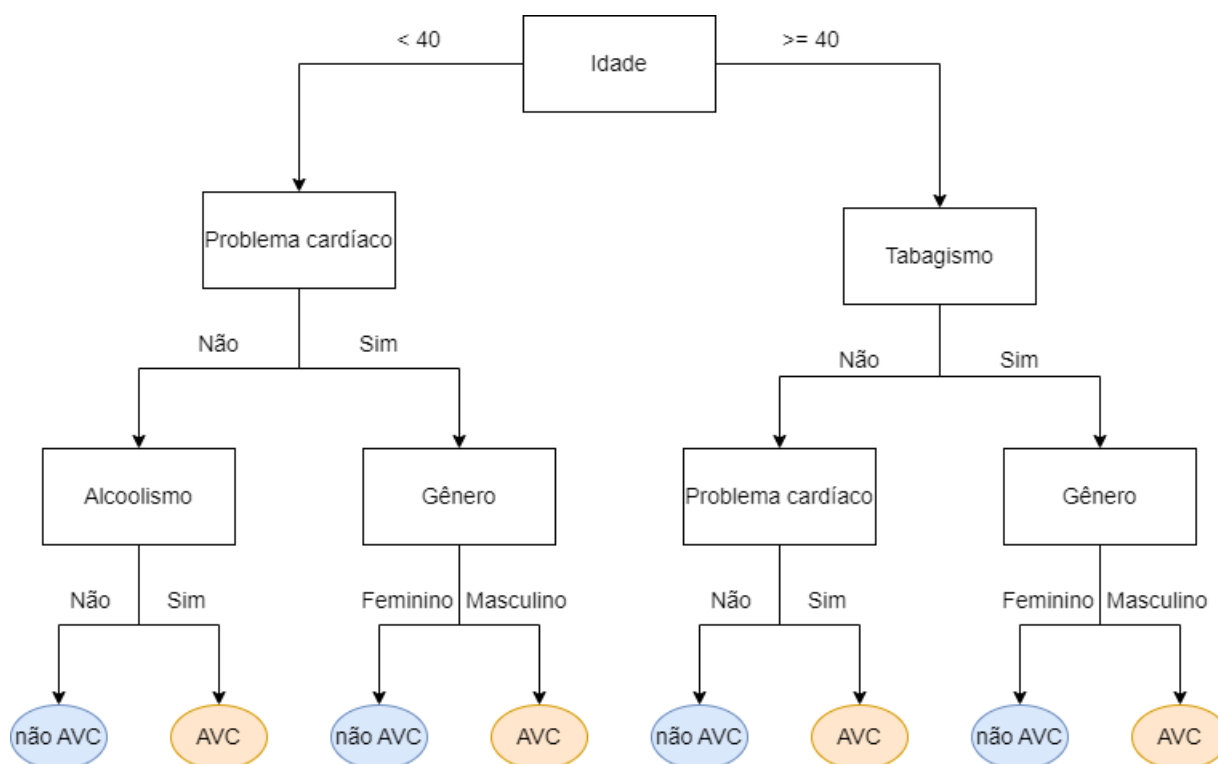
A estrutura de uma Árvore de Decisão é constituída por vários blocos de decisão, comumente denominados de nós. Cada nó funciona como um bloco condicional que encaminha a entrada recebida para um de seus filhos de acordo com o valor que ele analisa (regra de decisão). Tal processo é feito para todo nó pertencente à árvore até chegar em determinado nó folha que, por sua vez, representa uma das possíveis classes do problema.

Portanto, em uma árvore de decisão, uma entrada percorre um dos diferentes caminhos possíveis até chegar em um nó folha que a classifica como uma das possíveis rotulações da situação. Tal caminho acaba sendo definido de acordo com os valores que pertencem a entrada em questão, conforme [IBM \(2022c\)](#) explica. A Figura 3 ilustra a estrutura de uma possível Árvore de Decisão para variáveis ligadas a fatores de risco do AVC.

As regras de decisão utilizadas pelos nós para separar as entradas em diferentes caminhos são estabelecidas na etapa de treinamento do modelo de acordo com determinado conjunto de

dados utilizado nessa fase. Por isso, existem infinitas possibilidades de estruturas para árvores de decisão, sendo cada uma delas definidas de acordo com a situação e ao conjunto de dados abordado.

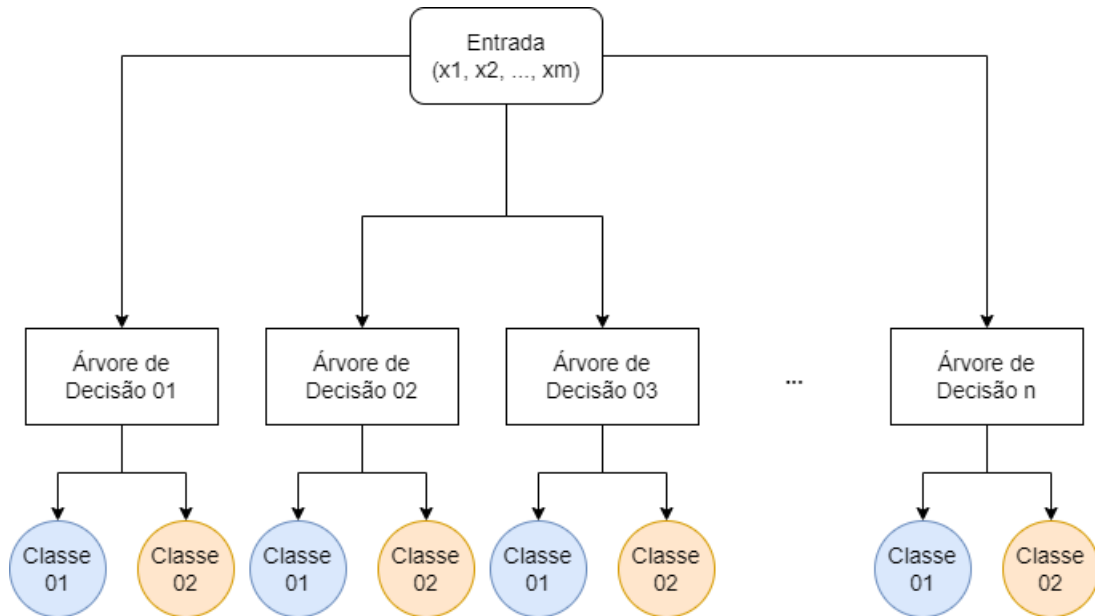
Figura 3 – Exemplo de uma Árvore de Decisão para um conjunto de dados que contenha as variáveis Idade, Problema Cardíaco, Alcoolismo, Gênero e Tabagismo, com as possíveis rotulações de classe não AVC e AVC.



Fonte: Elaborado pelo autor.

Segundo [Scikit-learn \(entre 2007 e 2022b\)](#), a ideia do algoritmo de Floresta Aleatória é justamente utilizar dessas diversas possibilidades de estrutura das Árvore de Decisão e criar assim uma floresta com várias árvores diferentes abordando o mesmo problema, conforme ilustrado na Figura 4, classificando determinada entrada recebida de acordo com a classificação mais frequente das árvores que compõem o modelo.

Figura 4 – Estrutura de um modelo de Floresta Aleatória para m variáveis, contendo n Árvores de Decisão com duas possibilidades de classificação: Classe 01 e Classe 02.



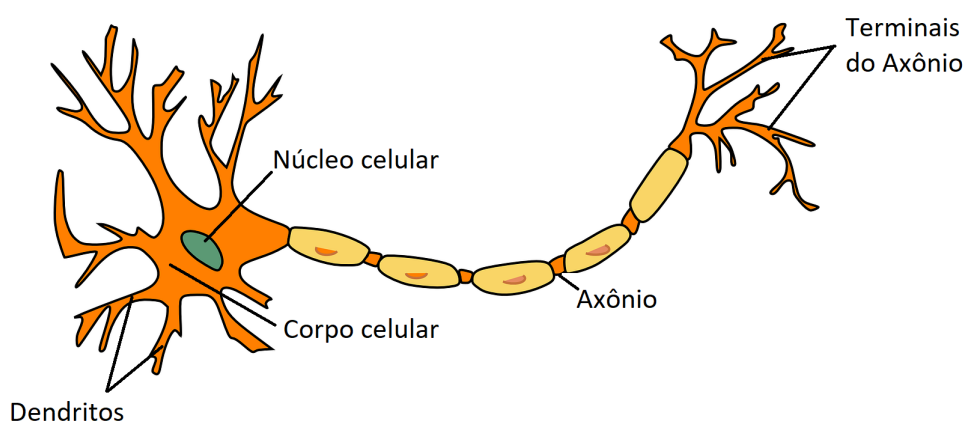
Fonte: Elaborado pelo autor.

2.3 Aprendizagem profunda

2.3.1 Neurônios artificiais

O neurônio artificial é o principal componente das redes neurais computacionais e seu funcionamento é baseado no neurônio real, ilustrado na Figura 5, que compõem o sistema nervoso. De acordo com [Varella \(2020\)](#), os neurônios reais, base do sistema nervoso, são células que estabelecem conexões entre si para transmitir impulsos nervosos pela região cerebral. O neurônio artificial é uma estrutura com funcionamento que se assemelha do neurônio real, mantendo a principal característica dessas estruturas de criarem conexões entre si, recebendo informações provenientes de outros neurônios e passando elas adiante.

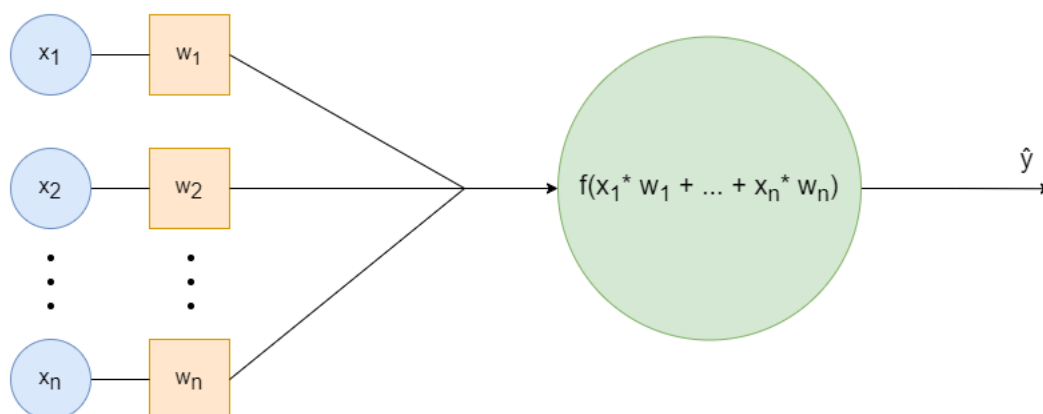
Figura 5 – Ilustração de um neurônio biológico



Fonte: Elaborado pelo autor.

Ilustrados na Figura 6, os principais componentes de um neurônio artificial são as entradas, os pesos, a função de ativação e saída. As entradas, geralmente representadas pela letra x , recebem a informação que chega para o neurônio. Tal informação é multiplicada por um determinado peso, comumente representado pela letra w , e o resultado disso é passado para uma função de ativação que processa os valores recebidos e retorna um ou mais valores, os quais servirão de entrada para um outro neurônio artificial conectado à rede, conforme explicado por [Yacim e Boshoff \(2018\)](#).

Figura 6 – Ilustração de um neurônio artificial simples



Fonte: Elaborado pelo autor.

2.3.2 Redes Neurais Artificiais

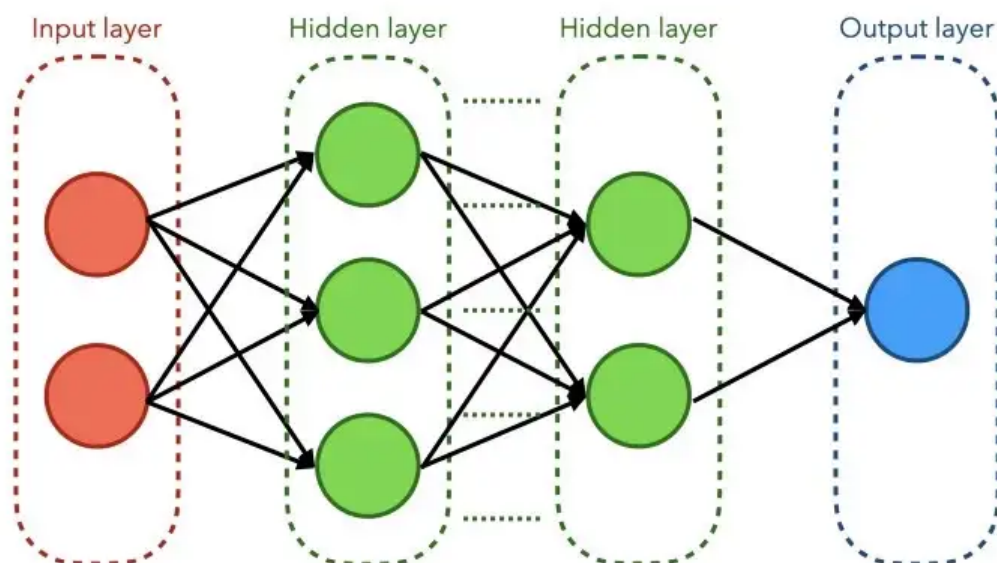
De acordo com [AWS \(2022\)](#), Redes Neurais Artificiais, ou *Artificial Neural Network* (ANN), são estruturas do campo da inteligência artificial que possibilitam os computadores a processar dados se baseando no comportamento do cérebro humano. Ao receber um conjunto de dados específico, uma rede neural artificial é capaz de se ajustar a este conjunto, aprendendo a reconhecer padrões sobre estes dados e se tornando mais eficiente à medida que os processa. A estrutura da rede é composta por um conjunto de neurônios artificiais agrupados em camadas interconectadas que possibilitam o computador a aprender, reconhecer padrões e tomar decisões inteligentes.

Como ilustrado na Figura 7, uma rede neural simples possui em sua estrutura camadas de três tipos: camada de entrada (*Input layer*), camada oculta (*Hidden layer*) e camada de saída (*Output layer*). A camada de entrada é onde se encontram os neurônios responsáveis por receber os dados que serão processados pela rede. Tal camada faz o primeiro processamento desses dados e os encaminham para o próximo conjunto de neurônios. As camadas do tipo oculta recebem dados da camada de entrada ou de outras camadas ocultas, processando novamente os dados e passando o resultado para outra camada. Redes neurais podem ter uma ou mais camadas desse tipo.

Por fim, o conjunto de neurônios da camada de saída faz o último processamento dos dados, fornecendo o resultado final de todos os dados processados pela rede. Para problemas de classificação binária, tem-se apenas um nó na camada de saída, o qual indica a probabilidade de determinada entrada analisada pela rede ser ou não da classe abordada. Já para problemas de várias classes diferentes, tem-se vários nós nessa camada, conforme também abordado por [AWS \(2022\)](#).

Além dessa estrutura de rede abordada, existem outros tipos mais específicos de redes neurais que apresentam melhores resultados quando aplicados em determinadas situações. No caso de problemas que envolvem a classificação de imagens e o reconhecimento de objetos, o uso de Redes Neurais Convolucionais é mais recorrente do que o uso da rede neural padrão, justamente pela boa capacidade que estas possuem de valorizar características da imagem.

Figura 7 – Ilustração simplificada de uma Rede Neural Artificial simples



Fonte: [Nurfikri \(2020\)](#)

2.3.3 Redes Neurais Convolucionais

Comumente chamadas de Conv-Net ou simplesmente CNN (acrônimo de *Convolutional Neural Network*), as redes neurais do tipo convolucional são muito utilizadas em tarefas que envolvem a identificação de objetos, reconhecimento de características e classificação de imagens. Esse tipo de rede tem a capacidade de valorizar detalhes da imagem que são relevantes para a classificação, realçando determinados agrupamentos de pixels através da aplicação de filtros matriciais na imagem para identificação de padrões, o que colabora para um processo de classificação mais eficiente.

As Redes Neurais Convolucionais apresentam três tipos principais de camadas em sua estrutura: camadas de convolução (*convolution*), camadas de agrupamento (*pooling*) e as camadas totalmente conectadas. Segundo [O'Shea e Nash \(2015\)](#), camadas de convolução, ou convolucionais, são as camadas responsáveis pelo processo de realce das características da imagem, mencionado no parágrafo anterior.

Esse processo se dá pela aplicação de um filtro, comumente chamado de *kernel*, na imagem, vista como uma matriz de pixels, que está sendo processada pela rede, funcionando como um detector de recursos. Esse filtro tem formato matricial com tamanho típico de 3x3, podendo ser maior dependendo da rede, e o processo de aplicação se dá de forma iterativa via a realização do produto escalar entre seus valores e os da matriz da imagem processada, finalizando-se com a soma dos resultados dessa multiplicação, como explica [IBM \(2022a\)](#).

Tal processo é iterativo, acontecendo até que o filtro multiplique toda a imagem. Alguns

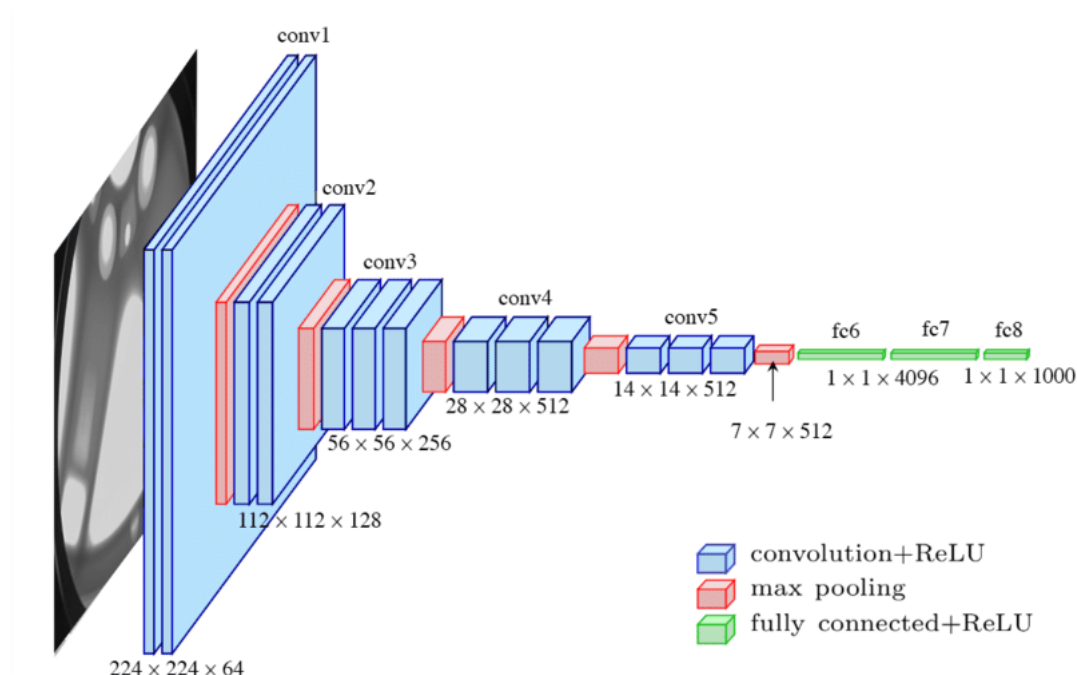
hiperparâmetros que podem variar o volume da imagem resultante são: o número de filtros que serão aplicados a imagem, o número de passos denominado *stride* que o *kernel* percorre para se deslocar de uma multiplicação para a outra e o *padding*, que adiciona zeros na matriz da imagem para chegar a determinado tamanho que favoreça a aplicação do *kernel*.

As camadas de agrupamento irão pegar o resultado da camada de convolução e reduzir o número de valores recebidos, mantendo as características realçadas pelos filtros da convolução, porém eliminando áreas vazias de informação relevante para a classificação. O processo de pooling destas camadas se dá por dois tipos diferentes que podem ser utilizados, o agrupamento máximo ou o agrupamento médio.

O agrupamento máximo é feito através da aplicação de um filtro que se desloca pela imagem que simplesmente pega o maior valor dentre os valores analisados na imagem e o coloca em uma matriz de saída, descartando os demais valores. O agrupamento médio se dá pela aplicação de um filtro que também se move pela imagem analisada, porém calculando o valor médio dos valores analisados pelo filtro, colocando-o em uma matriz de saída de forma semelhante ao outro agrupamento descrito. O agrupamento máximo tem maior frequência de utilização quando comparado com o agrupamento médio, de acordo com [IBM \(2022a\)](#).

Por fim, nas camadas totalmente conectadas tem-se diversos neurônios totalmente conectados que irão processar os valores recebidos aplicando alguma função de ativação, gerando valores de saída para as próximas camadas ou para a camada de saída. Toda rede neural convolucional tem seu formato base composto por essas camadas, porém existem diversos tipos de CNN diferentes que vão variar o número de instâncias dessas camadas mencionadas e os hiperparâmetros presentes nas mesmas. A arquitetura de rede convolucional denominada de VGG-16, utilizada no desenvolvimento deste trabalho e ilustrada na Figura 8, apresenta treze camadas de convolução e cinco camadas de agrupamento, assim como três camadas totalmente conectadas, segundo [JORDAN \(2018\)](#).

Figura 8 – Representação de uma rede da arquitetura de Rede Neural Convolutacional VGG-16



Fonte: [Ferguson et al. \(2017\)](#)

2.4 Métricas para avaliação do desempenho de modelos classificadores

Com o modelo classificador ajustado de acordo com a situação abordada, torna-se fundamental a realização de uma boa avaliação de seu desempenho. “O objetivo da fase de avaliação é estimar os resultados de mineração de dados de forma rigorosa e obter a confiança de que são válidos e confiáveis antes de avançar.” ([PROVOST; FAWCETT, 2016](#), p. 31). Avaliar de forma adequada permite a identificação de problemas, como sobreajuste ou excesso de generalização, por exemplo, assim como visualizar de forma geral como seria a performance do modelo quando colocado em produção na prática.

Para realizar o processo de avaliação de modelos classificadores, existem diversas métricas que abordam diferentes pontos da classificação, como quantidade de erros, acertos, de classificações falsas negativas e positivas, entre outras. Para o presente trabalho, foram utilizadas as métricas de Acurácia, Matriz de Confusão, Taxa de Falso Positivo (FPR), Taxa de Verdadeiro Positivo (TPR) e a Área da Curva de Característica de Operação do Receptor (AUC ROC).

2.4.1 Acurácia

A acurácia é uma das métricas mais simples de se medir e justamente por isso sua utilização é muito popular na avaliação de classificadores. Porém, por ser demasiadamente simples, tal métrica pode esconder deficiências do modelo analisado e transparecer um resultado não muito confiável. O cálculo da acurácia é dado pela divisão do número de acertos do modelo pelo número total de casos analisados, ou seja, pelo número de acertos somados com o número de erros, conforme explicado por [Google \(2022\)](#) e apresentado na Equação 2.2:

$$acuracia = \frac{acertos}{acertos + erros} \quad (2.2)$$

Como o cálculo dessa métrica não leva em conta fatores além de acertos e erros, como o tipo de erro, por exemplo, é bastante interessante utilizar a Acurácia em conjunto com outras métricas para fazer uma avaliação mais eficiente, levando em conta também os tipos de valores falsos que o modelo gera.

2.4.2 Matriz de confusão

Na fase de avaliação de um modelo também pode ser interessante olhar para o desempenho do modelo com relação a erros específicos, e a matriz de confusão apresenta de forma simplificada esse comportamento. “Uma matriz de confusão separa as decisões tomadas pelo classificador, tornando explícito como uma classe está sendo confundida com outra. Desta forma, diferentes tipos de erros podem ser tratados separadamente” ([PROVOST; FAWCETT, 2016](#), p. 189).

Na matriz de confusão tem-se dois tipos de acertos e erros, os verdadeiros positivos e negativos, e os falsos positivos e negativos, respectivamente. Valores verdadeiros ocorrem quando o modelo classifica corretamente valores negativos (0) ou positivos (1) e valores falsos ocorrem quando se tem a classificação incorreta desses mesmos valores, negativos ou positivos. Analisar a taxa de falso positivo e de falso negativo permite melhor visualização de qual tipo de erro é predominante na questão, assim como ter uma ideia percentual probabilística de como o modelo em questão reagiria em situações reais com relação aos tipos de erros.

É comumente utilizado para simplificar a nomenclatura das métricas de falsos e verdadeiros (erros e acertos, respectivamente) as seguintes siglas: FN para falso negativo (*false negative*), TP para verdadeiro positivo (*true positive*), TN para verdadeiro negativo (*true negative*) e FP para falso positivo (*false positive*).

A matriz de confusão é uma métrica que permite a análise desses valores de falsos e verdadeiros de uma forma bastante simplificada. Ela é frequentemente utilizada na avaliação de classificadores pois fornece uma visão sobre a confusão entre as classes que um determinado classificador realizou na análise de casos passados para ele. Uma matriz de confusão é uma matriz

de dimensão $n \times n$ sendo que n é o número de classes do problema abordado. A Figura 9 traz um exemplo dessa métrica para análise de um modelo de classificação binária, ou seja, com a classificação de apenas duas classes distintas.

Figura 9 – Exemplo de matriz de confusão para um classificador binário em que as classes são rotuladas como 0 (negativo) e 1 (positivo). No eixo vertical tem-se a rotulação 0 ou 1 para os valores previstos pelo modelo e no eixo horizontal, a rotulação dos valores reais

0	260	15
1	10	220
	0	1

Fonte: Elaborado pelo autor.

2.4.3 Taxa de FP, taxa de TP, Curva ROC e AUC ROC

A taxa de FP, ou FPR, e a taxa de TP, ou TPR, são métricas usadas em conjunto para avaliar e descrever a eficiência de um modelo classificador. A FPR mede a fração dos valores negativos que foram previstos incorretamente, julgados como positivos pelo classificador em questão. Para a realização do cálculo dessa métrica, é utilizada a Equação 2.3:

$$FPR = \frac{FP}{FP + TN} \quad (2.3)$$

Já a TPR, também conhecida como sensibilidade, mede a fração dos valores positivos que foram previstos de forma correta, julgados também como positivos pelo modelo. O cálculo da TPR é feito através da Equação 2.4, conforme explicado por Deepchecks (2022).

$$TPR = \frac{TP}{TP + FN} \quad (2.4)$$

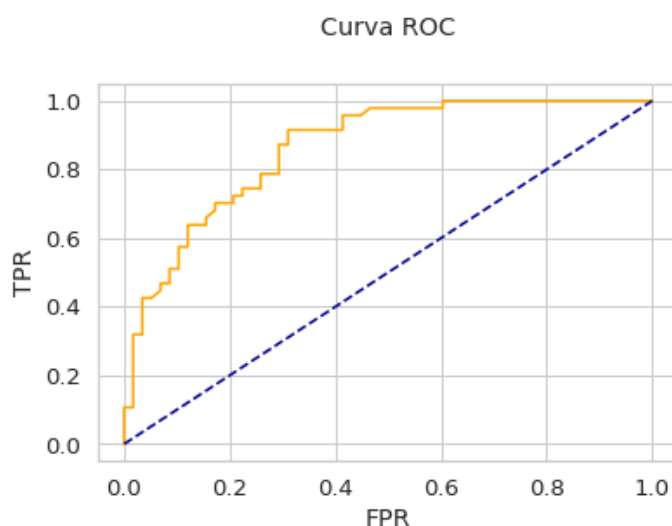
É de grande importância que sejam analisadas em conjunto ambas as métricas descritas, pois assim torna-se possível ter um diagnóstico confiável e saber se o modelo está extremamente tendencioso ou não. Por exemplo, em uma situação em que determinado classificador rotula

sempre todos os dados como sendo de uma classe positiva, sua TPR terá valor máximo (1.0), pois para todo valor positivo recebido, ele acertará todos. Porém, este mesmo modelo, como consequência, errará todos os casos negativos, e isso só poderá ser enxergado se a FPR for analisada, pois seu valor também será 1.0.

Nesse cenário descrito, analisar apenas a TPR indicaria um modelo muito eficiente, já que ele não errou nenhum caso positivo que recebeu. Dessa forma, estaria sendo omitido o alto viés e a baixa capacidade de generalização do classificador em questão. Um modelo tendencioso assim possivelmente apresentaria péssimos resultados na prática.

Para facilitar a análise em conjunto de ambas as métricas, existe uma outra métrica muito poderosa denominada de *Receiver Operating Characteristic Curve*, ou simplesmente *ROC Curve*. Em português, Curva de Característica de Operação do Receptor, a Curva ROC, de acordo com Bagheri (2019), é uma ferramenta gráfica que permite verificar de forma conjunta TPR e FPR à medida que o limiar entre as classes (probabilidade que separa as classes) do modelo varia. Tal variação é feita para diferentes valores do limiar que geram diferentes valores de TPR e FPR como resultado. Estes valores formam pontos em um gráfico TPR x FPR e a junção desses pontos formam a curva ROC, conforme mostrado na Figura 10:

Figura 10 – Exemplo Curva ROC (em laranja) constituída a partir da variação de TPR e FPR do modelo de Regressão Logística criado. A curva azul é uma curva de referência, indicando pontos de TPR e FPR de valores iguais.



Fonte: Elaborado pelo autor.

Outra forma de utilizar a curva ROC para avaliação de determinado modelo é calculando sua área. A partir deste processo, tem-se um número que pode variar entre 0 a 1.0, sendo que, quanto mais próximo de 1.0, mais eficiente o modelo seria, com máximo TPR e mínimo FPR. A área abaixo da curva ROC é uma métrica denominada de Pontuação AUC (*area under curve*) ROC.

3 Trabalhos Correlatos

O tema deste trabalho foi inspirado no trabalho de conclusão de curso ANÁLISE DE DADOS PARA AUXILIAR NO DIAGNÓSTICO PRECOCE DE ACIDENTE VASCULAR CEREBRAL, desenvolvido pela autora [TAMAKE \(2020\)](#) e também orientado pelo Professor Clayton, que orientou o presente trabalho. Tal trabalho correlato também aborda a classificação de casos de AVC através do processamento de dados ligados aos fatores de risco da doença, aplicando alguns modelos de Aprendizagem de Máquina, como feito aqui.

O trabalho de [TAMAKE \(2020\)](#), assim como o presente trabalho, também faz a etapa de pré-processamento, etapa que contém a análise e tratamento dos dados, antes de aplicar os modelos de Aprendizagem de Máquina. O que difere as etapas de pré-processamento dos trabalhos é que [TAMAKE \(2020\)](#) aborda um conjunto de dados com mais registros e menos desbalanceado do que o utilizado neste trabalho e, dessa forma, a autora optou por não realizar nenhum balanceamento nele.

Com relação à criação de fato das soluções (criação dos classificadores de casos de AVC), a etapa de modelagem, em que foram aplicados os modelos de Aprendizagem de Máquina para criar os classificadores de AVC, [TAMAKE \(2020\)](#) utiliza uma divisão simples dos dados, apenas separando 75% dos registros para treinamento dos modelos e 25% para avaliação dos mesmos. No presente trabalho, foi definido como método de divisão dos dados, para treino e teste, a Validação Cruzada, a qual divide o conjunto de dados total em dobras de treino e teste várias vezes, selecionando dobras diferentes do conjunto a cada iteração.

De maneira geral, o principal motivo que diferencia a escolha dos métodos, divisão simples ou Validação Cruzada, de divisão dos dados utilizados em cada um dos trabalhos é que aqui, no atual trabalho, tem-se um conjunto pequeno de dados, o que poderia gerar resultados pouco verídicos na avaliação dos modelos e poderia gerar também treinamentos não eficientes caso o método de divisão dos dados fosse uma divisão simples de treino e teste.

Isso ocorreria, pois, a divisão simples faz a separação aleatória dos dados e, em um conjunto pequeno, como o caso do atual trabalho, diferentes partes do conjunto podem ser mais ou menos significativas para o treinamento do modelo ou para a avaliação do mesmo. Dessa forma, torna-se mais apropriado fazer um processo iterativo de diversas divisões de treino e teste em todo o conjunto, como o que foi feito utilizando a Validação Cruzada, para que o modelo não fique sujeito ao acaso da aleatoriedade de uma divisão simples.

Outro fator que difere o presente trabalho do trabalho correlato citado é que [TAMAKE \(2020\)](#) aplica os modelos de Floresta Aleatória, K-Vizinhos Próximos (KNN) e Máquina de Vetores de Suporte (SVM), enquanto que no atual trabalho foram aplicados os modelos de Regressão Logística e Floresta Aleatória. O modelo de Regressão Logística foi escolhido pois é

um dos modelos mais populares da área de Aprendizado de Máquina em contextos didáticos, apresentando uma forma de classificação das entradas recebidas bastante distinta de como o modelo de Floresta Aleatória processa as informações. O modelo de Floresta Aleatória foi escolhido por conta de sua alta popularidade na área de Aprendizado de Máquina e pela sua capacidade de gerar bons resultados de classificação.

Ademais, o atual trabalho tratou também da classificação de imagens de Tomografia Computadorizada (TC) através de técnicas da área de Aprendizagem Profunda, em que foi construída uma Rede Neural Convolucional para a classificação das imagens de TC analisadas em imagens de vítimas de AVC ou de pessoas em condição normal (não AVC). O trabalho de [TAMAKE \(2020\)](#) faz a classificação dos casos de AVC, como mencionado no início da seção, apenas através do processamento de informações relacionadas a fatores de risco da doença e, portanto, não aborda o contexto de imagens de TC com Redes Neurais.

Todos esses fatores mencionados nesta seção resultaram em novos resultados e em diferentes abordagens do problema em questão, caracterizando um novo trabalho sobre o tema da classificação de casos de Acidente Vascular Cerebral, o qual é um tema de grande importância devido à gravidade e alta ocorrência dessa doença no mundo todo.

Portanto, mesmo que ambos os trabalhos abordem a classificação de casos de AVC, os dados para criação dos modelos dos trabalhos são bases diferentes, assim como as estratégias utilizadas no pré-processamento e no processamento dos dados, em que foram aplicados diferentes modelos classificadores da área de Aprendizado de Máquina e também diferentes métricas para avaliação desses modelos treinados.

Ademais, como mencionado anteriormente, o atual trabalho aborda o contexto de classificação de imagens de TC através de uma Rede Neural, o que está fora do escopo do trabalho correlato mencionado. Por isso, por se tratarem de trabalhos com abordagens diferentes da classificação preditiva de casos de AVC, não foi realizada uma análise comparativa de forma direta dos resultados gerados por ambos os trabalhos.

4 Metodologia

A metodologia deste trabalho se consistiu no primeiro momento a entender mais sobre os assuntos abordados, através de um levantamento bibliográfico sobre os temas de Aprendizado de Máquina e Aprendizado Profundo, estudando os principais modelos e técnicas dessas áreas para criar uma solução mais robusta e eficiente do problema abordado. Para ambos os modelos criados aqui, foram seguidas as etapas: levantamento e pré-processamento, aplicação de técnicas das áreas dos dados que envolvem o problema (fatores de risco e imagens de tomografia computadorizada), processamento de dados para criação do modelo classificador e fase final de avaliação do mesmo, para entender o desempenho do modelo criado.

4.1 Base de dados para os fatores de risco

Os dados utilizados na construção do modelo de Aprendizado de Máquina para fatores de risco foram retirados da plataforma *Kaggle*¹, sendo este um dos conjuntos de dados disponibilizados para propósitos educacionais. O conjunto original conta ao todo com 5110 registros, em que cada um contém informações de uma pessoa diferente, rotulada como vítima de AVC ou não. O conjunto original conta com uma grande quantidade de valores sendo de pessoas rotuladas como não AVC, ou seja, está originalmente desbalanceado.

4.1.1 Variáveis categóricas

Além da variável AVC, o conjunto de dados conta com várias outras variáveis categóricas que são fatores de risco da doença ou que estão ligadas a eles:

- Sexo: Masculino ou feminino, em que o sexo masculino é um fator de risco
- Hipertensão: Indivíduo é ou não hipertenso. Hipertensão é um fator de risco.
- Doença cardíaca: Indivíduo tem ou não alguma doença cardíaca, o que é um fator de risco.
- Estado civil de ser ou já ter sido casado: Sim ou não, tal variável pode estar ligada a estilo de vida e estresse, o qual é um fator de risco da doença.
- Tipo de residência: Urbana ou rural, tal variável está ligada a estilo de vida, que pode ser um fator considerável para a doença

¹<https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset>. Acesso em: 18 Outubro 2022.

- Vítima de AVC: Sim ou não, essa é a variável alvo do trabalho, sendo utilizada para rotular os dados para previsão.
- Tipo de emprego: empregado privado, funcionário público (governo), empreendedor ou não trabalha.
- Status com relação a tabagismo: indivíduo nunca fumou, fuma formalmente, fuma regularmente ou situação desconhecida.

4.1.2 Variáveis quantitativas

Existem também colunas compostas por variáveis quantitativas na base de dados utilizada. Tais variáveis são fatores de risco do AVC (com exceção da coluna ID) e compõem as seguintes colunas:

- Idade: idade avançada é um fator de risco para a doença
- Nível de glicose: diabetes é um fator de risco para a doença
- BMI: Índice de massa corpórea no padrão americano, utilizado para identificar obesidade.
- Id: Identificação de cada registro. Tal variável foi desconsiderada por não ter impacto na variável alvo do trabalho que rotula os casos de AVC.

4.2 Base de dados para as imagens de tomografia computadorizada

Para criação do modelo classificador de imagens de TC, foi utilizado um conjunto de dados com 473 imagens no total, divididas em três classes diferentes: AVC Isquêmico, AVC Hemorrágico ou não AVC. Entretanto, como o objetivo deste trabalho é de apenas identificar a ocorrência de Acidente Vascular Cerebral em determinada vítima, sem especificar o tipo de AVC de fato, este conjunto foi redistribuído para apenas duas classes. Dessa forma, o conjunto original passou a ter apenas as classes *AVC* e *não AVC*:

- *AVC* possui 299 imagens
- *não AVC* possui 174 imagens

4.3 Ferramentas utilizadas

4.3.1 Python

Python² é uma linguagem de programação de alto nível, gratuita e muito utilizada no desenvolvimento de aplicações modernas. Criada por volta de 1990, tal linguagem apresenta

²<https://www.python.org/>. Acesso em: 19 Dezembro 2022.

um uso mais simples quando comparado com o uso de outras linguagens do mercado, trazendo mais facilidade para quem não é necessariamente do ramo da programação. Ela prioriza a legibilidade do código, combinando uma sintaxe clara com recursos poderosos de sua biblioteca padrão (recursos nativos) e de outras bibliotecas desenvolvidas pela comunidade.

Por ser acessível e amigável, a linguagem de programação Python se tornou bastante popular no mercado, com muitos *frameworks* e módulos prontos desenvolvidos e disponibilizados para os usuários da ferramenta. Bibliotecas de áreas como Ciência de Dados, Automação, Aprendizado de Máquina, Aprendizado profundo, Redes Neurais, são bastante famosas e utilizadas pelos usuários da linguagem, podendo ser facilmente importadas para que não haja a necessidade de implementar funções complexas e componentes desde o início.

A alta disponibilidade de bibliotecas e recursos para as áreas de Aprendizado de Máquina e Aprendizado profundo, assim como a facilidade de uso da linguagem foram os principais critérios para escolha do uso dessa linguagem para o desenvolvimento deste trabalho. Abaixo estão listadas as principais bibliotecas Python criadas por terceiros e utilizadas para o desenvolvimento da solução proposta.

4.3.2 Google Colab

O Google Colab³ é uma plataforma desenvolvida pela Google para gerenciamento e execução de códigos online. Ele permite escrever e executar Python no navegador de maneira bastante facilitada, permitindo fácil compartilhamento dos códigos e disponibilidade de uso, gratuitamente, de processamento de CPU e GPU do próprio servidor da Google.

4.3.3 Bibliotecas utilizadas da linguagem

- Numpy⁴: utilizada para trabalhar com vetores e matrizes de forma eficiente e facilitada, possibilitando operações matemáticas velozes com essas estruturas em grandes dimensões e complexidades. Otimizada com a linguagem de programação C, Numpy traz velocidade e performance bastante agradável para essas operações no Python.
- Pandas⁵: biblioteca Python para análise e manipulação de dados. Permite gerenciar dados em grandes volumes, aplicando operações matemáticas de diferentes graus de complexidade de uma forma bastante flexível, fácil, rápida e poderosa. Diversas operações presentes nesta biblioteca são construídas em cima de componentes Numpy.
- Matplotlib⁶: biblioteca base para muitas bibliotecas de visualização de dados gráfica, compreende a criação de gráficos estáticos, animados e interativos em Python.

³<https://colab.research.google.com/>. Acesso em: 15 Janeiro 2023.

⁴<https://numpy.org/>. Acesso em: 20 Dezembro 2022.

⁵<https://pandas.pydata.org/>. Acesso em: 20 Dezembro 2022.

⁶<https://matplotlib.org/>. Acesso em: 20 Dezembro 2022.

- Seaborn⁷: biblioteca Python para visualização de dados em alto nível. Baseada na biblioteca Matplotlib, Seaborn traz muitas implementações de gráficos prontas como: gráfico de linhas, gráfico de barras, gráfico de pontos de calor, etc.
- Tensor Flow⁸: plataforma de código aberta que facilita a criação de modelos de Aprendizado de Máquina e Aprendizagem profunda. Tensor Flow traz possibilidades para que o usuário do *framework* possa trabalhar com modelos prontos ou cria-los de acordo com as necessidades. Com esta biblioteca é possível facilmente estruturar redes neurais simples e complexas, alterando o número de neurônios, função de ativação, adicionando diferentes tipos de camadas, camadas de pré-processamento, entre outros.
- Scikit-learn⁹: também conhecida como Sk Learn, é uma das principais ferramentas para se trabalhar com Aprendizado de Máquina na linguagem Python. Tal biblioteca traz diversas implementações de modelos de regressão, classificação e agrupamento, assim como a implementação de várias métricas de avaliação bastante conhecidas como: acurácia, precisão, matriz de confusão, erro quadrado médio, curva ROC, etc. Dessa forma, o programador que utiliza o Sk Learn pode se preocupar com o uso dos modelos e seus resultados mais do que com a criação do modelo em si.

Construída a partir da junção de outras bibliotecas como Numpy, Matplotlib e SciPy¹⁰ (biblioteca que implementa métodos estatísticos em Python), a Scikit-learn traz consigo também diversos componentes para otimização/ajuste de hiperparâmetros, possibilitando que modelos eficientes sejam criados e aperfeiçoados para contextos mais específicos.

Sk Learn traz também algumas implementações de componentes para divisão de dados inteligente de forma mais inteligente e elaborada, como algoritmos para validação cruzada de dados, e algoritmos para normalização dos dados, como normalização padrão, normalização min-max, entre outras. Com essa biblioteca também é possível gerar *datasets* e aplicar técnicas para correção de conjuntos já existentes, por meio de técnicas como a *OneHotEncoder*, por exemplo, a qual é facilmente aplicada com a Sk Learn.

⁷<https://seaborn.pydata.org/>. Acesso em: 20 Dezembro 2022.

⁸<https://www.tensorflow.org/>. Acesso em: 21 Dezembro 2022.

⁹<https://scikit-learn.org/stable/index.html>. Acesso em: 21 Dezembro 2022.

¹⁰<https://scipy.org/>. Acesso em: 21 Dezembro 2022.

5 Desenvolvimento

5.1 Classificador para fatores de risco

5.1.1 Pré-processamento: preparação dos dados

Antes de iniciar o processo de criação de um modelo preditivo, tanto para classificação quanto para regressão, é necessário preparar os dados que serão utilizados nesse processo. Isto é necessário para que o melhor desempenho deste modelo a ser criado seja alcançado.

As tecnologias analíticas que podemos utilizar são poderosas, mas impõem determinados requisitos sobre os dados que usam. Com frequência, elas exigem que os dados estejam em uma forma diferente de como são fornecidos naturalmente, e alguma conversão será necessária. ([PROVOST; FAWCETT, 2016](#), p. 29-30).

Dessa forma, algumas alterações nos dados foram feitas com intuito de prepará-los para uso correto na criação do classificador. Ao analisar a distribuição da variável alvo do conjunto de dados, foi encontrado a princípio desbalanceamento dos valores. Com 4861 casos de não ocorrência de AVC e apenas 249 casos para ocorrência da doença, seria sintetizado um classificador enviesado para casos de não AVC, classificando dados que receberia de forma tendenciosa.

Para realizar o balanceamento dessas duas classes (avc e não avc) foi realizado o seguinte procedimento:

1. Separação dos elementos do conjunto original em duas classes: avc e não avc.
2. Redução (subamostragem) do conjunto de dados da classe não avc: de 4861 elementos para 275.
3. Junção do conjunto de dados da classe avc (249 elementos) com o novo conjunto de dados reduzido da classe não avc (276 elementos), totalizando 525 elementos.

Portanto, com o balanceamento das classes presentes na variável alvo, tem-se um novo conjunto de dados de 525 elementos com 52.6% dos dados sendo casos de AVC e o restante para casos de condição normal (não AVC). Esse novo conjunto de dados criado foi utilizado ao decorrer de todo este trabalho para desenvolver o modelo classificador.

Com relação ao formato de cada atributo que compõe o conjunto de dados, alguns não apresentavam formato numérico e, portanto, tiveram que ser corrigidos para que não houvesse erros na criação do classificador. As variáveis contínuas presentes no conjunto de dados já

estavam com o formato adequado numérico e não precisaram de correção. Mais informações a respeito da distribuição dos valores de cada coluna podem ser visualizadas na Tabela 1.

Tabela 1 – Frequência de cada possível valor dos atributos categóricos presentes no conjunto de dados.

Atributo	Distribuição das frequências
Gênero	Masculino: 41.71% Feminino: 58.29%
Hipertensão	0: 82.48% 1: 17.52%
Doença cardíaca	0: 88.76% 1: 11.24%
É / já foi casado	Sim: 80% Não: 20%
Tipo de emprego	Setor privado: 62.86% Autônomo: 22.09% Cargo público: 14.29% Nunca trabalhou: 0.76%
Tipo de residência	Urbano: 50.86% Rural: 49.14%
Tabagismo	Nunca fumou: 37.33% Fuma formalmente: 24.57% Fuma constantemente: 22.48% Situação desconhecida: 15.62%
AVC	0: 52.57% 1: 47.43%

Fonte: Elaborado pelo autor.

Para corrigir o formato textual de variáveis categóricas que apresentam apenas duas opções de valores, basta apenas substituir uma delas pelo dígito “1” e a outra pelo “0”, conforme demonstrado na Figura 11. Já para casos em que as colunas apresentam mais do que dois possíveis valores, é possível realizar a correção das mesmas através da técnica de *One-Hot Encoding*¹. Tal transformador recebe uma coluna com diferentes possíveis valores textuais e

¹<https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.OneHotEncoder.html>. Acesso em: 18 Dezembro 2022.

cria novas colunas binárias para cada valor diferente.

Figura 11 – Exemplo de correção para a coluna com apenas duas opções de valores (variável binária).

Gênero	Masculino
Masculino	1
Masculino	1
Feminino	0
Feminino	1
Feminino	0
Masculino	0
Feminino	0

Fonte: Elaborado pelo autor.

Figura 12 – Exemplo de aplicação da técnica *One-Hot Encoding* para correção de uma coluna com quatro opções de valores diferentes.

Tipo de emprego	Privado	Autônomo	Cargo público	Criança
Privado	1	0	0	0
Privado	1	0	0	0
Autônomo	0	1	0	0
Privado	1	0	0	0
Criança	0	0	0	1
Cargo público	0	0	1	0
Autônomo	0	1	0	0

Fonte: Elaborado pelo autor.

Como mostrado na Figura 12, toda linha obrigatoriamente contém o dígito “1” em pelo menos uma das novas colunas criadas. Isso ocorre, pois, a criação dessas colunas é feita a partir da ocorrência de todos os diferentes possíveis valores, de forma a ter uma nova coluna para cada valor diferente presente. Por essa razão, pode-se deduzir o valor de uma coluna a partir dos valores das demais, já que os valores são exclusivos.

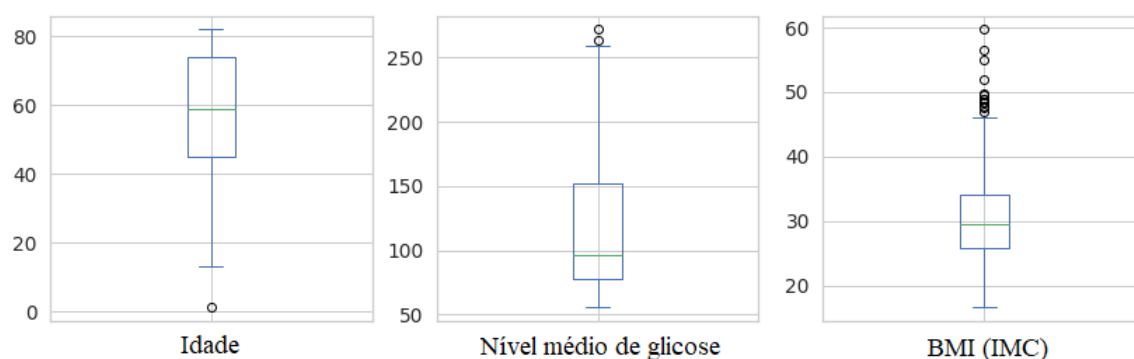
Além do formato numérico das variáveis, também é necessário lidar com a presença de dados nulos no conjunto de dados para evitar erros durante o processo de ajuste do modelo. Dessa forma, cada variável presente no conjunto de dados foi analisada com relação a presença desse tipo de dado e a única que apresentou os dados nulos foi BMI, com 7.6% (40 elementos em um total de 525 registros) de valores desse tipo.

Para correção dos valores nulos encontrados em BMI foi calculada a mediana dos valores não nulos dessa coluna e depois feita a substituição dos valores nulos por essa mediana calculada. A mediana foi escolhida para substituição dos valores nulos pois gerou melhores

resultados em comparação com a substituição desses mesmos valores pela média da coluna. Após tal procedimento, todos os valores nulos do conjunto de dados foram corrigidos, já que só existiam valores assim nessa coluna.

Outro aspecto importante que pode acabar atrapalhando a criação do modelo é a presença de anomalias, ou *outliers*, no conjunto de dados. Tais valores podem atrapalhar a análise de determinado conjunto, gerando ruídos às medidas estatísticas da respectiva distribuição de valores que pertencem. Para identificação desse tipo de valor, foi realizada a análise gráfica das distribuições das variáveis quantitativas do conjunto de dados utilizado, conforme é mostrado na Figura 13.

Figura 13 – Gráfico de distribuição dos atributos contínuos presentes no conjunto de dados com a presença de anomalias, representadas por círculos pretos.



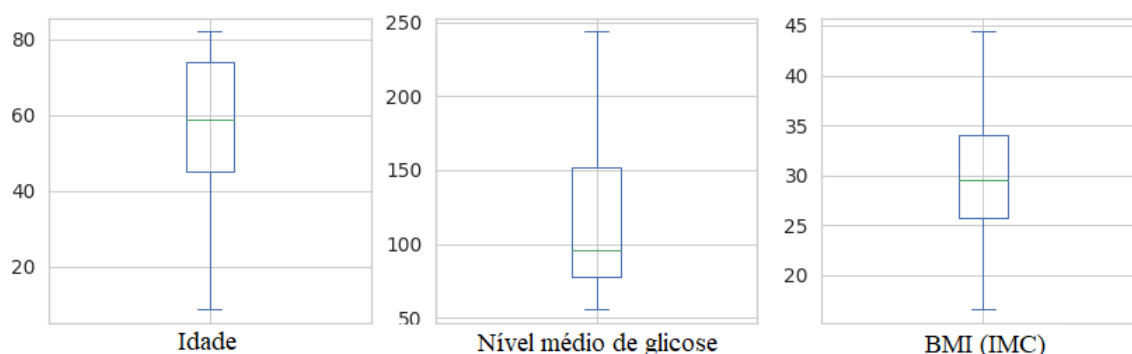
Fonte: Elaborado pelo autor.

Com a identificação das anomalias, agora pode-se iniciar a etapa de tratamento desses valores. Para isso, foi calculado o *Interquartile range* (Amplitude interquartil, em português), ou simplesmente IQR, dos atributos quantitativos. O IQR é utilizado para avaliar o grau de dispersão dos valores em torno da medida de centralidade do conjunto. Com tal valor calculado ($IQR = 3^{\circ} \text{ quartil} - 1^{\circ} \text{ quartil}$), foram determinados o limite superior e o limite inferior dos valores de cada atributo da seguinte forma:

1. Limite superior: $\text{Mediana} + 1.5 \times IQR$
2. Limite inferior: $\text{Mediana} - 1.5 \times IQR$

Com tais medidas calculadas, tornou-se possível substituir todo valor maior que o limite superior pelo próprio limite superior e todo valor menor que o limite inferior pelo próprio limite inferior. Após a realização desses procedimentos, as anomalias foram corrigidas e o resultado obtido pode ser visualizado na Figura 14.

Figura 14 – Gráfico de distribuição dos atributos contínuos presentes no conjunto de dados sem a presença de anomalias, após as alterações realizadas nos valores.



Fonte: Elaborado pelo autor.

Como algumas variáveis numéricas apresentaram ainda intervalos diferentes, a normalização foi realizada para que tal fato não prejudicasse a criação do classificador. Modelos de aprendizado de máquina lineares podem acabar sofrendo com a falta de padronização dos intervalos. Uma variável que varia de 0 a 80, enquanto outra varia de 50 a 200, por exemplo, pode acabar prejudicando o modelo linear. Por essa razão, nesses casos é interessante se realizar a normalização.

A normalização é uma técnica geralmente aplicada como parte da preparação de dados para o aprendizado de máquina. O objetivo da normalização é mudar os valores das colunas numéricas no conjunto de dados para usar uma escala comum, sem distorcer as diferenças nos intervalos de valores nem perder informações. A normalização também é necessária para alguns algoritmos para modelar os dados corretamente. (MICROSOFT, 2022, p. 1).

A normalização evita na modelagem problemas com tais diferenças nos intervalos das variáveis. Ela mantém os valores originais, porém em uma escala padrão. Como a normalização utilizada aqui foi a de máximos e mínimos, através da função `sklearn.preprocessing.MinMaxScaler()`², o padrão de escala ficou entre 0 e 1.

5.1.2 Modelagem: criação do classificador

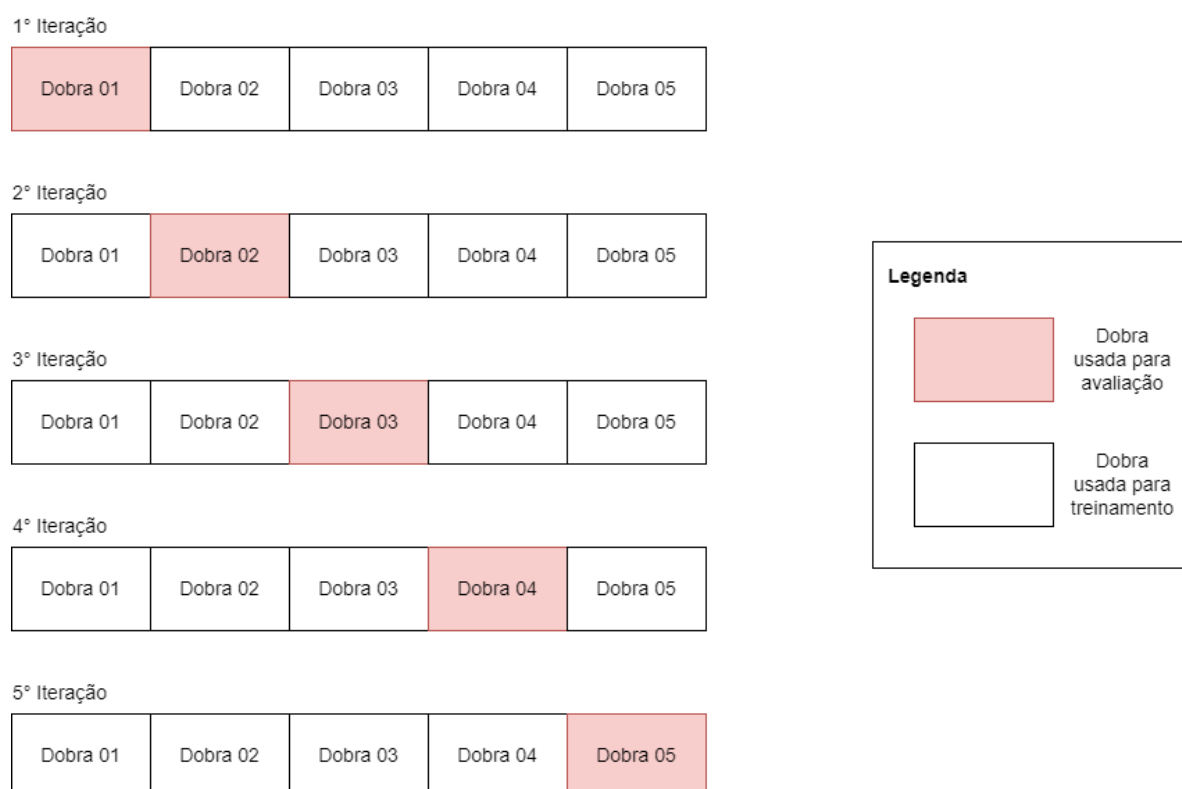
Com a etapa de preparação dos dados feita, o conjunto de dados está pronto para ser processado. Para iniciar o processo de criação do modelo, a modelagem, o conjunto total de dados disponível foi dividido em partes distintas de treino e teste, como forma de evitar o sobreajuste do modelo criado a partir da análise de tais dados. Utilizando uma parte diferente para treino e uma diferente para teste, a avaliação do modelo é feita a partir de dados desconhecidos pelo modelo, mostrando um resultado mais verdadeiro e confiável.

²<https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.MinMaxScaler.html>. Acesso em: 18 Dezembro 2022.

Porém, apenas essa divisão pode não ser suficiente para conjuntos pequenos, pois mesmo criando os conjuntos de treino e teste de forma aleatória, pode ocorrer a seleção de valores mais representativos dependendo da semente definida para aleatoriedade. Por essa razão, uma boa maneira de lidar com conjuntos pequenos é fazer várias divisões de treino e teste no conjunto, pois assim, treina-se o modelo e o avalia com diferentes subconjuntos de dados, trazendo resultados ainda mais confiáveis. Uma forma bastante comum de se fazer isso é através do uso da Validação Cruzada e essa foi a forma utilizada aqui.

A divisão dos elementos foi feita com cinco dobras: os 525 elementos presentes no conjunto de dados total foram divididos em cinco grupos de 125 elementos cada e, então, foram separados quatro deles para treinamento do modelo e um para teste. Esse processo de separação foi repetido por cinco vezes de forma que o modelo fosse novamente treinado e testado com diversas combinações de conjuntos, como ilustrado na Figura 15. Essa validação cruzada foi utilizada para criação de dois modelos classificadores diferentes, um a partir do algoritmo de Regressão Logística e outro a partir do algoritmo de Floresta Aleatória.

Figura 15 – Ilustração do processo de validação cruzada, com cinco dobras, feito para treinamento e avaliação dos dois classificadores criados. A junção das cinco dobras é o conjunto de dados total utilizado.



Fonte: Elaborado pelo autor.

O classificador com Regressão Logística foi criado a partir da classe *LogisticRegression()*³

³https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression. Acesso em: 19 Dezembro

do Scikit-Learn, com o valor 7 definido como semente para geração de valores aleatórios e os demais parâmetros com valores padrões da biblioteca. O modelo para classificação com Floresta Aleatória foi criado a partir da classe *RandomForestClassifier()*⁴, com a semente para geração de valores aleatórios também sendo 7, com duzentas árvores de decisão e com os demais parâmetros com valores padrões da biblioteca.

5.2 Rede Neural classificadora de imagens de tomografia computadorizada

5.2.1 Correção da rotulação dos dados e criação dos conjuntos de treino e de validação

O conjunto de dados original divide as imagens de tomografia computadorizada nas categorias de condição normal, condição de AVC Hemorrágico e condição de AVC Isquêmico, conforme é possível visualizar na Tabela 2. Porém, para abordar de forma mais condizente com o trabalho proposto, uma nova rotulação foi feita. As imagens foram movidas para uma nova distribuição de duas classes: AVC (Hemorrágico e Isquêmico) e não AVC (condição normal), como é mostrado na Tabela 3.

Sobre o tamanho desse novo conjunto gerado a partir da nova rotulação das classes feita, a classe avc possui 299 imagens enquanto que a classe não avc possui 174, sendo que a dimensão de cada imagem do conjunto é de 512 x 512 pixels. A respeito da distribuição dessas imagens para treino e validação, com intuito de evitar sobre ajuste, o conjunto de imagens foi distribuído como 65% (309 imagens) dos dados para treinamento e 35% (166 imagens) para validação. Dessa forma, a rede pode aprender com os dados de treinamento e ser avaliada com os dados de validação, os quais ela ainda não conhece por não terem sido usados na etapa de treino.

2022.

⁴<https://scikit-learn.org/stable/modules/ensemble.html#forests-of-randomized-trees>. Acesso em: 19 Dezembro 2022.

Tabela 2 – Distribuição das imagens do conjunto de dados antes da nova rotulação.

	Condição	Imagens
Rotulação original	Normal	174 (36.8%)
	AVC Hemorrágico	142 (30%)
	AVC Isquêmico	157 (33.2%)

Fonte: Elaborado pelo autor.

Tabela 3 – Distribuição das imagens do conjunto de dados depois da nova rotulação.

	Condição	Imagens
Nova rotulação	AVC	299 (63.2%)
	não AVC	174 (36.8%)

Fonte: Elaborado pelo autor.

5.2.2 Pré-processamento: redimensionamento dos valores e aumento dos dados do conjunto de dados original

Como o conjunto de dados utilizado aqui é composto por imagens, tem-se então um agrupamento de pixels, uma matriz, para cada imagem. Um pixel possui sua coloração definida a partir de um valor, para cada canal de cor, que varia de 0 a 255, totalizando 256 elementos. Embora seja possível realizar o treinamento de uma rede com imagens compostas por pixels de valores entre 0 a 255, redimensionar esse intervalo para valores de 0 a 1 trazem melhores resultados para o modelo, pois os valores que serão multiplicados pelos pesos dos nós acaba sendo menor.

Portanto, foi realizado o redimensionamento dos dados por meio da adição de uma camada inicial na rede do tipo *Rescaling*⁵ com valor para divisão de 255. Dessa forma, toda imagem (matriz) que entra na rede já tem seus pixels redimensionados, com seus valores divididos por 255, que é o maior valor possível presente nos pixels.

⁵https://www.tensorflow.org/api_docs/python/tf/keras/layers/Rescaling. Acesso em: 21 Dezembro 2022.

Outra etapa do pré-processamento realizado foi o aumento de informação do conjunto de dados para que melhores resultados fossem obtidos. Como mencionado na seção anterior, o conjunto de dados disponível conta com 473 imagens no total. Embora seja possível realizar o treinamento de uma rede neural com esse número de registros, aumentar o número de imagens, e consequentemente aumentar a quantidade de informação passada para a rede, pode gerar resultados mais satisfatórios e uma rede neural mais eficiente.

Uma forma bastante comum no campo do Aprendizado Profundo de aumentar os dados de um conjunto, quando não se tem a possibilidade de adicionar novas imagens externas ao conjunto em questão, é utilizar técnicas de aumento dos dados (*data augmentation*⁶), de modo que a quantidade de informação no conjunto cresça a partir das imagens já presentes no conjunto, sem que de fato sejam adicionadas novas imagens no conjunto.

Para realizar o aumento da informação do conjunto de dados foram adicionadas camadas no início da rede neural que fazem algumas alterações no conteúdo da imagem, caracterizando uma nova imagem, porém sem prejudicar as informações usadas para a classificação. Ao todo foram três camadas que aplicam inversão, rotação e ajustes de contraste nas imagens recebidas, respectivamente. Detalhes mais específicos sobre essas camadas citadas podem ser encontrados na Tabela 4.

Tabela 4 – Alterações feitas para aumento dos dados do conjunto de imagens utilizado.

Alteração aplicada	Fator	Método utilizado
Inversão (<i>flip</i>)	Horizontal/ vertical	<code>tensorflow.keras.layers.RandomFlip("horizontal_and_vertical", seed=10)</code>
Rotação	0.2	<code>tensorflow.keras.layers.RandomRotation(0.2, seed=7)</code>
Contraste	Valor entre [0.3, 0.5]	<code>tensorflow.keras.layers.RandomContrast(factor=(0.3, 0.5), seed=(7,17))</code>

Fonte: Elaborado pelo autor.

5.2.3 Estrutura da Rede Neural Convolutacional desenvolvida

Para desenvolver a solução proposta neste trabalho, foi criada uma Rede Neural Convolutacional para classificação binária das imagens de TC nas classes de *avc* e *não avc*. A estrutura dessa rede conta com as camadas que implementam o pré-processamento, comentado

⁶https://www.tensorflow.org/tutorials/images/data_augmentation. Acesso em: 21 Dezembro 2022.

nas seções anteriores, e com a transferência de aprendizado de uma rede VGG-16 pré-treinada, contendo as camadas de convolução e de agrupamento dessa arquitetura, com parâmetros já definidos e pesos já gerados a partir do conjunto de dados *imagenet*⁷, que contém diversas imagens pré-rotuladas para treinamentos avançados de Aprendizado Profundo.

Além da VGG-16, a rede desenvolvida neste trabalho conta também com uma camada de achatamento (*flatten*), com duas camadas densas de neurônios com ativação *ReLU*, com 64 e 32 neurônios, respectivamente, e com uma última camada contendo apenas um neurônio, tendo este como função de ativação a função sigmoide, a qual gera o valor usado para classificação final de determinada imagem processada pela rede. Portanto, a estrutura completa da rede desenvolvida é composta por todos esses componentes citados, conforme mostra a Tabela 5.

⁷<https://www.image-net.org/>. Acesso em: 22 Dezembro 2022.

Tabela 5 – Estrutura final da Rede Neural Convolutacional criada para imagens RGB de dimensão 512x512 pixels.

	Componentes da rede	Principal função do componente	Tamanho da saída
1	Camada de entrada e redimensionamento	Redimensionamento de cada pixel de 0 a 255 para 0 a 1	(512, 512, 3)
2	Camada de inversão	Inversão aleatória (horizontal/vertical)	(512, 512, 3)
3	Camada de rotação	Rotação de 0.2 com sentido aleatório	(512, 512, 3)
4	Camada de contraste	Contraste aleatório com fator de 0.3 a 0.5	(512, 512, 3)
5	Camadas de convolução e agrupamento VGG-16	Transferência de aprendizado de uma CNN VGG, com pesos fixos pré-gerados a partir do conjunto <i>imagenet</i>	(16, 16, 512)
6	Camada de achatamento	Redimensiona saída para apenas um valor	(131072)
7	Camada densa de neurônios	64 neurônios com função de ativação ReLU	(64)
8	Camada densa de neurônios	32 neurônios com função de ativação ReLU	(32)
9	Camada densa de neurônios	Apenas 1 neurônio com função de ativação sigmoide	(1)

Fonte: Elaborado pelo autor.

5.2.4 Treinamento da rede estruturada

Para realizar o treinamento da rede criada, fora feita antes a compilação do modelo, definindo a forma de cálculo do erro como entropia cruzada binária e uma taxa de aprendizado de 0.001. Ademais, fora definida uma função *call-back* para interromper o treinamento caso o modelo alcançasse acurácia de validação de 95% e AUC ROC de validação de 0.9.

O treinamento foi feito a partir do conjunto de dados de treino, composto por dados desconhecidos pelo modelo, com 12 épocas (iterações que percorrem todo o conjunto de dados) e com tamanho do lote (*batch size*) de 32 imagens por iteração do treinamento de cada época. As camadas da rede pré-gerada VGG-16, integrada na rede via transferência de aprendizado, não tiveram seus pesos alterados durante o treinamento feito, pois manter os pesos originais resultou em melhores resultados em comparação com treiná-los novamente a partir do conjunto de dados utilizado neste trabalho.

6 Resultados

O processo de avaliação é uma das principais partes de qualquer projeto que envolva a aplicação de modelos preditivos, pois é preciso analisar desempenhos obtidos e, assim, medir a eficiência de cada um. É de grande importância que este processo seja feito de forma correta para que os resultados da avaliação não sejam enganosos e não escondam falhas do modelo.

6.1 Resultados dos modelos criados para classificação de fatores de risco

Como forma de medir o desempenho geral dos classificadores criados, a acurácia de cada um foi calculada a partir da função `accuracy_score()`¹, da biblioteca Scikit-Learn, para cada dobra de teste utilizada na avaliação de ambos os modelos. Como consequência, tem-se, portanto, uma sequência de cinco valores e a acurácia final foi definida como a média desses valores, como mostra a Tabela 6.

Tabela 6 – Valor médio e desvio padrão das acurácias calculadas em cada iteração de teste da validação cruzada, para ambos os classificadores criados.

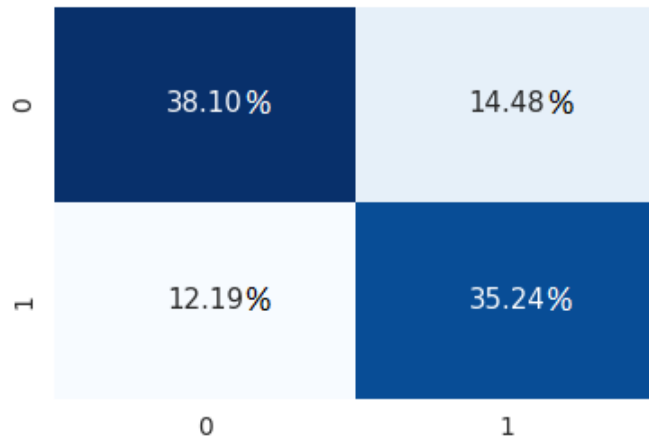
Acurácia do classificador	Regressão Logística	Média: 73.33% Desvio padrão: 2.48%
	Floresta Aleatória	Média: 73.14% Desvio padrão: 1.40%

Fonte: Elaborado pelo autor.

Também foi gerada e armazenada, para cada dobra de teste da Validação Cruzada feita, uma matriz de confusão. Após a execução de todas as iterações da validação, a média e o desvio dos valores foram calculados e os resultados, para o modelo de Regressão Logística e para o modelo de Floresta Aleatória, podem ser vistos nas Figuras 16 e 17, respectivamente.

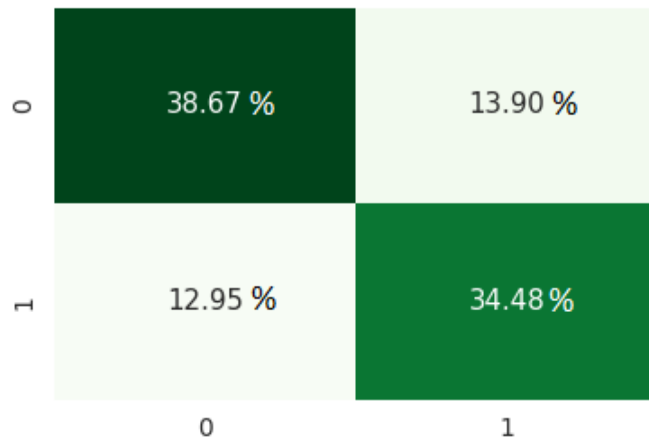
¹https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html. Acesso em: 22 Dezembro 2022.

Figura 16 – Matriz de confusão para o modelo de Regressão Logística, criado a partir da média dos valores obtidos nos testes da validação cruzada.



Fonte: Elaborado pelo autor.

Figura 17 – Matriz de confusão para o modelo de Floresta Aleatória, criado a partir da média dos valores obtidos nos testes da validação cruzada.



Fonte: Elaborado pelo autor.

Além da visualização da frequência de cada tipo de erro e acerto feita pela geração das matrizes de confusão de cada modelo, foi feita também a análise das taxas de verdadeiro positivo e falso positivo para os dois modelos, conforme é mostrado, respectivamente, nas Tabelas 7 e 8.

Tabela 7 – Média e desvio padrão da taxa de valores verdadeiros positivos gerados em cada iteração de teste da validação cruzada.

TPR	Regressão Logística	Média: 74.53% Desvio padrão: 4.10%
	Floresta Aleatória	Média: 72.95% Desvio padrão: 4.11%

Fonte: Elaborado pelo autor.

Tabela 8 – Média e desvio padrão da taxa de valores falsos positivos gerados em cada iteração de teste da validação cruzada.

FPR	Regressão Logística	Média: 27.28% Desvio padrão: 5.27%
	Floresta Aleatória	Média: 26.23% Desvio padrão: 3.99%

Fonte: Elaborado pelo autor.

Para saber sobre o desempenho dos modelos com relação a TPR e FPR a partir da variação de limiares de classificação para o modelo, ou seja, se o modelo apresenta boa capacidade de classificação com relação a dados positivos e negativos, foi calculada a área sob a curva ROC de cada modelo, através da função `roc_auc_score()`², após a validação cruzada feita, e o resultado pode ser visto na Tabela 9.

²https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_auc_score.html. Acesso em: 22 Dezembro 2022.

Tabela 9 – Média e desvio padrão dos valores de AUC ROC calculados em cada iteração de teste da validação cruzada para ambos os modelos.

AUC ROC	Regressão Logística	Média: 0.8029 Desvio padrão: 0.0506
	Floresta Aleatória	Média: 0.7982 Desvio padrão: 0.0349

Fonte: Elaborado pelo autor.

6.2 Resultados da rede neural criada para classificação de imagens de tomografia computadorizada

Para medir a eficiência da rede neural desenvolvida, as mesmas métricas da seção anterior foram calculadas novamente para o novo cenário. A acurácia atingida pelo classificador demonstrou boa capacidade de acerto geral do mesmo, com relação a ambos os conjuntos, de treino e de validação, processados, como mostra a Tabela 10:

Tabela 10 – Acurácia alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.

Acurácia da Rede	Dados de treino	93.06%
	Dados de validação	97.39%

Fonte: Elaborado pelo autor.

Com relação aos tipos de erros, o classificador de imagens de TC apresentou bons resultados quanto ao equilíbrio de ambos, mantendo casos de falso positivo e falso negativo com baixas frequências de ocorrência. Duas matrizes de confusão foram calculadas para análise desses valores, uma para dados de treinamento, já conhecidos pelo modelo por conta da fase de ajuste, e uma para dados de validação, com as rotulações não conhecidas pela rede neural. As Figuras 18 e 19 trazem os resultados dessas matrizes calculadas.

Figura 18 – Matriz de confusão com os valores gerados na classificação feita pela rede a partir do conjunto de dados de treino, com valores conhecidos pelo classificador.

0	111	3
1	12	90
	0	1

Fonte: Elaborado pelo autor.

Figura 19 – Matriz de confusão com os valores gerados na classificação feita pela rede a partir do conjunto de dados de validação, com valores desconhecidos pelo classificador.

0	59	1
1	2	53
	0	1

Fonte: Elaborado pelo autor.

As taxas de verdadeiro positivo e falso positivo também foram calculadas, conforme é mostrado nas Tabelas 11 e 12 respectivamente, para os dados de treinamento e para os dados de validação, com intuito de fornecer uma análise mais precisa desses casos. Os valores obtidos indicam bons resultados, com grandes taxas de acertos para casos positivos e baixas taxas de falsos positivos.

Tabela 11 – Taxa de valores verdadeiros positivos alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.

TPR	Dados de treino	88.24%
	Dados de validação	96.36%

Fonte: Elaborado pelo autor.

Tabela 12 – Taxa de valores falsos positivos alcançada na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida

FPR	Dados de treino	2.63%
	Dados de validação	1.67%

Fonte: Elaborado pelo autor.

Para analisar o desempenho da rede neural criada com relação ao cálculo de TPR e FPR a partir da variação de limiares de classificação, foi calculada a área sob a curva ROC para os dados, conhecidos, de treino e também para os dados, não conhecidos, de validação. Os resultados, mostrados na Tabela 13, indicam que a rede apresenta boa quantidade de acertos mesmo com a variação do limiar, indicando boa capacidade de classificação da mesma de acordo com os dados utilizados.

Tabela 13 – AUC ROC calculadas na classificação dos dados utilizados para treinamento e validação da rede neural desenvolvida.

AUC ROC	Dados de treino	0.9802
	Dados de validação	0.9944

Fonte: Elaborado pelo autor.

7 Conclusão

O trabalho de conclusão de curso desenvolvido neste documento teve como principal intuito demonstrar formas de classificar casos de Acidente Vascular Cerebral com Inteligência Artificial, a partir da análise de informações, fatores de risco e imagens de tomografia computadorizada (TC) sobre vítimas da doença. Para isso, em um primeiro momento foi realizado o levantamento bibliográfico do problema, realizando estudos sobre a doença em si para saber como ela se caracteriza no indivíduo, quais são os motivos para sua ocorrência, qual o comportamento das vítimas e os principais fatores de risco englobados na doença.

Além de conhecer sobre o Acidente Vascular Cerebral, também foram realizadas pesquisas bibliográficas sobre as áreas de Aprendizagem de Máquina (*Machine Learning*) e Aprendizagem Profunda (*Deep Learning*), subáreas da Inteligência Artificial, conhecendo sobre técnicas importantes dessas áreas para o desenvolvimento da solução proposta. Ademais, foram pesquisados conjuntos de dados sobre fatores de risco do AVC e imagens de TC para serem utilizados no treinamento dos modelos classificadores desenvolvidos.

Para o caso da classificação de fatores de risco, dois modelos foram desenvolvidos através da aplicação dos modelos de Regressão Logística e Floresta Aleatória, e seus resultados foram bastante próximos em suas avaliações de desempenho. A acurácia calculada foi de 73.33% para o modelo de Regressão Logística com desvio padrão de 2.48% para cada caso de teste avaliado. Já para o modelo de Floresta Aleatória, a acurácia calculada foi de 73.14%, com desvio de 1.4% para cada teste avaliado.

Esses valores indicam que os classificadores acertam uma quantidade significativa dos casos que avaliam, com boa capacidade de generalização e com performances semelhantes. O valor médio da *AUC ROC* de 0.8 calculado para ambos os modelos demonstra que as decisões de classificação feitas pelos modelos são eficientes mesmo com a variação dos limiares de classificação.

Os valores próximos das Taxas de Verdadeiros Positivos, *TPR*, e de Falsos positivos, *FPR*, calculados para os classificadores, sugerem que as frequências de ocorrências de *TP*, *FP*, *TN* e *TP* foram parecidas para ambos os modelos, e as matrizes de confusão geradas na avaliação dos dois classificadores comprovam isso. Além disso, a *TPR* média calculada de 74.5% para o classificador de Regressão Logística e a média de 73% para o modelo de Floresta Aleatória demonstram bons desempenhos dos modelos criados em relação ao acerto de casos positivos, o que é bastante interessante no contexto do problema abordado.

Para o caso de processamento e classificação de imagens de TC, a Rede Neural do tipo Convolucional desenvolvida apresentou bons resultados na classificação dos dados processados. A acurácia de 96, 76% para dados de treinamento, conhecidos, e 99, 13% para dados de teste,

não conhecidos, indicam que a rede criada consegue acertar a grande maioria dos dados que analisa. O valor calculado de AUC ROC de 0.9874 para treinamento e 0.9997 para teste indicam que o alto desempenho da rede se mantém mesmo com a variação do limiar de classificação, indicando que a mesma consegue distinguir bem imagens de TC dos casos de *avc* e dos casos de *não avc*. A análise dessas métricas indicam que a rede possui uma taxa consideravelmente alta de acertos, acertando praticamente todas as entradas que analisa, o que pode ser um indício de uma rede excessivamente ajustada aos dados utilizados em seu desenvolvimento.

Como o problema em questão, o AVC, é ligado diretamente a questões de medicina e saúde, os classificadores criados não são soluções para autodiagnóstico da doença e nem são substitutos de profissionais da saúde. É extremamente recomendado que em caso de qualquer suspeita, o indivíduo procure assistência médica e não negligencie a situação, para que assim problemas e prejuízos possam ser evitados.

Portanto, foram desenvolvidos modelos classificadores de fatores de risco e de imagens de tomografia computadorizada que apresentaram resultados satisfatórios para fins educacionais, abrangendo técnicas da área de Inteligência Artificial e cumprindo com o intuito deste trabalho. Embora os resultados gerados neste trabalho foram de qualidade considerável, vale ressaltar que os classificadores desenvolvidos são apenas demonstrações de formas de diagnóstico da doença abordada, existindo apenas para fins demonstrativos e não de fato como soluções do problema em si.

Com relação a trabalhos futuros, é interessante realizar a criação das soluções a partir de bases de dados mais completas, com mais registros e variáveis disponíveis para aplicação dos modelos classificadores. Além disso, utilizar, em trabalhos futuros, bases de dados com mais informações a respeito de como os dados foram coletados possivelmente trará resultados mais robustos e sólidos para serem utilizados na prática em contextos reais.

Referências

AWS. *O que é uma rede neural?* 2022. Disponível em: <<https://aws.amazon.com/pt/what-is/neural-network>>. Acesso em: 28 novembro 2022.

BAGHERI, R. *ROC Curve, a Complete Introduction*. 2019. Disponível em: <<https://towardsdatascience.com/roc-curve-a-complete-introduction-2f2da2e0434c>>. Acesso em: 10 novembro 2022.

COPELAND, M. *Qual é a Diferença entre Inteligência Artificial, Machine Learning e Deep Learning?* NVIDIA, 2021. Disponível em: <<https://blog.nvidia.com.br/2021/03/10/qual-e-a-diferenca-entre-inteligencia-artificial-machine-learning-e-deep-learning/>>. Acesso em: 24 dezembro 2022.

DEEPCHECKS. *False Positive Rate: What is the false-positive rate in machine learning?* 2022. Disponível em: <<https://deepchecks.com/glossary/false-positive-rate/>>. Acesso em: 23 dezembro 2022.

FERGUSON, M.; AK, R.; LEE, Y.-T.; LAW, K. Automatic localization of casting defects with convolutional neural networks. In: . [S.l.: s.n.], 2017. p. 1726–1735.

GOOGLE. *Classificação: acurácia*. 2022. Disponível em: <<https://developers.google.com/machine-learning/crash-course/classification/accuracy>>. Acesso em: 23 dezembro 2022.

HONDA, H.; FACURE, M.; YAOHAO, P. *Os Três Tipos de Aprendizado de Máquina*. 2017. Disponível em: <<https://lamfo-unb.github.io/2017/07/27/tres-tipos-am/>>. Acesso em: 10 novembro 2022.

IBM. *Convolutional Neural Networks*. 2022a. Disponível em: <<https://www.ibm.com/topics/convolutional-neural-networks>>. Acesso em: 6 dezembro 2022.

IBM. *What is logistic regression?* 2022b. Disponível em: <<https://www.ibm.com/topics/logistic-regression>>. Acesso em: 10 novembro 2022.

IBM. *What is a Decision Tree?* 2022c. Disponível em: <<https://www.ibm.com/topics/decision-trees>>. Acesso em: 20 dezembro 2022.

JORDAN, J. *Common architectures in convolutional neural networks*. 2018. Disponível em: <<https://www.jeremyjordan.me/convnet-architectures/>>. Acesso em: 6 dezembro 2022.

MICROSOFT. *Componente Normalizar Dados*. 2022. Disponível em: <<https://learn.microsoft.com/pt-br/azure/machine-learning/component-reference/normalize-data>>. Acesso em: 11 dezembro 2022.

MONARD, M. C.; BARANAUSKAS, J. A. Conceitos sobre aprendizado de máquina. In: *Sistemas Inteligentes Fundamentos e Aplicações*. 1. ed. Barueri-SP: Manole Ltda, 2003. p. 89–114. ISBN 85-204-168.

NAVLANI, A. *Understanding Random Forests Classifiers in Python Tutorial: Learn about random forests and build your own model in python, for both classification and regression*. 2018. Disponível em: <<https://www.datacamp.com/tutorial/random-forests-classifier-python>>. Acesso em: 17 novembro 2022.

NURFIKRI, F. *An Illustrated Guide to Artificial Neural Networks*. 2020. Disponível em: <<https://towardsdatascience.com/an-illustrated-guide-to-artificial-neural-networks-f149a549ba74>>. Acesso em: 28 novembro 2022.

OLIVEIRA, R. de Magalhães Carneiro de; ANDRADE, L. A. F. de. Acidente vascular cerebral. *Rev Bras Hipertens*, v. 8, n. 10, p. 280–285, 2001. Disponível em: <<http://departamentos.cardiol.br/dha/revista/8-3/acidente.pdf>>. Acesso em: 25 outubro 2022.

O'SHEA, K.; NASH, R. *An Introduction to Convolutional Neural Networks*. arXiv, 2015. Disponível em: <<https://arxiv.org/abs/1511.08458>>. Acesso em: 5 dezembro 2022.

PROVOST, F.; FAWCETT, T. *Data Science para Negócios: O que Você Precisa Saber Sobre Mineração de Dados e Pensamento Analítico de Dados*. 1. ed. Rio de Janeiro: Alta Books, 2016. 383 p.

RUEDA, F. *Como identificar o AVC isquêmico na tomografia computadorizada? [vídeo]*. 2022. Disponível em: <<https://pebmed.com.br/como-identificar-o-avc-isquemico-na-tomografia-computadorizada/>>. Acesso em: 25 outubro 2022.

SAÚDE, B. V. em. *Acidente vascular cerebral (AVC)*. Ministério da Saúde, 2015. Disponível em: <<https://bvsmis.saude.gov.br/avc-acidente-vascular-cerebral/>>. Acesso em: 27 outubro 2022.

SCIKIT-LEARN. *Linear Models: Logistic regression*. entre 2007 e 2022. Disponível em: <https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression>. Acesso em: 12 novembro 2022.

SCIKIT-LEARN. *Random Forest Classifier*. entre 2007 e 2022. Disponível em: <<https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>>. Acesso em: 20 novembro 2022.

TAMAKE, B. L. *Análise de dados para auxiliar no diagnóstico precoce de Acidente Vascular Cerebral - AVC*. 52 f. Monografia (Graduação) — Faculdade de Ciências, Universidade Estadual Paulista Júlio de Mesquita Filho UNESP, Bauru, 2020.

TIWARI, T.; TIWARI, T.; TIWARI, S. How artificial intelligence, machine learning and deep learning are radically different? *International Journal of Advanced Research in Computer Science and Software Engineering*, v. 8, p. 1, 03 2018.

VARELLA, M. *NEURÔNIO*. 2020. Disponível em: <<https://drauziovarella.uol.com.br/corpo-humano/neuronio/>>. Acesso em: 27 novembro 2022.

YACIM, J.; BOSHOFF, D. Impact of artificial neural networks training algorithms on accurate prediction of property values. *Journal of Real Estate Research*, v. 40, p. 375–418, 11 2018.