

**CENTRO FEDERAL DE EDUCAÇÃO TECNOLÓGICA DE MINAS GERAIS
CAMPUS DIVINÓPOLIS**

Vitor Martins Soares

**INTELIGÊNCIA ARTIFICIAL NA SAÚDE
Da Previsão de Doenças à Classificação em Sistemas de Triagem**

Divinópolis-MG

2023

VITOR MARTINS SOARES

INTELIGÊNCIA ARTIFICIAL NA SAÚDE

Da Previsão de Doenças à Classificação em Sistemas de Triagem

Trabalho de Conclusão de Curso apresentado no curso de Graduação em Engenharia de Computação do Centro Federal de Educação Tecnológica de Minas Gerais como requisito parcial para obtenção do título de Bacharel em Engenharia de Computação.

Orientador: Prof. Me. Michel Pires da Silva

Coorientador: Prof. Dr. Alisson Marques da Silva

DIVINÓPOLIS-MG

2023

VITOR MARTINS SOARES

INTELIGÊNCIA ARTIFICIAL NA SAÚDE

Da Previsão de Doenças à Classificação em Sistemas de Triagem

Trabalho de Conclusão de Curso apresentado no curso de Graduação em Engenharia de Computação do Centro Federal de Educação Tecnológica de Minas Gerais como requisito parcial para obtenção do título de Bacharel em Engenharia de Computação.

Aprovado em 12/12/2023.

Prof. Me. Michel Pires da Silva

Centro Federal de Educação Tecnológica de Minas Gerais

Prof. Dr. Alisson Marques da Silva

Centro Federal de Educação Tecnológica de Minas Gerais

Prof. Me. Tiago Alves de Oliveira

Centro Federal de Educação Tecnológica de Minas Gerais

AGRADECIMENTOS

Iniciando com minha família, o pilar de minha vida, expresso minha gratidão imensa. À minha mãe, Sandra Vaz Soares Martins, que me introduziu ao mundo acadêmico e sempre valorizou o estudo como uma ferramenta de transformação. Ao meu pai, Rildo Martins Ferreira, que sempre me apoiou em minha jornada acadêmica, mostrando entusiasmo por cada nova aprendizagem minha. A ambos, por compartilharem comigo e meu irmão Lucas Martins Soares, também colega de turma, os desafios e conquistas dessa jornada.

Agradeço aos meus avós, Ari Soares de Oliveira e Marina Vaz Soares, que mesmo em tempos de recursos limitados para o estudo, souberam valorizar e transmitir a importância do conhecimento. Aos meus primos e amigos mais próximos, sempre presentes, oferecendo apoio e momentos de descontração necessários para manter o equilíbrio.

Aos meus colegas de classe, cuja companhia e apoio mútuo foram essenciais, e em especial a Gabriel de Souza Rosa e Yuri Dimitre Dias de Faria, com quem compartilhei o entusiasmo pelas competições acadêmicas, abrindo novos horizontes. A Jorge Vitor Gonçalves de Souza e Gabriel Mesquita Pereira, por sua amizade e apoio inestimável, especialmente nos momentos fora da faculdade e em tempo de pandemia. A Lucas Martins Soares, que além de irmão, foi um companheiro que mesmo com nossas discussões, sempre trabalhamos com excelência.

Um agradecimento especial aos professores que marcaram minha trajetória acadêmica, não só pelo conhecimento transmitido, mas também pelas lições de vida. Ao professor Tiago Alves, por despertar meu interesse pela programação; ao professor Marlon, que me fez explorar o fascinante mundo do hardware; ao professor Guilherme, por aprofundar meu apreço pelas ciências exatas; ao professor Nestor, cuja carisma tornava cada aula uma experiência única; ao professor Raulivan, por sua paixão contagiante pelo ensino; à professora Thabatta, nossa guia e apoio constante; ao professor José Geraldo Pedrosa, por conectar tecnologia e sociedade; ao professor Rafael Marcelino, por reacender meu interesse pela física; ao professor Eduardo Habib, por sua dedicação e empatia como coordenador do curso; e ao meu coorientador professor Alisson, por abrir meus olhos para o vasto universo da inteligência artificial. Um

agradecimento final e muito especial ao professor Michel, meu orientador, cujo entusiasmo e orientação foram cruciais em cada etapa deste trabalho.

Por último, mas não menos importante, agradeço à minha psicóloga, cuja orientação foi vital para manter minha saúde mental e foco durante esta importante fase da minha vida.

Dedico este trabalho àqueles cuja força e amor foram a base da minha jornada. À minha família, cujo apoio incondicional e crença nos meus sonhos me impulsionaram a seguir em frente, mesmo nos momentos mais desafiadores. Vocês são a minha inspiração e o meu porto seguro. Aos meus amigos, por estarem sempre ao meu lado, compartilhando risadas e conselhos valiosos. Vocês transformaram os desafios em momentos mais leves e a jornada acadêmica em uma experiência enriquecedora. A todos vocês, meu mais profundo agradecimento. Este trabalho é um reflexo do amor, da sabedoria e da força que cada um de vocês me proporcionou.

RESUMO

A monografia investiga a aplicabilidade de algoritmos de Inteligência Artificial para prever doenças e categorizar urgências em sistemas de triagem. Utilizando um conjunto de dados de sintomas e doenças, o estudo analisa técnicas de codificação de dados e algoritmos de Inteligência Artificial, aplicando métodos de validação robustos para assegurar a confiabilidade dos resultados. Os algoritmos escolhidos - Árvore de Decisão, Redes Neurais Recorrentes e Redes Neurais Convolucionais - demonstraram ser adequados para o tratamento dos dados disponíveis. A combinação das codificações *One Hot* e *TF-IDF* com Árvore de Decisão e as Rede Neurais como a recorrente e a convolucional mostraram-se eficazes na diferenciação de sintomas comuns e na classificação das doenças. Os resultados indicam que o modelo proposto é capaz de classificar categorias de urgência e prever diagnósticos de múltiplas doenças de maneira eficiente. O estudo enfrentou desafios como o risco de *underfitting* e *overfitting*, e a falta de um grande volume de dados variados para treinar eficientemente as Redes Neurais Artificiais. A pesquisa contribui para a área da saúde, fornecendo novas percepções sobre o uso de algoritmos de Inteligência Artificial profundo para a previsão de diagnósticos e direcionamento no sistema de triagem. Ressalta-se a importância de abordagens inovadoras para lidar com a escassez de dados reais e massivos, mostrando como conjuntos de dados artificiais podem ser utilizados para evitar a utilização de dados sensíveis para treinamento.

Palavras-chave: Inteligência Artificial, Aprendizado de Máquina, Saúde, Doença, Triagem.

ABSTRACT

The thesis investigates the applicability of Artificial Intelligence algorithms for predicting diseases and categorizing urgencies in triage systems. Using a dataset of symptoms and diseases, the study examines data encoding techniques and Artificial Intelligence algorithms, applying robust validation methods to ensure the reliability of the results. The chosen algorithms - Decision Trees, Recurrent Neural Networks, and Convolutional Neural Networks - have proven to be suitable for processing the available data. The combination of *One Hot* and *TF-IDF* encodings with Decision Trees and Neural Networks, such as recurrent and convolutional, proved effective in differentiating common symptoms and classifying diseases. The results indicate that the proposed model is capable of classifying urgency categories and predicting diagnoses of multiple diseases efficiently. The study faced challenges such as the risk of *underfitting* and *overfitting*, and the lack of a large volume of varied data to efficiently train the Artificial Neural Networks. The research contributes to the health field, providing new insights into the use of deep Artificial Intelligence algorithms for the prediction of diagnoses and guidance in the triage system. It highlights the importance of innovative approaches to deal with the scarcity of real and massive data, showing how artificial datasets can be used to avoid using sensitive data for training.

Keywords: Artificial Intelligence, Machine Learning, Health, Disease, Triage.

LISTA DE ILUSTRAÇÕES

Figura 1 – Exemplo de Árvore de Decisão.	9
Figura 2 – Exemplo de Rede Neural Artificial.	10
Figura 3 – Codificação <i>One-Hot Encoding</i> (OH).	13
Figura 4 – Matriz de Confusão.	14
Figura 5 – Arquitetura Projetada para o Projeto.	19

LISTA DE ABREVIATURAS E SIGLAS

IA	Inteligência Artificial
ML	Aprendizado de Máquina do inglês <i>Machine Learning</i>
MLP	Múltiplas Camadas de Perceptron do inglês <i>Multi Layer Perceptron</i>
KNN	K vizinhos mais próximos do inglês <i>K - Nearest Neighbors</i>
CNN	Rede Neural Convolutacional do inglês <i>Convolutional Neural Network</i>
P	Positivo
N	Negativo
V	Verdadeiro
F	Falso
SVM	Máquina de Vetores de Suporte do inglês <i>Support Vector Machine</i>
PLN	Processamento de Linguagem Natural
RNN	Rede Neural Recorrente do inglês <i>Recurrent Neural Networks</i>
GNN	Redes Grafos Neurais do inglês <i>Graph Neural Networks</i>
GCN	Redes Grafos Convolucionais do inglês <i>Graph Convolutional Network</i>
RNKN	Rede de Conhecimento Neural Recursivo do inglês <i>Recurrent Neural Knowledge Network</i>
P@10	Precisão em 10 do inglês <i>Precision at 10</i>
DCG	Ganho Cumulativo Descontado do inglês <i>Discounted Cumulative Gain</i>
ANN	Rede Neural Artificial do inglês <i>Artificial Neural Network</i>
LGPD	Lei Geral de Proteção de Dados Pessoais
BoW	Bolsa de palavras do inglês <i>Bag-of-Words</i>
TF-IDF	Frequência do Termo-Inverso da Frequência nos Documentos do inglês <i>Term Frequency Inverse Document Frequency</i>
OH	<i>One-Hot Encoding</i>
STM	Sistema de Triagem de Manchester

SUMÁRIO

1	INTRODUÇÃO	1
1.1	Contextualização	1
1.2	Motivação e Relevância	3
1.3	Objetivo	4
1.3.1	Objetivo Geral	5
1.3.2	Objetivos Específicos	5
1.4	Organização do Trabalho	6
2	REFERÊNCIAL TEÓRICO	7
2.1	Triagem	7
2.2	Inteligência Artificial	8
2.2.1	Árvore de Decisão	8
2.2.2	Redes Neurais Artificiais	9
2.3	Lei Geral de Proteção de Dados	11
2.4	Codificação	12
2.5	Métricas	13
3	TRABALHOS RELACIONADOS	15
3.1	Inteligência Artificial na Saúde	15
4	METODOLOGIA	18
4.1	Contextualização	18
4.2	Extração e Raspagem de Dados	18
4.3	Pré-processamento	20
4.4	Algoritmos de Inteligência Artificial	21
4.5	Considerações	22
5	ANÁLISES E RESULTADOS	23
6	CONCLUSÃO	28
	REFERÊNCIAS	29

1 INTRODUÇÃO

Nesta monografia, é abordada a crescente importância da computação no âmbito da saúde, uma inovação crucial para agilizar o atendimento médico. Desde a chegada do paciente à clínica ou hospital até o diagnóstico e tratamento, a tecnologia tem um papel vital. Será destacado como os sistemas de triagem computadorizados em clínicas e hospitais e previsão de doenças são fundamentais, discutindo os fatores motivacionais e a relevância deles para a elaboração deste Trabalho de Conclusão de Curso. Além disso, serão apresentados os objetivos que guiam a pesquisa e a estrutura proposta para os capítulos subsequentes desta monografia.

1.1 Contextualização

Compreender a relação entre comportamento humano e sintomas é vital para a promoção do bem-estar e prevenção de doenças. Estabelecer uma relação precisa entre doenças e seus sintomas é crucial para assegurar diagnósticos corretos. Muitas vezes, sintomas aparentes podem estar associados a condições de saúde distintas ou ser manifestações normais do corpo, como uma dor nas costas ocasional. Portanto, uma avaliação detalhada do comportamento das enfermidades é vital para uma prática clínica adequada e eficiente (Waddell *et al.*, 1984).

Além disso, é essencial avaliar e priorizar pacientes com base na gravidade e urgência de suas condições. Esse processo, conhecido como sistema de triagem que é um método criado pelos militares para ao apoio à guerra, é atribuído ao cirurgião do exército de Napoleão Bonaparte, chamado Jean Dominique Larrey, que tinha como ideia principal, a separação de soldados feridos com atenção médica urgente e priorizar tratamento médico para a recuperação rápida destes (Coutinho; Cecilio; Mota, 2012).

O sistema de triagem é de muita importância em departamentos de emergência, onde realiza uma avaliação inicial, abrangendo sinais vitais, sintomas e as principais queixas dos pacientes. Essa fase é determinante para estabelecer a prioridade de atendimento, focando especialmente em casos graves ou de risco iminente.

Além disso, o sistema desempenha um papel chave na gestão do fluxo de pacientes, especialmente em ambientes com alta demanda e variabilidade, como é o

caso das emergências hospitalares. Garantindo atenção imediata aos casos mais urgentes, a triagem contribui significativamente para a eficiência do atendimento médico, reduzindo tempos de espera e otimizando a utilização dos recursos de saúde.

Nos cenários de emergência, a eficácia da triagem é indispensável. Estes sistemas não só diminuem a mortalidade e complicações, direcionando os pacientes rapidamente para o cuidado apropriado, mas também são ferramentas valiosas na coleta e análise de dados, essenciais para a melhoria contínua da qualidade e eficiência dos serviços de saúde.

Portanto, implementar sistemas de triagem robustos e eficientes é fundamental para assegurar um atendimento médico de alta qualidade, otimizando recursos e melhorando os resultados de saúde dos pacientes.

No campo dos diagnósticos e da relação entre doença e sintomas, a computação tem apresentado avanços significativos, permitindo a automação de processos que facilitam a identificação precoce de diversas doenças. Com o rápido progresso tecnológico, conceitos como Inteligência Artificial (IA) têm ganhado destaque novamente, tornando-se componentes essenciais para inúmeras soluções emergentes na área da saúde como o reposicionamento e a descoberta de medicamentos, a definição de novos materiais, dispositivos e técnicas, entre outros (Nascimento Neto *et al.*, 2020).

A IA, um vasto campo da computação, caracteriza-se por realizar ações que incluem aprendizado, percepção, raciocínio e tomada de decisões. Em outras palavras, suas estratégias buscam reproduzir conceitos da inteligência humana, visando superá-la em velocidade, precisão e capacidade de processamento (Russell; Norvig, 2009). Essa abordagem representa uma evolução da tecnologia, pois as limitações da mente humana tornam-se evidentes à medida que a complexidade dos problemas aumenta.

Os algoritmos de IA podem "aprender" a partir de coleções de dados de entrada, otimizando seu desempenho por meio de um processo iterativo, o que os capacita a fazer previsões ou tomar decisões sem serem explicitamente programados para essas tarefas. Com isso já é possível implementar algoritmos eficientes para aplicações e domínios de aprendizado específicos (Mitchell *et al.*, 2007).

Diante do exposto, este trabalho visa explorar a aplicação do IA no diagnóstico de doenças e classificação de urgência em sistemas de triagem. A proposta foi utilizar as vantagens da agilidade dos computadores para correlacionar os sintomas físicos com a

necessidade de urgência de atendimento e com doenças, para então auxiliar a detecção do nível de importância dos sintomas e das possíveis enfermidades de maneira mais eficaz e rápida do que seria possível apenas com a intervenção humana, pois, a IA mostrou um potencial enorme em diversos contextos diferentes. Espera-se que, ao aproveitar o potencial dessas tecnologias na área da saúde, seja possível melhorar a agilidade no atendimento emergencial, a eficácia dos diagnósticos e contribuir para a prevenção de doenças e melhor classificação de sistemas de triagem.

1.2 Motivação e Relevância

A Inteligência Artificial tem sido amplamente utilizada na área da saúde, facilitando a descoberta de novos medicamentos e auxiliando na identificação de possíveis doenças por meio da análise de dados. Além do mais, com a utilização de modelos de IA, pode-se melhorar a segurança do paciente, qualidade do atendimento e reduzir os custos de saúde (Waring; Lindvall; Umeton, 2020). Atualmente, no Brasil, cerca de 10 milhões de pessoas pesquisam e se informam sobre saúde na internet regularmente (Moretti; Oliveira; Silva, 2012).

Embora a inteligência artificial esteja em crescente desenvolvimento no mundo acadêmico, apenas 15% dos hospitais estão atualmente o utilizando (Waring; Lindvall; Umeton, 2020). Além do fato de que os dados na área da saúde geralmente não estão disponíveis abertamente e disponibilizados, um dos problemas da utilização de Inteligência Artificial é também a falta de transparência destes sistemas, como configurações do modelo e tomadas de decisão, o que leva na falta de confiabilidade no sistema.

Já a correta classificação de risco depende do treinamento e experiência da enfermeira na aplicação sistema de triagem. Os estudos que mostraram baixa sensibilidade para detectar o paciente emergente e muito urgente relacionaram o resultado também com a capacidade do enfermeiro em classificar corretamente, o que interfere na validação do protocolo. Demonstra, também, a importância da auditoria como elemento de aprimoramento e melhoria do acerto da categoria de prioridade da triagem (Coutinho; Cecilio; Mota, 2012). Isso mostra a necessidade da análise e previsão de maneira correta no sistema de triagem em clínicas e hospitais.

Diante deste cenário, a proposta desta monografia é aplicar IA como Rede Neural Artificial do inglês *Artificial Neural Network* (ANN) e Árvores de Decisão para detectar doenças a partir dos sintomas físicos apresentados pelos pacientes e classificar a urgência do atendimento em sistemas de triagem, que é pouco estudada com a utilização de IA.

A relevância desta proposta reside no potencial para contribuir para a saúde em geral, à medida que sistemas de previsão de doença ou de triagem, podem auxiliar especialistas e hospitais para melhor classificação e redirecionamento de seus pacientes. Especificamente, além deste estudo examinar os algoritmos de IA, explorará também métodos eficazes de codificação de palavras para abordar este caso, e analisará dados relacionados a doenças e sintomas sobre os quais os algoritmos serão aplicados.

Essa abordagem pode redirecionar os pacientes de maneira efetiva com o sistema de triagem e identificar de maneira precoce as doenças, possibilitando tratamentos mais rápidos, eficazes e assim melhorando, conseqüentemente, a qualidade de vida dos pacientes. Além disso, proporciona diagnósticos mais precisos e confiáveis, apoiando os médicos na tomada de decisões clínicas.

Esta pesquisa se justifica devido ao crescente interesse e necessidade de soluções automatizadas na área da saúde, que possam apoiar na detecção precoce de doenças, aprimorar a precisão diagnóstica e potencialmente melhorar os desfechos de saúde. Além disso, a aplicação de algoritmos de IA na saúde é um campo de pesquisa ainda emergente, e esta proposta pode contribuir para a compreensão de seus potenciais aplicações e limitações. Em termos práticos, a implementação de sistemas como este trabalho poderia auxiliar os profissionais de saúde em sua prática clínica, ao oferecer uma ferramenta adicional para o diagnóstico de enfermidades baseado em sintomas.

1.3 Objetivo

A utilização de algoritmos de Inteligência Artificial (IA) para a predição de doenças através de sintomas físicos é uma inovação que promete transformar o diagnóstico médico, proporcionando respostas rápidas e precisas. Paralelamente, a aplicação destes algoritmos em sistemas de triagem tem o potencial de melhorar significativamente a classificação e priorização de pacientes em ambientes de saúde, especialmente em departamentos de emergência.

Para ambos os usos - predição de doenças e triagem eficiente - a seleção criteriosa e análise de dados são fundamentais. Assegurar a qualidade dos dados de treinamento é crucial para evitar erros que possam afetar tanto a precisão do diagnóstico quanto a eficácia da triagem. Além disso, é vital considerar a organização dos dados coletados para maximizar a efetividade dos algoritmos, tanto na identificação precisa de doenças quanto na classificação apropriada dos pacientes no sistema de triagem.

1.3.1 Objetivo Geral

O objetivo deste estudo é treinar modelos de inteligência artificial para que possa melhorar a priorização de pacientes e prever enfermidades com base em sintomas físicos, contribuindo assim, para diagnósticos mais rápidos e precisos na prática clínica.

1.3.2 Objetivos Específicos

Para alcançar o objetivo principal, são necessários os seguintes objetivos específicos:

- a) **Buscar e analisar dados de doenças e sintomas:** Contribui para o objetivo geral, fornecendo a base de dados necessária para a aplicação dos algoritmos de IA.
- b) **Revisar a literatura sobre as aplicações de técnicas de IA aplicadas na área da saúde:** É crucial para entender como a IA tem sido utilizada no contexto da saúde e para identificar as melhores práticas e abordagens que podem ser adaptadas e aplicadas ao conjunto de dados coletado.
- c) **Revisar métodos de codificação de palavras para melhor aplicação das técnicas de IA:** É importante para garantir que as técnicas de IA selecionadas sejam capazes de interpretar e processar adequadamente os dados coletados. A codificação de palavras é um aspecto crucial do processamento de linguagem natural Processamento de Linguagem Natural (PLN), uma subárea da IA que será amplamente utilizada neste trabalho.
- d) **Implementar e testar algoritmos de IA para diagnóstico de doenças e classificação em sistemas de triagem:** A implementação efetiva das

abordagens de IA selecionadas permitirá melhor classificar a priorização e prever as enfermidades a partir dos sintomas físicos apresentados.

- e) **Avaliar diversas destas técnicas e propor a utilização da melhor:** A identificação da melhor técnica para recomendação para uso futuro. A avaliação correta dos métodos contribui para a obtenção de resultados mais precisos e eficazes, apoiando assim o objetivo geral do trabalho.

1.4 Organização do Trabalho

No Capítulo 2, é introduzido um referencial teórico, discutindo em detalhes o funcionamento dos sistemas de triagem, os principais métodos de codificação, algoritmos, técnicas de inteligência artificial e métricas que serão utilizados e abordados no trabalho atual. Esta discussão proporciona o fundamento teórico necessário para o entendimento da abordagem proposta neste trabalho.

Para explorar de maneira abrangente os conceitos discutidos, correlacionando-os ao trabalho proposto, o Capítulo 3 apresenta uma revisão de literatura de trabalhos que aplicam IA na previsão de doenças na área da saúde. Estes trabalhos correlatos são cruciais para entender o atual estado da arte, bem como para identificar oportunidades de pesquisa e contribuições inovadoras.

No Capítulo 4, é apresentado a metodologia planejada e utilizada nesta pesquisa. Isso inclui a descrição do processo de obtenção de dados, o pré-processamento realizado, informações relevantes sobre os dados adotados e os modelos de IA que podem ser o objetivo proposto.

No Capítulo 5, é demonstrado os resultados obtidos por meio de diferentes métodos, seguidos de uma discussão sobre esses resultados e suas implicações. É realizado comparações com os trabalhos revisados na literatura e avaliado a eficiência das previsões realizadas.

Por fim, no Capítulo 6, é descrito as conclusões, reflexões finais e sugestões para trabalhos futuros, com base nos resultados obtidos e nas experiências adquiridas durante a realização desta pesquisa.

2 REFERÊNCIAL TEÓRICO

Neste capítulo, é apresentado os principais conceitos envolvidos na elaboração deste trabalho, desde explicação e funcionamento do sistema de triagem e o funcionamento da proteção de dados reais de enfermos à métodos de codificação, algoritmos e técnicas de aprendizado de máquina, além de métricas empregadas em trabalhos similares. O objetivo é consolidar um alicerce teórico sólido que embasará a metodologia aplicada na avaliação da eficácia do trabalho proposto. Serão discutidos em detalhes informações sobre triagem e sobre a Lei Geral de Proteção de Dados Pessoais (LGPD) e as peculiaridades, vantagens, desvantagens e circunstâncias de aplicação de cada técnica de pré-processamento, previsão e métricas. Este referencial teórico proporciona a base necessária para a compreensão, implementação e avaliação da abordagem proposta neste trabalho.

2.1 Triagem

A triagem tem como objetivo priorizar aqueles que necessitam de atenção mais imediata. Em hospitais e clínicas, a triagem é usada para determinar a gravidade das condições dos pacientes. Isso ajuda a decidir quem precisa de cuidados urgentes e quem pode esperar.

Desenvolvido em 1994, o Sistema de Triagem de Manchester (STM) representa um marco no campo da triagem médica, incorporando 52 fluxogramas diferenciados e uma escala de risco detalhada. Este sistema classifica os pacientes em cinco categorias de urgência, identificadas por cores: vermelho para casos emergenciais, laranja para muito urgente, amarelo para urgente, verde para pouco urgente e azul para situações não urgentes. O STM é notável por sua abordagem objetiva e metodológica na determinação da gravidade dos casos, estabelecendo prioridades clínicas e definindo o tempo ideal de espera para o atendimento, desde a admissão do paciente até a consulta médica.

É essencial ressaltar que o foco do STM é na identificação precisa da principal queixa do paciente, e não no diagnóstico médico em si. No entanto, a adoção do STM em unidades de emergência no Brasil gerou debates. Surgiram preocupações relativas à implementação deste protocolo sem uma validação prévia completa. Isso levanta questões

sobre a eficácia e adequação de um instrumento não testado e potencialmente desalinhado com as nuances culturais e práticas brasileiras (Guedes; Martins; Chianca, 2015).

2.2 Inteligência Artificial

A inteligência artificial desempenha um papel crucial na saúde, como destacado pelos trabalhos discutidos na Seção 3.1. Diversos algoritmos foram empregados nesses estudos para abordar a classificação e/ou previsão de doenças múltiplas classes, demonstrando a amplitude e a profundidade do campo.

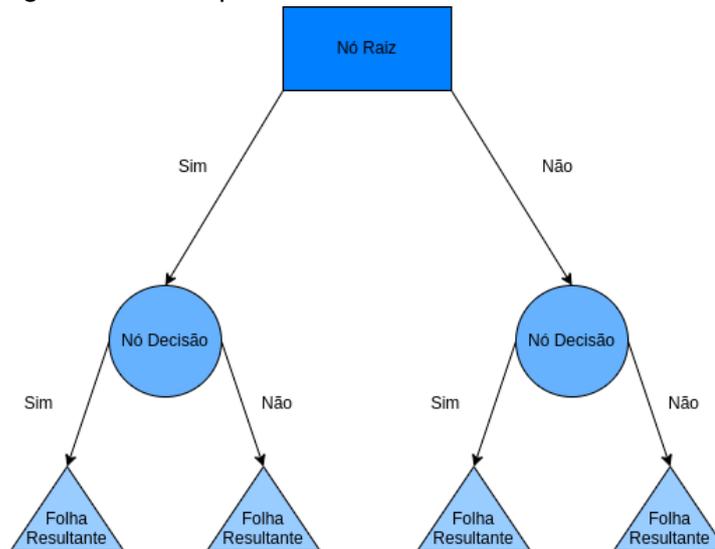
Nesta Seção, faremos um levantamento dos principais algoritmos de aprendizado de máquina utilizados nos trabalhos correlatos. O objetivo é fornecer uma compreensão mais aprofundada dessas técnicas e destacar suas particularidades, vantagens e limitações, além de discutir como elas têm sido aplicadas para resolver problemas semelhantes ao que este trabalho se propõe a enfrentar.

2.2.1 Árvore de Decisão

A Árvore de Decisão (Decision Tree, em são uma classe de redes neurais profundas, comumente aplicadas para analisar imagens. No entanto, estudos recentes têm demonstrado que as CNNs também podem ser eficazes na previsão de doenças(i)nglês) é uma técnica do aprendizado de máquina supervisionado que tem como objetivo criar um modelo de treinamento que é baseado no caminho de árvore, começando da raiz em que cada nó existe uma tomada de decisão, arbitrariamente booleana e com uma sequência de separação de dados até se chegar em uma folha da árvore que é o resultado final do algoritmo como na Figura 1.

Este método utiliza uma estrutura de árvore baseada em decisões para categorizar um conjunto de dados ou para calcular valores específicos vinculados a essas decisões. As principais vantagens incluem a facilidade de implementação, a versatilidade para classificar dados tanto em categorias quanto em valores numéricos. Contudo, apresenta desafios como erros nas decisões, o aumento da complexidade com a expansão da árvore de decisão, e a escalada nos cálculos quando novos dados de treinamento adicionam mais "nós"ou "decisões" (Charbuty; Abdulazeez, 2021).

Figura 1 – Exemplo de Árvore de Decisão.



Fonte: Elaborado pelo autor, 2023

2.2.2 Redes Neurais Artificiais

As Rede Neurais Artificiais são algoritmos capazes de lidar tanto com problemas de classificação quanto de regressão. Estes visam encontrar uma linha que separe os dados, podendo manter no plano de seus dados ou criar novos planos para melhor separá-los. São diversos tipos de ANN, como a RNN, CNN, GNN, GCN e outras, cada uma com sua peculiaridade e aplicação, mas não é do escopo deste artigo discutir isto.

As ANN operam em épocas durante o processo de treinamento. O conceito de épocas refere-se a uma única passagem completa de todos os dados de treinamento pelo modelo neural. Em outras palavras, uma época é concluída quando cada exemplo de treinamento foi apresentado à rede uma vez.

Durante cada época, os pesos e os vies da rede neural são ajustados para minimizar a função de perda, que representa o quão distantes as previsões do modelo estão dos valores reais. Esse ajuste é realizado através de algoritmos de otimização, como o gradiente descendente. O objetivo é aprimorar a capacidade da rede neural de fazer previsões precisas.

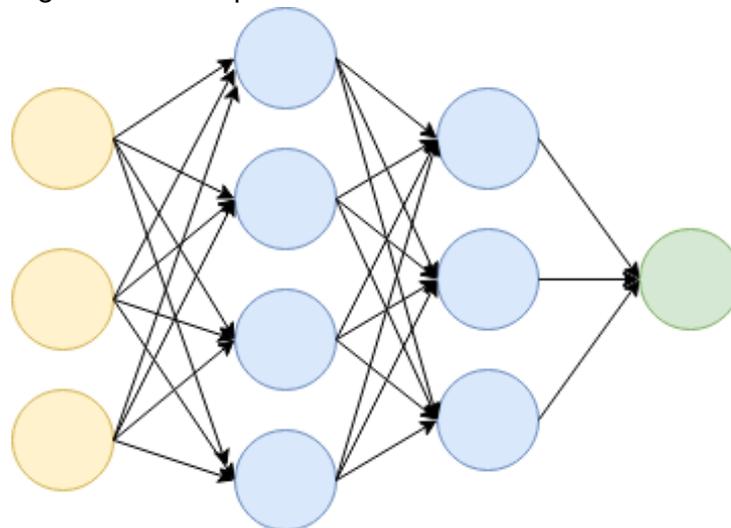
As vantagens das ANN são diversas, como a capacidade de modelar relacionamentos complexos e não lineares de dados, a adaptabilidade para grande variedade de dados, altas quantidades de entradas de dados, o que acrescenta bastante em dimensões a ser processada, que podem ser úteis para manipular dados como

imagens e texto.

Entre as desvantagens das Redes Neurais Artificiais (ANN), destacam-se a exigência de uma ampla variedade e grande volume de dados para o treinamento eficaz. Além disso, sua natureza de “caixa-preta” resulta em falta de transparência, tornando as operações internas destes modelos complexas e desafiadoras de interpretar. Essa característica pode dificultar a compreensão dos processos de tomada de decisão e dos padrões aprendidos pela rede.

Ademais, eles podem enfrentar desafios relacionados ao *'underfitting'* e *'overfitting'*. O *'underfitting'* ocorre quando uma Rede Neural Artificial (ANN) é treinada com dados insuficientes, resultando em um modelo que não aprende adequadamente os padrões dos dados. Por outro lado, o *'overfitting'* é uma situação em que a ANN é treinada excessivamente com os dados, fazendo com que o modelo se ajuste demasiadamente a esses dados e perca sua capacidade de generalização para novos dados não vistos antes (Ray, 2019).

Figura 2 – Exemplo de Rede Neural Artificial.



Fonte: Elaborado pelo autor, 2023

As derivações de redes neurais mais importante que serão abordados nesse trabalho são:

- a) **Redes Neurais Recorrentes (RNN):** As RNNs são uma classe de redes neurais artificiais onde as conexões entre os nós formam um grafo direcionado ao longo de uma sequência temporal. Este fato permite que as RNNs usem sua

memória interna para processar sequências de entradas. Estas redes são comumente usadas no processamento de linguagem natural, e espera-se que sua capacidade de lidar com sequências de dados possa ser útil para o nosso problema de previsão de doenças.

- b) **Redes Neurais Convolucionais (CNN):** As CNNs são uma classe de redes neurais profundas, comumente aplicadas para analisar imagens. No entanto, estudos recentes têm demonstrado que as CNNs também podem ser eficazes na previsão de doenças (Dahiwade; Patle; Meshram, 2019).

2.3 Lei Geral de Proteção de Dados

O trabalho apesar de muito importante na prática ele se esbarra em alguns problemas burocráticos que impossibilita de gerar dados e treinamentos com informações reais e verdadeiras, isto acontece por causa da LGPD. Esta lei visa medidas preventivas, proativas na manutenção e privacidade dos dados de terceiro (Rapôso *et al.*, 2019).

A Lei Geral de Proteção de Dados (LGPD) visa proteger as informações de indivíduos identificáveis, tratadas por entidades durante diversas operações como coleta, armazenamento e processamento. Seu foco principal é assegurar a transparência no manejo desses dados, proteger os direitos dos usuários e prevenir o uso inadequado ou discriminatório das informações pessoais. Um aspecto inovador da LGPD é a introdução do conceito de dados pessoais sensíveis, reforçando seu caráter conceitual e ampliando a proteção no âmbito da privacidade e dos direitos humanos (Leme; Blank, 2020).

Especificamente falando da LGPD na área da saúde, existe diversos problemas até mesmo quando se trata do Sistema Único de Saúde (SUS), que contém uma infinidade de dados ainda não anonimizados, muito menos utilizando pseudônimos que transitam via sistemas digitais e meio físico entre instituições de todas as instâncias federativas. Isso se torna problemático principalmente em municípios ou grupos populacionais pequenos, em que, pelo número restrito de titulares, os riscos de identificação dos indivíduos são mais importantes (Aragão; Schiocchet *et al.*, 2020).

Com isso existe diversos problemas em manter e realizar estudos com dados por causa LGPD obrigando a utilização de dados gerados artificialmente, gerando também dados enviesado e padronizado, o que pode repercutir negativamente nos resultados

gerados pelas Inteligências Artificiais.

2.4 Codificação

A codificação é necessária porque algoritmos são baseados em modelos matemáticos como os de Inteligência Artificial operam com valores numéricos, que significa a não compreensão de letras. Portanto, é necessário transformar os sintomas, que são originalmente expressos textualmente, em representações numéricas. Diversas técnicas podem ser empregadas para essa conversão, incluindo:

- a) **OH:** A codificação *One Hot (OH)* gera um vetor de palavras onde a presença de um valor particular nos dados é marcada como 1 e a sua ausência como 0 como na Figura 3. Este método é eficiente e de fácil implementação, geralmente utilizado para conjuntos de dados com categorias finitas. Entretanto, quando se lida com vetores muito grandes, a seleção de variáveis se torna necessária. A codificação OH é comumente utilizada em ANN devido à sua eficácia e simplicidade (Bagui *et al.*, 2021)
- b) **Bolsa de palavras do inglês *Bag-of-Words (BoW)*:** A codificação *Bag of Words* é uma abordagem que possui semelhanças com o OH, mas difere em um aspecto fundamental: em vez de simplesmente marcar a presença ou ausência de palavras com 0 e 1, o BoW conta a frequência das palavras presentes no texto de entrada que correspondem aos elementos do vetor de palavras, o que possibilita a existência de valores diferentes de 0 e 1. (Goldberg; Hirst, 2017).
- c) **Frequência do Termo-Inverso da Frequência nos Documentos do inglês *Term Frequency Inverse Document Frequency (TF-IDF)*:** O TF-IDF é uma medida estatística que determina a importância de uma palavra em um documento, dentro de um conjunto de documentos. '*Term Frequency*' refere-se à frequência com que um termo ocorre em um documento específico, enquanto '*Inverse Document Frequency*' avalia a importância de um termo no contexto geral, considerando todos os documentos. Assim, um termo que aparece frequentemente em um documento específico é considerado relevante, mas se também for comum em muitos outros documentos, sua relevância é diminuída.

Isso se baseia na ideia de que termos que ocorrem frequentemente em um único documento, mas raramente em outros, são mais informativos.

Figura 3 – Codificação OH.

	Sintoma 1	Sintoma 2	Sintoma 3	Doença
Amostra 1	0	1	1	1
Amostra 2	1	1	0	0

Fonte: Elaborado pelo autor, 2023

2.5 Métricas

As métricas são maneiras de avaliar e analisar os resultados de algoritmos de aprendizado para que a interpretação seja efetuada corretamente. Algumas métricas como precisão, sensibilidade e especificidade são muito sensíveis a dados desbalanceados e pode não ser a melhor maneira de se estimar. Com o fornecimento de avaliações gráficas se tem diversas interpretações do desempenho da classificação. Para calcular estas métricas é necessário entender sobre a matriz de confusão. Para isso, vamos supor que existam duas classes para prever, em que uma amostra pode ser prevista pelo algoritmo, que como Positivo (P) ou Negativo (N), as previsões corretas são aquelas chamadas de Verdadeiro (V) e as incorretas são chamadas de Falso (F). Quando temos um positivo previsto como positivo então é VP, e quando negativo como negativo então é VN. Agora caso o negativo seja previsto como positivo então é FP e caso seja um positivo previsto como negativo então é FN como na figura a seguir:

Com este conhecimento é possível calcular as seguintes métricas:

- a) **Acurácia:** A acurácia diz respeito a previsões V, sendo ela N ou P, em meio a todas as observações.

$$\text{Acurácia} = (\text{VP} + \text{VN}) / (\text{VP} + \text{VN} + \text{FP} + \text{FN}) \quad (2.1)$$

- b) **Precisão:** A precisão diz respeito a previsões positivas quando são realmente

Figura 4 – Matriz de Confusão.

		Classe Verdadeira	
		Positivo (P)	Negativo (N)
Classe Prevista	Verdadeiro (V)	Verdadeiro Positivo (VP)	Falso Positivo (FP)
	Falso (F)	Falso Negativo (FN)	Verdadeiro Negativo (VN)

Fonte: Elaborado pelo autor, 2023

corretas, realizando uma divisão de todas as previsões positivas reais por todas as positivas falsas somadas com verdadeiras.

$$\text{Precisão} = \text{VP} / (\text{VP} + \text{FP}) \quad (2.2)$$

- c) **Recall:** A sensibilidade ou Recall é a proporção de instâncias positivas reais que foram identificadas corretamente, é calculada utilizando a divisão dos verdadeiros positivos por verdadeiros positivos e falsos negativos.

$$\text{Sensibilidade} = (\text{VP}) / (\text{VP} + \text{FN}) \quad (2.3)$$

- d) **F1-Score:** O F1-Score é a média harmônica da precisão e sensibilidade para realizar o balanceamento de ambos.

$$\text{F1-Score} = 2 * \text{Precisão} * \text{Recall} / (\text{Precision} + \text{Recall}) \quad (2.4)$$

Estas métricas são utilizadas principalmente em problemas de classificação binária, mas podem se estender para problemas de classificação de múltiplas classes. Para isso, basta calcular os valores de cada classe e realizar uma média de todas as classes.

3 TRABALHOS RELACIONADOS

Neste capítulo, apresentamos o estado da arte das aplicações de inteligência artificiais na área da saúde. Nosso objetivo é complementar e respaldar as afirmações feitas anteriormente, alinhando-se aos objetivos deste trabalho. Foi revisado os impactos da inteligência artificial na área da saúde e abordagens emergentes. Durante esta discussão, é destacado as vantagens e desvantagens, bem como os resultados significativos dos trabalhos relacionados.

3.1 Inteligência Artificial na Saúde

Os conceitos de inteligência artificial são amplamente abordados nas mais diversas áreas do conhecimento através da previsão ou classificação de dados, tendo como objetivo maximizar o desempenho e minimizar o tempo na tomada de diferentes decisões. Nas múltiplas áreas da saúde é possível observar essa ênfase em metodologias e soluções como nos artigos de (Sun *et al.*, 2020), (Dahiwade; Patle; Meshram, 2019), (Jiang *et al.*, 2020).

Em (Sun *et al.*, 2020) foi proposto um novo modelo como Redes Grafos Neurais para a previsão de doenças, este modelo aprende indutivamente, utilizando-se grafos que contém amostras de conceito médico e de registro de paciente, sendo capaz de lidar com novos pacientes e identificar os sintomas mais relevantes para a previsão da doença. Utilizando os dois grafos, o treinamento é feito baseado em Redes Grafos Convolucionais que consegue prever as doenças para cada novo paciente testado, o adicional dele é a possibilidade de previsão de doenças raras além das comuns. Este modelo oferece uma solução intuitiva e precisa para a previsão de doenças, abordando o problema de escassez de dados e a dificuldade de diagnosticar doenças raras ao mesmo tempo. É utilizado a função de perda como uma métrica de treinamento que ajuda a otimizar o modelo para se ajustar melhor aos dados de treinamento. Como método de avaliação foi utilizado as métricas *Precision*, *Recall*, *F1-Score*.

Em (Dahiwade; Patle; Meshram, 2019) é apresentado duas soluções, utilizando o K vizinhos mais próximos do inglês *K - Nearest Neighbors* (KNN) e a Rede Neural Convolucional (Rede Neural Convolucional do inglês *Convolutional Neural Network* (CNN)),

utilizando um conjunto de dados de sintomas, a partir de hábitos da vida de pessoas e informações de exames de rotina.

Na utilização do algoritmo de CNN o conjunto de dados inicial é primeiro transformado em uma forma vetorial. Neste processo, as palavras, ou no nosso caso, os sintomas, são incorporadas ao vetor. Quando um dado está faltando, o valor é preenchido com 0. O vetor, agora preenchido com a representação das palavras, é então enviado para a CNN. A primeira etapa na CNN é a camada convolucional, seguida pela camada de *pooling*, que utiliza o método de '*pooling* máximo'. Esta camada está conectada a uma rede neural totalmente conectada. Finalmente, a classificação do valor é realizada através da função *softmax*. A métrica utilizada neste trabalho foi a acurácia. A CNN obteve a acurácia melhor para este tipo de estudo comparado á KNN, obtendo também um tempo de execução menor.

Em (Jiang *et al.*, 2020) é apresentada uma solução baseada em uma Rede Neural Recorrente do inglês *Recurrent Neural Networks* (RNN) para o diagnóstico de múltiplas doenças. A RNN proposta utiliza informações de prontuários médicos de pacientes diferentes para criar um modelo generalista capaz de fornecer resultados mais precisos no diagnóstico de doenças em comparação com abordagens clássicas do estado da arte, como redes neurais artificiais de múltiplas camadas e cadeias de Markov. O modelo conhecido como Rede de Conhecimento Neural Recursivo do inglês *Recurrent Neural Knowledge Network* (RNKN) tem um desafio intrínseco: quanto mais parâmetros são aprendidos, maior é a dificuldade de alcançar a solução ótima global. No entanto, o modelo ainda apresenta uma valiosa capacidade de auxiliar profissionais médicos na previsão de diagnósticos para uma gama variada de doenças. Dois pontos essenciais são abordados neste trabalho: o impacto do pré-processamento dos dados no treinamento e a determinação do número adequado de épocas de treinamento. Os autores concluem que o nível de detalhes nos prontuários e o tratamento adequado desses dados influenciam diretamente o número de épocas necessárias para obter um modelo de boa qualidade para o problema. Além disso, eles destacam que a utilização de estratégias auxiliares, como árvore de Huffman, para a composição de assinaturas dos dados pode contribuir positivamente para melhorar a etapa de pré-processamento e, conseqüentemente, para o desempenho final do modelo. Para avaliar e analisar o desempenho do modelo são utilizados métricas como Precisão em 10 do inglês *Precision at 10* (P@10), Ganho

Cumulativo Descontado do inglês *Discounted Cumulative Gain* (DCG) e função de perda. Este método foi comparado com modelo de probabilidade condicional (naive Bayes), modelo de diagnóstico comum baseado em árvore de decisão, floresta randômica, Máquina de Vetores de Suporte do inglês *Support Vector Machine* (SVM), RNN, CNN e entre outros modelos clássicos e bem estabelecidos.

É crucial destacar que os artigos analisados propõem novos métodos para prever doenças, empregando algoritmos inovadores e abordagens diferenciadas no tratamento de dados. Essas novas técnicas aumentam a precisão das previsões e a assertividade dos algoritmos. O presente trabalho, por sua vez, se distingue por focar na escassez de dados reais e massivos, diferentemente dos estudos anteriores que dispõem de tais dados. Nosso diferencial reside na validação das diferentes estratégias de IA como auxílio no processo de triagem, possibilitando amplificar as vantagens e benefícios já encontrados em modelos como o protocolo de Manchester. Esta abordagem pode otimizar a assistência médica, permitindo que, durante a triagem, o paciente seja direcionado à fila de prioridade apropriada, ao mesmo tempo em que fornece ao médico uma previsão preliminar da doença do paciente, baseada no algoritmo. É importante salientar que o objetivo deste estudo não é substituir o profissional especializado, mas sim oferecer suporte na agilização do atendimento e na melhor priorização dos pacientes.

4 METODOLOGIA

Para alcançar os objetivos propostos para este trabalho, é apresentado neste capítulo a metodologia utilizada para a composição de uma solução para a previsão de categorias de urgência ou de doenças com base na análise de sintomas físicos e técnicas de inteligência artificial. Para tanto, tem-se na Seção 4.2 a apresentação do processo de extração do conjunto de dados a ser utilizado como parte da solução abordada. Na Seção 4.3, são detalhados os passos de pré-processamento aplicados aos dados extraídos. Por fim, as técnicas de inteligência artificiais abordadas são apresentadas em detalhes na Seção 4.4.

4.1 Contextualização

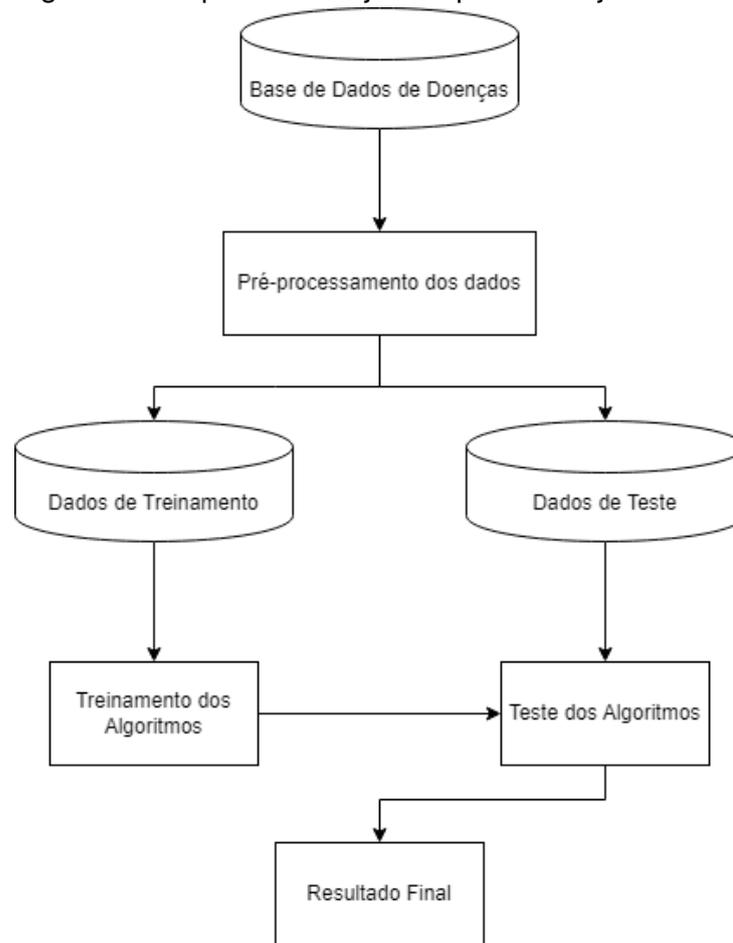
A metodologia segue a seguinte arquitetura como na Figura 5, em que é inicialmente feito a extração do conjunto de dados em, com a base de dados e os dados estabelecidos, após isto, é feito o pré-processamento dos dados, em que os dados serão manipulados e estabelecido os conjuntos de entradas e saídas do sistema, dados de treinamento e teste e como será feito a validação, em seguida é selecionado alguns algoritmos de Inteligência Artificial que serão utilizados para a previsão e classificação, seja das doenças ou do sistema de triagem.

4.2 Extração e Raspagem de Dados

A extração de dados refere-se ao ato de identificar e coletar informações específicas de uma fonte de dados *online*, enquanto a raspagem de dados é o processo de extrair e copiar dados estruturados ou não estruturados de um site ou página web. Para a realização do trabalho tornou-se necessário a coleta dos dados de doenças e seus sintomas de fontes *online*.

Os dados da área da saúde podem ser difíceis de se ter acesso principalmente pela necessidade de se preservar os dados dos pacientes e se respeitar a LGPD. Diversos artigos utilizam dados públicos e específicos para cada doenças, porém, o ideal seria um conjunto de dados com doenças e seus possíveis sintomas.

Figura 5 – Arquitetura Projetada para o Projeto.



Fonte: Elaborado pelo autor, 2023

Para isso, foi selecionada uma base de dados no site *Kaggle*, uma das maiores comunidade de ciência de dados do mundo com ferramentas poderosas para auxiliar no desenvolvimento estes tipos de projetos e estudos (Patil, 2019). Essa base de dados é composta por cerca de 4920 entradas de doenças, contendo 41 doenças diferentes e 131 sintomas diferentes. Cada doença tem cerca de 120 amostras de dados, sendo que cada entrada de doença tem no máximo 17 sintomas, com a possibilidade de obter-se um número menor de sintomas. Os dados dessa base foram gerados artificialmente.

Para a previsão de categorias de urgência do Sistema de Triagem de Manchester (STM) foi utilizado a mesma base de doenças, porém foi substituído as doenças por valores arbitrários de 1 a 5, sendo 1 para situações não urgentes, 2 para pouco urgente, 3 para urgente, 4 para muito urgente e 5 para casos emergenciais. Para simular foi atribuído esses valores de maneira aleatória. Vale salientar a quantidade pequena de classes para se prever, o que diminui bastante a complexidade em relação as doenças.

4.3 Pré-processamento

O pré-processamento de dados é uma etapa crucial na previsão de dados através de inteligência artificial. Esta fase consiste na limpeza e transformação dos dados brutos a fim de torná-los mais adequados para a modelagem. Diversas estratégias de pré-processamento podem ser utilizadas, dependendo da natureza dos dados e do contexto do problema. Técnicas comuns incluem a imputação de dados ausentes, a remoção de *outliers*, a normalização de dados, a codificação de variáveis categóricas, entre outros.

No pré-processamento dos dados efetuado, é categorizado todas as doenças e sintomas identificados, gerando dois conjuntos de atributos distintos. É importante observar que os dados desta base são inicialmente representados como conjuntos de caracteres, uma forma de representação que os algoritmos de aprendizado de máquina não reconhecem diretamente. Assim, foi necessário converter todos os sintomas em representações numéricas ou booleanas para permitir a aplicação eficaz dos algoritmos de Inteligência Artificial.

A conversão foi realizada através da codificação dos sintomas em vetores de palavras. Conforme discutido anteriormente, existem três métodos de codificação potencialmente úteis para essa tarefa: *Bag of Words* (BoW), *One Hot Encoding* (OH), e *Term Frequency-Inverse Document Frequency* (TF-IDF). Contudo, o método BoW não é adequado para este trabalho, pois cada sintoma aparece no máximo uma vez em cada registro. Assim, o vetor gerado por BoW seria idêntico ao gerado por OH.

No entanto, TF-IDF oferece uma abordagem diferente, pois diminui a importância dos sintomas que aparecem frequentemente em várias doenças, como "dor de cabeça". Isso pode ser útil, pois sintomas comuns não necessariamente ajudam a distinguir uma doença de outra. Portanto, optamos por usar a codificação OH e TF-IDF, dependendo das especificidades do modelo aplicado.

Para uma verificação completa dos algoritmos foi estipulado três divisões de conjuntos de dados, cada um com uma parte para treinamento e a outra para testes, sendo elas respectivamente 50% e 50%, 75% e 25%, 90% e 10%. Estas porcentagens foram selecionadas para cobrir diferentes escalas de treinamento, incluindo pequeno, médio e grande. Ao analisar estas três divisões, é possível determinar a quantidade ideal

de variáveis de treinamento necessárias. Para uma validação mais robusta, será empregada a técnica de Validação Cruzada K-Fold, na qual os conjuntos de dados são alternados entre treinamento e teste. Para isso foi utilizado 3 variações de K, sendo K=50, K=25 e K=10, isso para atender as divisões de dados proposta anteriormente.

4.4 Algoritmos de Inteligência Artificial

Os algoritmos de Inteligência Artificial (IA) têm se mostrado extremamente eficazes em uma variedade de tarefas, desde o processamento de linguagem natural até a previsão de doenças. Diferentes algoritmos possuem forças e fraquezas distintas e sua escolha adequadamente é crucial para garantir o melhor desempenho do modelo.

Neste capítulo, é apresentado e três algoritmos de IA que foram selecionados para o treinamento do modelo de previsão de doenças e feito uma discussão sobre eles. Esses algoritmos são:

- a) **Árvore de Decisão:** Este é um tipo de algoritmo de aprendizado supervisionado que é usado para classificação e regressão. Uma árvore de decisão aprende a partir de dados para aproximar uma função de decisão com estruturas de árvore. A opção por utilizar este algoritmo em combinação com a codificação One Hot (OH), uma vez que suas decisões são binárias (sim ou não), assim como os dados codificados pelo método OH (0 ou 1).
- b) **Redes Neurais Recorrentes (RNN):** A configuração utilizada para os testes foi uma camada de *Embedding*, uma camada recorrente de 32 unidades que a função de ativação é tangente hiperbólica, uma camada densa em que a função de ativação é a *sigmoid*.
- c) **Redes Neurais Convolucionais (CNN):** A configuração utilizada para os testes foi uma camada de redimensionamento para ser processado pela camada convolucional, uma camada convolucional de uma dimensão de 32 unidades que a função de ativação é a função relu, uma camada de *pooling* de uma dimensão, uma camada de *Flatten* que converte a saída multidimensional e uma camada densa em que a função de ativação é a *sigmoid*.

Apesar de árvore de decisão não necessitar de uma configuração prévia de número de camadas, de neurônios e o tipo de camada e épocas de treinamento, as ANN precisam,

para isso as redes RNN e CNN. A árvore de decisão é uma ótima combinação com a codificação OH porque suas decisões são resumidas em sim ou não, e a codificação OH resumem os dados em 0 e 1. Já para as ANN a utilização do TF-IDF pode ser benéfico, pois irá diferenciar os sintomas comuns para os não comuns, então para o RNN e CNN será utilizado o TF-IDF.

4.5 Considerações

Será possível ver no próximo capítulo como alguns métodos simples podem ser melhores em determinadas aplicações, fazendo desnecessário a utilização de algoritmos complexos.

É importante estar ciente dos desafios lidados, como o *underfitting* e *overfitting*, que foram monitorados através de métricas de avaliação de modelo apropriadas.

5 ANÁLISES E RESULTADOS

Ao final dos experimentos, foram obtidos uma série de resultados. Inicialmente, destaca-se a análise de dois conjuntos de dados: os dados da Tabela 1, que utiliza OH como método de codificação, e os da Tabela 2, que aplica o TF-IDF. Ambos os conjuntos foram processados com 75% dos dados destinados ao treinamento e os 25% restantes para testes, além da implementação da validação cruzada *K-FOLD* para verificar a capacidade de generalização dos modelos.

Tabela 1 – Resultado utilizando OH.

Triagem					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,99	0,78	1,00	0,75	1,00
<i>Precision</i>	1,00	0,87	1,00	0,81	1,00
<i>Recall</i>	0,99	0,67	1,00	0,63	1,00
<i>F1-Score</i>	0,99	0,75	1,00	0,71	1,00
Doenças					
<i>Accuracy</i>	0,99	0,00	1,00	0,60	1,00
<i>Precision</i>	1,00	0,01	1,00	0,93	1,00
<i>Recall</i>	0,99	0,00	1,00	0,45	1,00
<i>F1-Score</i>	0,99	0,00	1,00	0,61	1,00

Fonte: Elaborado pelo autor, 2023.

A Tabela 1 revela que o OH proporcionou desempenhos notavelmente bons nos algoritmos de Árvore de Decisão e CNN. Por outro lado, a RNN só alcançou um desempenho aceitável com um volume de treinamento mais elevado, especificamente com 100 épocas.

Observando a Tabela 2, percebe-se que o TF-IDF não foi eficaz para o RNN, mesmo após um maior número de épocas de treinamento, enquanto a Árvore de Decisão e a CNN apresentaram desempenhos comparáveis. Isso sugere que, dadas as características e o volume dos dados, a codificação OH pode ser a opção mais adequada.

Na avaliação inicial, surgiu a preocupação de um potencial *Overfitting*, evidenciado por taxas de acurácia atingindo 100%, um fenômeno raro. Isso pode ser atribuído tanto à limitação no volume de dados quanto ao elevado número de sintomas, facilitando o

Tabela 2 – Resultado utilizando TF-IDF.

Triagem					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,99	0,29	1,00	0,75	1,00
<i>Precision</i>	1,00	0,00	1,00	0,81	1,00
<i>Recall</i>	0,99	0,00	1,00	0,63	1,00
<i>F1-Score</i>	0,99	0,00	1,00	0,71	1,00
Doenças					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,99	0,00	1,00	0,00	1,00
<i>Precision</i>	1,00	0,01	1,00	0,02	1,00
<i>Recall</i>	0,99	0,00	1,00	0,00	1,00
<i>F1-Score</i>	0,99	0,00	1,00	0,00	1,00

Fonte: Elaborado pelo autor, 2023.

processo de aprendizagem. Para uma análise mais refinada, introduziu-se um método de seleção de características, o *Variance Threshold*, que elimina atributos com baixa variância, definindo o limiar em 0.10. Essa abordagem de seleção, quando combinada com a codificação TF-IDF, mostrou-se dispensável devido à natureza única desta codificação, que representa o mesmo atributo com valores variados, diferentemente do observado com o OH.

Tabela 3 – Resultado utilizando OH e seleção de características.

Triagem					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,78	0,86	0,86	0,87	0,81
<i>Precision</i>	0,85	0,94	0,90	0,93	0,97
<i>Recall</i>	0,50	0,80	0,81	0,83	0,80
<i>F1-Score</i>	0,63	0,86	0,85	0,88	0,88
Doenças					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,80	0,00	0,81	0,76	0,77
<i>Precision</i>	0,99	0,01	0,93	0,92	0,77
<i>Recall</i>	0,36	0,00	0,77	0,68	0,77
<i>F1-Score</i>	0,53	0,00	0,84	0,78	0,77

Fonte: Elaborado pelo autor, 2023.

Ao observar a Tabela 3, fica evidente a relevância da seleção de atributos em algoritmos como o RNN, que opera com sequências de dados em vez de entradas simultâneas. Em comparação com outros modelos, o RNN exibiu um aprimoramento notável na triagem e na identificação de doenças ao incorporar essa seleção. Contrariamente, os demais algoritmos sofreram uma redução na eficácia, possivelmente devido à exclusão de certos atributos críticos da base de dados. Estes resultados atuais fornecem uma visão mais precisa do impacto de ajustes sutis. Originalmente, 75% da base de dados foi utilizada para treinamento e 25% para testes, mantendo-se a codificação OH e a seleção de características via *Variance Threshold*. Para um entendimento mais aprofundado, serão exploradas proporções diferentes de treinamento e teste, sendo 50% para cada na Tabela 4, e 90% para treinamento contra 10% para teste na Tabela 5.

Tabela 4 – Resultado com 50% de treinamento e teste.

Triagem					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,64	0,83	0,85	0,87	0,80
<i>Precision</i>	0,92	0,90	0,89	0,97	0,98
<i>Recall</i>	0,30	0,80	0,80	0,80	0,79
<i>F1-Score</i>	0,45	0,85	0,84	0,88	0,87
Doenças					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,70	0,00	0,82	0,76	0,77
<i>Precision</i>	0,50	0,02	0,93	0,94	0,77
<i>Recall</i>	0,01	0,00	0,77	0,67	0,77
<i>F1-Score</i>	0,02	0,00	0,84	0,78	0,77

Fonte: Elaborado pelo autor, 2023.

Analisando as Tabelas 3, 4 e 5, é possível ver que a CNN teve um resultado melhor com maior porcentagem de treinamento quando se tinha poucas épocas de treinamento, com muitas não alterou tanto, entretanto para rede neural recorrente o desfecho com porcentagens maiores foi razoavelmente melhor, não o suficiente para influenciar negativamente em uma previsão real, contudo a árvore de decisão obteve resultados semelhantes em todas divisões de treinamento e teste. Para analisar os dados

Tabela 5 – Resultados com 90% de treinamento e 10% de teste.

Triagem					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,82	0,87	0,87	0,88	0,81
<i>Precision</i>	0,87	0,94	0,89	0,95	0,98
<i>Recall</i>	0,61	0,81	0,82	0,82	0,79
<i>F1-Score</i>	0,72	0,87	0,85	0,88	0,87
Doenças					
	CNN(10 Ep)	RNN(10 Ep)	CNN(100 Ep)	RNN(100 Ep)	Árvore de Decisão
<i>Accuracy</i>	0,80	0,00	0,81	0,79	0,77
<i>Precision</i>	0,98	0,01	0,94	0,92	0,77
<i>Recall</i>	0,57	0,00	0,77	0,74	0,77
<i>F1-Score</i>	0,72	0,00	0,85	0,82	0,77

Fonte: Elaborado pelo autor, 2023.

de cada algoritmo, a Tabela 3 será referenciada. Conclui-se que a árvore de decisão é eficaz em problemas mais simples, especialmente sem a seleção de características, onde apresentou desempenho impecável devido à sua habilidade de correlacionar facilmente características e desfechos. Contudo, com a seleção, a performance na previsão de doenças não foi tão robusta, embora tenha se mantido consistente em todas as métricas. Na previsão de triagem, destaca-se com a melhor precisão (98%) e *Recall* e *F1-Score* comparáveis aos outros algoritmos. Quanto à CNN, mostrou-se o método mais eficiente para a previsão de doenças, independente da quantidade de épocas de treinamento. Treinamentos mais longos melhoraram notavelmente o *Recall* e, por extensão, o *F1-Score*, embora com uma ligeira redução na precisão. Esta tendência também foi observada na classificação de triagem, com desempenhos excelentes em ambas as configurações de treinamento. Por sua vez, a RNN exibiu comportamento distinto dos demais modelos, com desempenho insatisfatório na previsão de doenças em treinamentos curtos, mas melhorias significativas com treinamentos mais extensos, superando até mesmo a árvore de decisão em precisão e *F1-Score*. Na triagem, as performances com diferentes durações de treinamento foram satisfatórias, ultrapassando a CNN nas mesmas configurações.

Observa-se que, com dados artificialmente gerados e adequadamente pré-processados, tanto a árvore de decisão quanto a CNN apresentam desempenhos

excelentes. No entanto, a eficácia da árvore de decisão diminui significativamente com a implementação da seleção de características, uma limitação atribuída à sua capacidade restrita de generalização. Em contraste, a RNN demonstra melhorias notáveis com a aplicação desta seleção, especialmente na classificação de triagem, onde superou as outras duas abordagens. Contudo, na previsão de doenças, seu desempenho ainda é insuficiente. A CNN destaca-se por sua consistência, mantendo um nível elevado de eficácia independentemente da codificação (seja OH ou TF-IDF) e das proporções de treinamento e teste utilizadas. Assim, recomenda-se a utilização da CNN em ambas as aplicações, conforme evidenciado pelos aspectos mencionados.

6 CONCLUSÃO

Este estudo explorou a aplicação de algoritmos de Inteligência Artificial e Aprendizado de Máquina para a previsão de doenças e categorização de urgências em sistemas de triagem. O foco foi dado à análise de dados de sintomas e doenças, utilizando técnicas de codificação avançadas e validação robusta para assegurar a confiabilidade dos resultados.

Dentre os algoritmos avaliados, a Rede Neural Convolutiva (CNN) destacou-se como a mais eficiente. Sua versatilidade foi comprovada pela sua capacidade de adaptar-se a diferentes técnicas de codificação, tanto OH quanto TF-IDF, mantendo um desempenho superior em todas as configurações de teste. Este resultado reforça a habilidade da CNN em processar e extrair características significativas dos dados, um aspecto crucial para a eficiência em diagnósticos de saúde e triagens.

A CNN mostrou não apenas uma alta capacidade de aprendizado em volumes variados de dados, mas também uma notável resistência ao *overfitting*, o que a torna uma ferramenta confiável e precisa para previsões em cenários reais. Sua aplicabilidade estende-se além da previsão de doenças, sendo útil também na análise de imagens médicas e outras aplicações na saúde.

Este trabalho contribui para o campo do aprendizado de máquina na saúde, demonstrando o potencial da CNN em melhorar processos de triagem e diagnóstico. Apesar dos desafios enfrentados, como o risco de *overfitting* e *underfitting* em outros modelos, e a necessidade de grandes volumes de dados diversificados, os resultados aqui apresentados são promissores.

Para trabalhos futuros, sugere-se a expansão do conjunto de dados com informações mais abrangentes, experimentação com outros modelos de IA e estratégias para aprimorar a eficácia dos modelos. Além disso, propõe-se a aplicação desses modelos como ferramentas de apoio para profissionais de saúde, visando agilizar o atendimento e melhorar a priorização dos pacientes em sistemas de saúde. Este estudo não visa substituir o julgamento clínico, mas oferecer um recurso complementar que pode enriquecer o processo de tomada de decisão médica.

REFERÊNCIAS

ARAGÃO, Suélyn Mattos de; SCHIOCCHET, Taysa *et al.* Lei geral de proteção de dados: desafio do sistema único de saúde, 2020.

BAGUI, Sikha *et al.* Machine learning and deep learning for phishing email classification using one-hot encoding. **Journal of Computer Science**, Science Publications, v. 17, n. 7, p. 610–623, 2021.

CHARBUTY, Bahzad; ABDULAZEEZ, Adnan. Classification based on decision tree algorithm for machine learning. **Journal of Applied Science and Technology Trends**, v. 2, n. 01, p. 20–28, 2021.

COUTINHO, Ana Augusta Pires; CECILIO, Luiz Carlos de Oliveira; MOTA, Joaquim Antônio César. Classificação de risco em serviços de emergência: uma discussão da literatura sobre o Sistema de Triagem de Manchester. **Rev Med Minas Gerais**, v. 22, n. 2, p. 188–98, 2012.

DAHIWADE, Dhiraj; PATLE, Gajanan; MESHARAM, Ektaa. Designing disease prediction model using machine learning approach. *In*: IEEE. 2019 3rd International Conference on Computing Methodologies and Communication (ICCMC). [S. l.: s. n.], 2019. p. 1211–1215.

GOLDBERG, Yoav; HIRST, Graeme. Neural network methods in natural language processing. morgan & claypool publishers (2017). **zitiert auf**, p. 69, 2017.

GUEDES, Helisamara Mota; MARTINS, José Carlos Amado; CHIANCA, Tânia Couto Machado. Valor de predição do Sistema de Triagem de Manchester: avaliação dos desfechos clínicos de pacientes. **Revista Brasileira de Enfermagem**, SciELO Brasil, v. 68, p. 45–51, 2015.

JIANG, Jingchi *et al.* Medical knowledge embedding based on recursive neural network for multi-disease diagnosis. **Artificial Intelligence in Medicine**, Elsevier, v. 103, p. 101772, 2020.

LEME, Renata Salgado; BLANK, Marcelo. Lei Geral de Proteção de Dados e segurança da informação na área da saúde. **Cadernos Ibero-Americanos de Direito Sanitário**, v. 9, n. 3, p. 210–224, 2020.

MITCHELL, Tom Michael *et al.* **Machine learning**. [S. l.]: McGraw-hill New York, 2007. v. 1.

MORETTI, Felipe Azevedo; OLIVEIRA, Vanessa Elias de; SILVA, Edina Mariko Koga da. Acesso a informações de saúde na internet: uma questão de saúde pública? **Revista da Associação Médica Brasileira**, SciELO Brasil, v. 58, p. 650–658, 2012.

NASCIMENTO NETO, Conrado Dias do *et al.* Inteligência artificial e novas tecnologias em saúde: desafios e perspectivas. **Brazilian Journal of Development**, v. 6, n. 2, p. 9431–9445, 2020.

PATIL, Pranay. **Disease Symptom Description Dataset**. Accessed: 25/10/2023. 2019. Disponível em: %5Curl%7Bhttps://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset%7D.

RAPÔSO, Cláudio Filipe Lima *et al.* Lgpd-lei geral de proteção de dados pessoais em tecnologia da informação: Revisão sistemática. **RACE-Revista de Administração do Cesmac**, v. 4, p. 58–67, 2019.

RAY, Susmita. A quick review of machine learning algorithms. *In: IEEE*. 2019 International conference on machine learning, big data, cloud and parallel computing (COMITCon). [S. l.: s. n.], 2019. p. 35–39.

RUSSELL, Stuart; NORVIG, Peter. **Artificial intelligence**. 3. ed. Upper Saddle River, NJ: Pearson, dez. 2009.

SUN, Zhenchao *et al.* Disease prediction via graph neural networks. **IEEE Journal of Biomedical and Health Informatics**, IEEE, v. 25, n. 3, p. 818–826, 2020.

WADDELL, Gordon *et al.* Symptoms and signs: physical disease or illness behaviour? **Br Med J (Clin Res Ed)**, British Medical Journal Publishing Group, v. 289, n. 6447, p. 739–741, 1984.

WARING, Jonathan; LINDVALL, Charlotta; UMETON, Renato. Automated machine learning: Review of the state-of-the-art and opportunities for healthcare. **Artificial intelligence in medicine**, Elsevier, v. 104, p. 101822, 2020.