



Universidade Federal do Rio de Janeiro

Escola Politécnica

MBA em Engenharia de Dados

(MBED)

**REDES NEURAIAS APLICADO EM EXAMES DE
ELETROCARDIOGRAMA**

Autor:

Carlos Alberto Barbosa Filho

Orientador:

Cláudio Luiz Latta de Souza, M.Sc.

Coorientador:

Manuel Villas Boas Junior, M.Sc.

Examinador:

Norberto Bellas, M.Sc.

Examinador:

Vinicius Drumond Gonzaga, M.Sc.

Rio de Janeiro

Agosto de 2023

Declaração de Autoria e de Direitos

Eu, **Carlos Alberto Barbosa Filho** CPF 216933498-09, autor da monografia **REDES NEURAIIS APLICADA EM EXAMES DE ELETROCARDIOGRAMA**, subscrevo para os devidos fins, as seguintes informações:

1. O autor declara que o trabalho apresentado na defesa da monografia do curso de Pós-Graduação, Especialização MBA em Big Data, Business Intelligence e Business Analytics da Escola Politécnica da UFRJ é de sua autoria, sendo original em forma e conteúdo.
2. Excetuam-se do item 1 eventuais transcrições de texto, figuras, tabelas, conceitos e ideias, que identifiquem claramente a fonte original, explicitando as autorizações obtidas dos respectivos proprietários, quando necessárias.
3. O autor permite que a UFRJ, por um prazo indeterminado, efetue em qualquer mídia de divulgação, a publicação do trabalho acadêmico em sua totalidade, ou em parte. Essa autorização não envolve ônus de qualquer natureza à UFRJ, ou aos seus representantes.
4. O autor declara, ainda, ter a capacidade jurídica para a prática do presente ato, assim como ter conhecimento do teor da presente Declaração, estando ciente das sanções e punições legais, no que tange a cópia parcial, ou total, de obra intelectual, o que se configura como violação do direito autoral previsto no Código Penal Brasileiro no art.184 e art.299, bem como na Lei 9.610.
5. O autor é o único responsável pelo conteúdo apresentado nos trabalhos acadêmicos publicados, não cabendo à UFRJ, aos seus representantes, ou ao(s) orientador(es), qualquer responsabilização/ indenização nesse sentido.
6. Por ser verdade, firmo a presente declaração.

Rio de Janeiro, 05 de agosto de 2023.

Carlos Alberto Barbosa Filho

UNIVERSIDADE FEDERAL DO RIO DE JANEIRO

Av. Athos da Silveira, 149 - Centro de Tecnologia, Bloco H, sala - 212,
Cidade Universitária Rio de Janeiro – RJ - CEP 21949-900.

Este exemplar é de propriedade Escola Politécnica da Universidade Federal do Rio de Janeiro, que poderá incluí-lo em base de dados, armazenar em computador, microfilmear ou adotar qualquer forma de arquivamento.

Permitida a menção, reprodução parcial ou integral e a transmissão entre bibliotecas deste trabalho, sem modificação de seu texto, em qualquer meio que esteja ou venha a ser fixado, para pesquisa acadêmica, comentários e citações, desde que sem finalidade comercial e que seja feita a referência bibliográfica completa.

Os conceitos expressos neste trabalho são de responsabilidade do(s) autor(es).

DEDICATÓRIA

Dedico este trabalho a toda comunidade acadêmica da UFRJ.

Sinto muito orgulho de fazer este MBA nesta instituição, tendo contato com professores e alunos brilhantes, que muito contribuem para o desenvolvimento.

AGRADECIMENTO

Agradeço a Deus, e minha família, que sempre se empenharam, para que eu pudesse ter acesso a educação, incentivando e proporcionando condições para tal.

Agradeço a minha companheira Angélica, pelo carinho, apoio e compreensão, em tantos momentos em que preciso me privar de algumas coisas, para me dedicar exclusiva aos estudos.

Agradeço ao professor Claudio Latta, que sempre me deu toda atenção necessária e colaborou muito neste MBA, contribuindo para meu aprendizado.

Por fim, agradeço a todos da comunidade UFRJ que estiveram juntos nesta jornada.

RESUMO

A cada dia, os recursos de inteligência artificial têm conquistado maior importância na área da saúde, tornando-se uma ferramenta indispensável, para análises mais rápidas e assertivas, principalmente em exames complexos. O eletrocardiograma (ECG) embora seja um exame simples de aplicar, a interpretação dos gráficos não é tão simples. Com a utilização de inteligência artificial (IA), a interpretação fica mais fácil, acessível a qualquer profissional de saúde, além de melhorar a assertividade. Este trabalho tem por objetivo demonstrar como a IA pode ser utilizada nesses exames, para classificação de arritmias cardíacas, através da utilização de redes neurais. Para isso foi usado uma base com dados de exames de ECG, para o treinamento de duas redes neurais, sendo uma abordagem tradicional, de classificação através de dados numéricos de ECG, e a outra de classificação através de imagens com rede neural convolucional, onde as imagens utilizadas são referentes aos conjuntos numéricos da mesma base de dados.

Palavras-Chave: (Arritmias, Redes Neurais, CNN, Machine Learning, Inteligência Artificial)

ABSTRACT

Every day, artificial intelligence resources have gained greater importance in the healthcare sector, becoming an indispensable tool for faster and more assertive analyses, especially in complex exams. The electrocardiogram (ECG), although it is a simple test to apply, interpreting the graphs is not so simple. With the use of artificial intelligence (AI), interpretation becomes easier, accessible to any healthcare professional, in addition to improving assertiveness. This work aims to demonstrate how AI can be used in these exams, to classify cardiac arrhythmias, through the use of neural networks. For this, a dataset was used, with data from ECG exams, for the training of two neural networks, one being a traditional approach, classification through numerical ECG data, and the other classification using images with a convolutional neural network, where the images used refer to numerical sets from the same dataset.

Keywords: (Arrhythmias, Neural Networks, CNN, Machine Learning, Artificial Intelligence)

SIGLAS

AVC	Acidente Vascular Cerebral
CNN	Convolucional neural network
DCV	Doença Cardiovascular
ECG	Eletrocardiograma
GBD	Global Burden of Disease
IA	Inteligência Artificial
IBGE	Instituto Brasileiro de Geografia e Estatística
IOT	Internet of Things
KNN	K-Nearest-Neighbor
LABDAPS	Laboratório de big data e análise preditiva em saúde
LGPD	Lei Geral de Proteção de Dados
OMS	Organização Mundial de Saúde
PNS	Pesquisa Nacional de Saúde
RBEM	Revista Brasileira de Educação Médica
SUS	Sistema Único de Saúde
UFRJ	Universidade Federal do Rio de Janeiro
UNICAMP	Universidade de Campinas

LISTA DE FIGURAS

Figura 2.1.2.1 – Posicionamento das derivações precordiais.	7
Figura 2.1.2.2 – Corte transversal do torax, correlações dos eletrodos e a região do coração.	7
Figura 2.1.2.3 – Representação gráfica do batimento cardíaco no ECG.	8
Figura 2.1.2.4 – Batimento cardíaco normal no ECG.	9
Figura 2.1.2.5 – Batimento cardíaco extrassístole ventricular no ECG.	9
Figura 2.1.2.6 – Batimento cardíaco extrassístole atrial no ECG.	10
Figura 2.1.2.7 – Batimento de fusão ventricular no ECG.	10
Figura 2.5.5.1 – Underfitting x Overfitting	17
Figura 2.5.6.1 – Funcionamento básico do neurônio artificial	18
Figura 2.5.6.2 – Funcionamento básico das redes neurais	19
Figura 2.5.7.1 – Arquitetura de rede neural convolucional	20
Figura 3.2.1 – Organização dos dados	24
Figura 3.2.1.1 – Distribuição das amostras de treino	25
Figura 3.2.1.2 – Distribuição das amostras de teste	26
Figura 3.2.1.3 – Amostras de treino (barras)	26
Figura 3.2.1.4 – Amostras de teste	27
Figura 3.2.1.5 – Distribuição do conjunto de treino balanceado	28
Figura 3.3.1 – Plotagem das batidas cardíacas	29
Figura 3.3.2 – Atuação do filtro gaussiano	29
Figura 3.4.1 – Dados de treino e teste consolidados	31
Figura 4.4.1 – Desempenho de treinamento CNN	34
Figura 4.4.2 – Desempenho de treinamento CNN	35
Figura 4.5.1 – Previsões erradas	36

LISTA DE TABELAS

Tabela 3.3.1 – Desempenho treinamento rede neural convolucional	30
Tabela 3.4.1 – Desempenho treinamento rede neural	31
Tabela 4.3.1 – Desempenho conjunto de teste	33
Tabela 4.4.1 – Performance conjunto de teste	33
Tabela 4.5.1 – Previsões erradas	37

LISTA DE QUADROS

Quadro 3.2.1.1 – Distribuição do conjunto de dados	25
Quadro 3.2.1.2 – Dados ausentes	27

Sumário

Capítulo 1: Introdução.....	1
1.1 Tema.....	1
1.2 – Justificativa.....	1
1.3 – Objetivos.....	2
1.3.1 – Objetivo Geral.....	2
1.3.2 – Objetivo Específicos.....	2
1.4 – Delimitação.....	2
1.5 – Metodologia.....	2
1.6 – Descrição.....	3
Capítulo 2: Referencial Teórico.....	4
2.1 – A saúde pública no Brasil.....	4
2.1.1 – Doenças cardiovasculares.....	5
2.1.2 – Eletrocardiograma.....	5
2.2 – Big Data.....	10
2.3 - DATASUS.....	11
2.4 – Ciência de dados.....	11
2.4.1 – Análise exploratória.....	12
2.4.2 – Análise diagnóstica.....	12
2.4.3 – Análise preditiva.....	12
2.4.4 – Análise prescritiva.....	12
2.4.5 – Google Colab.....	13
2.4.6 – Python.....	13
2.4.7 – Pandas.....	13
2.4.8 – Numpy.....	13
2.4.9 – Seaborn.....	14
2.4.10 – Matplotlib.....	14
2.4.11 – Scikit-learn.....	14
2.4.12 – TensorFlow.....	14
2.4.13 – Inteligência artificial.....	14
2.5 – Machine learning.....	15
2.5.1 – Aprendizado Supervisionado.....	15
2.5.2 – Aprendizado Não Supervisionado.....	16
2.5.3 – Aprendizado semi-supervisionado.....	16
2.5.4 – Machine learning de reforço.....	16
2.5.5 – Underfitting X Overfitting.....	17
2.5.6 – Rede Neural.....	17
2.5.7 – Rede neural convolucional (CNN).....	19

Capítulo 3: Propostas Tecnológicas	21
3.1 – Pesquisas sobre inteligência artificial na área da saúde	21
3.2 – Base de dados	23
3.2.1 – Análise da base de dados	24
3.3 – Modelagem rede neural Convolutacional (CNN).....	28
3.4 – Modelagem rede neural	30
Capítulo 4: Resultados Obtidos ou Esperados.....	32
4.1 – Vantagens da aplicação da IA na área da saúde	32
4.2 – Análise da base de dados	32
4.3 – Treinamento da rede neural	32
4.4 – Treinamento da rede neural convolutacional	33
4.5 – Análise de erros na previsão.....	36
Capítulo 5: Conclusão e Trabalhos Futuros	38
5.1 – Conclusão	38
5.2 – Trabalhos Futuros	39
Referências Bibliográficas	40

CAPÍTULO 1

Introdução

Muitos são os benefícios da utilização do *big data* e *machine learning* na área da saúde. Importantes institutos como a FIO CRUZ e UNICAMP estão desenvolvendo pesquisas nesse campo, e resultados bastante animadores estão sendo obtidos, com isso, soluções baseadas nessas tecnologias vem sendo empregadas com bastante sucesso.

Como exemplo pode-se citar algoritmos de predição, com a finalidade de identificar doenças. Neste trabalho, uma base de dados com exames de eletrocardiogramas, é utilizado, e duas redes neurais são parametrizadas, e testadas nesse conjunto. É utilizado o Tensor Flow e Keras, em ambiente de linguagem Python.

Ambas as redes neurais apresentaram excelentes resultados, com boa precisão. Esta pesquisa, é continuação do trabalho de conclusão de curso da turma MB3B, apresentado em 03 de junho de 2023, cujo objetivo foi avaliar a aplicação dos modelos de machine learning KNN e regressão logística para a previsão de doenças do coração e diabetes.

1.1 – Tema

Este trabalho trata de questões de saúde e tecnologia, abordando os recursos de inteligência artificial, aplicados na área da saúde, com ênfase em redes neurais, aplicados na classificação de exames de eletrocardiograma.

1.2 – Justificativa

Com a utilização de novas tecnologias, em especial IA, é possível prever doenças em segundos, e esses recursos, tem se tornado um grande aliado no campo da saúde. Os problemas cardíacos têm alta incidência na população global, sendo a principal causa de óbitos no mundo.

Muitos são os benefícios que a IA oferece na área da saúde, ajudando a diminuir custos, dando celeridade e precisão nas avaliações clínicas.

1.3 – Objetivos

Este trabalho aborda como as redes neurais podem ser utilizados na classificação de doenças, baseados nos padrões ocultos em bases de dados de exames clínicos.

Através da pesquisa bibliográfica, apresenta-se alguns conceitos relacionados ao tema e pesquisas em andamento, com elevado grau de relevância no campo tecnológico abordado.

1.3.1 – Objetivo Geral

Demonstrar a importância da inteligência artificial na área da saúde, e a aplicação de rede neural na classificação de arritmias cardíacas através do eletrocardiograma.

1.3.2 – Objetivo Específicos

1. Análise do panorama geral da saúde pública no Brasil.
2. Apresentação de recursos tecnológicos e fundamentos para compreensão deste trabalho.
3. Análise exploratória no conjunto de dados.
4. Aplicação de rede neural, para a classificação de doenças do coração, em um conjunto de dados de eletrocardiogramas e avaliação da performance.

1.4 – Delimitação

Os conceitos apresentados podem ser aplicados de modo amplo dentro da área da saúde, porém neste trabalho foi limitado a classificação de doenças cardíacas, mais específico as arritmias, através da análise de exames de eletrocardiograma. Para isso foi utilizado uma rede neural.

Foi utilizada, uma base de dados de exames de eletrocardiograma.

O desenvolvimento do modelo foi elaborado na linguagem Python.

1.5 – Metodologia

Para a apresentação do panorama atual da saúde pública, foram feitas pesquisas em sites oficiais de institutos de pesquisa e artigos. A fundamentação teórica, baseia-se em pesquisa

bibliográfica. Os conjuntos de dados utilizados, foram selecionados na plataforma Kaggle. Foram feitas análises exploratórias e higienização nos dados selecionados. Após esse processo, foi efetuado o treinamento e análise de performance dos modelos de *machine learning*. Após análise dos resultados, foi elaborada a conclusão sobre a performance atingida.

1.6 – Descrição

O presente trabalho está estruturado da seguinte maneira:

O Capítulo 2, apresenta as fundamentações teóricas que sustentam os conceitos acerca da tecnologia proposta. Também é exposto os dados quantitativos e qualitativos que justificam a realidade da saúde pública no Brasil, bem como as doenças de maior incidência e fatores associados a essas.

O capítulo 3, apresenta a análise exploratória e tratamento dos dados, além da utilização dos modelos de *machine learning*.

No capítulo 4, é analisado os resultados obtidos, e os fatores que colaboraram com esses.

O Capítulo 5 trata da conclusão.

CAPÍTULO 2

Referencial Teórico

Nesta seção são apresentados as definições e conceitos relacionados ao estudo: REDES NEURAIS APLICADO EM EXAMES DE ELETROCARDIOGRAMA. Para tanto, são apresentados conceitos importantes para o entendimento do trabalho.

2.1 – A saúde pública no Brasil

O Brasil tem o maior sistema de saúde pública do mundo. Ele foi criado a partir da forte pressão de movimentos sociais na década de 80, e na constituição federal de 1988, foi dedicado um capítulo inteiro a saúde. Está escrito que o acesso a ela deve ser universal, gratuito e igualitário a todos. O Sistema Único de Saúde (SUS), foi criado através do artigo 196 da Constituição Federal.

Segundo um levantamento feito pelo Instituto Brasileiro de Geografia e Estatística (IBGE) em 2019, cerca de 150 milhões de pessoas, que corresponde a 70% da população Brasileira, dependem exclusivamente do Sistema Único de Saúde (SUS), para obterem tratamento médico (IBGE, 2019).

O orçamento da saúde em 2022, foi previsto em R\$153.31 Bilhões, dos quais são destinados para assistência hospitalar e ambulatorial R\$36.741 Bilhões, segundo o portal da transparência¹.

Apesar do montante aplicado, a saúde pública enfrenta muitos problemas, onde podemos destacar a falta de profissionais qualificados e longas filas de espera.

A principal causa de mortes no país, e uma das principais causas no mundo, são as doenças cardíacas. Elas são a causa número 1 de mortes no País, segundo dados do SUS.

De acordo com o estudo da “Institute for Health Metrics and Evaluation” (GBD, 2019), estima-se que 6,1% da população (13.237 milhões de pessoas) possam vir a ter algum tipo de problema cardíaco.

¹ PORTAL DA TRANSPARÊNCIA. **Função 10 - Saúde. 2022.** Disponível em: <https://www.portalttransparencia.gov.br/funcoes/10-saude?ano=2022>. Acesso em: 18 jul. 2022.

2.1.1 – Doenças cardiovasculares

As doenças cardiovasculares são a principal causa de morte no Brasil, e no mundo.

É um termo que define um conjunto de doenças do coração e vasos sanguíneos:

- a) Doença cerebrovascular: doença dos vasos sanguíneos que irrigam o cérebro.
- b) Doença coronariana: doença dos vasos sanguíneos que irrigam o músculo cardíaco.
- c) Doença arterial periférica: doença dos vasos sanguíneos que irrigam os membros superiores e inferiores.
- d) Doença cardíaca reumática: danos nas válvulas cardíacas e músculo do coração, causados pela febre reumática.
- e) Trombose venosa profunda e embolia pulmonar: coágulos sanguíneos nas veias das pernas, que podem se deslocar para os pulmões e coração.
- f) Cardiopatia congênita: malformação no desenvolvimento do feto, que atinge a estrutura do coração.
- g) Arritmias cardíacas: Batimentos irregulares do coração.

Acidentes vasculares cerebrais (AVC) e ataques cardíacos normalmente são eventos agudos ocasionados na maioria das vezes por um bloqueio que não deixa o sangue fluir para o cérebro. Na maioria das vezes está relacionado ao acúmulo de gordura nos vasos sanguíneos. Os AVCs, também podem ter como causa hemorragias em vasos sanguíneos do cérebro.

Segundo a Organização Pan-Americana de saúde (PAHO, 2016), a maioria dos problemas cardíacos estão relacionados a fatores de risco como: pressão alta, cigarros, sedentarismo, alimentação ruim, álcool, entre outros.

2.1.2 – Eletrocardiograma

Trata-se de um método diagnóstico de simples execução, bastante útil nos diagnósticos das doenças cardiovasculares, principalmente as agudas, como o infarto agudo do miocárdio e arritmias. A base da explicação desse exame é voltada para conceitos básicos da física².

Isso se dá, porque o eletrocardiograma, também chamado de ECG ou eletrocardiografia é um exame que avalia a atividade elétrica do coração por meio de eletrodos fixados na pele.

² RBEM. **Nova Metodologia de ensino do ECG: Desmitificando a Teoria na Prática – Ensino Prático do ECG.** Disponível em: <https://www.scielo.br/j/rbem/a/RXbsLmvxHH9jG7H9NFWRdJB/?lang=pt>. Acesso em 15 jul 2023

Segundo o laboratório Lavoisier³, através desse exame, pode-se diagnosticar vários problemas, dos quais é citado:

- Irregularidades no ritmo cardíaco (arritmia), seja por um coração acelerado (taquicardia), devagar (bradicardia) ou fora do ritmo;
- Aumento de cavidades cardíacas;
- Patologias coronarianas;
- Infarto do miocárdio;
- Distúrbios na condução elétrica do órgão;
- Problemas nas válvulas do coração;
- Pericardite - Inflamação da membrana que envolve o coração;
- Hipertrofia das câmaras cardíacas - átrios e ventrículos;
- Doenças que isolam o coração - derrame pericárdico ou pneumotórax;
- Infarto em situações emergenciais;
- Doenças genéticas;
- Doenças transmissíveis (Doença de Chagas).

Através desse exame, também é possível o monitoramento de dispositivos como o marca-passos e avaliar se algum medicamento está causando efeito no coração⁴.

O eletrocardiógrafo, é um galvanômetro que registra a atividade elétrica em um papel milimetrado, através de vários eletrodos, posicionados sob a pele, em pontos específicos do corpo, registrando a atividade elétrica, no domínio do tempo.

Analisando os padrões desta atividade, é possível detectar anomalias no funcionamento do órgão.

Basicamente, o aparelho de ECG tem 10 eletrodos, onde 04 são colocados nas extremidade dos braços e pernas, e os outros 06, denominados precordiais (V1 a V6), são colocados ao longo do peito, em posições pré-determinadas⁵. As figuras 2.1.2.1 e 2.1.2.2, demonstram os pontos que os eletrodos precordiais devem ficar e a correlação destes com o coração.

³ ELETROCARDIOGRAMA. **O que é e como é feito o exame ECG.** Disponível em: <https://lavoisier.com.br/saude/blog/eletrocardiograma>. Acesso em 15 jul 2023

⁴ Idem.

⁵ ECG. **Você sabe colocar os eletrodos do ECG de forma correta? Tem certeza? Disponível em:** <https://cardiopapers.com.br/curso-basico-de-eletrocardiograma-parte-04/>. Acesso em 20 jul 2023.

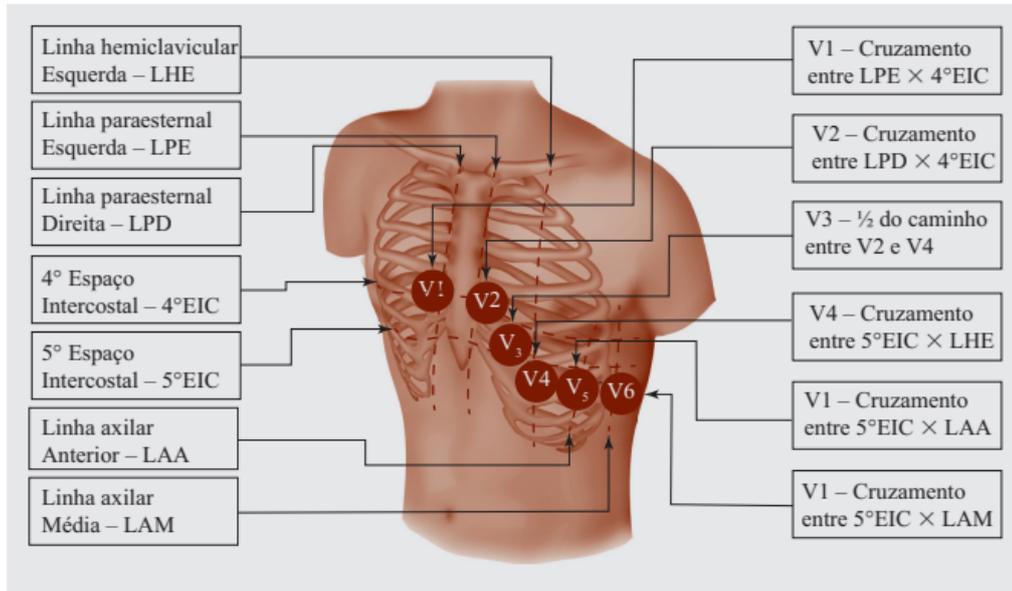


Figura 2.1.2.1 – Posicionamento das derivações precordiais.

Fonte: ECG-Manual prático de eletrocardiograma – Hcor, p. 21, 2021.

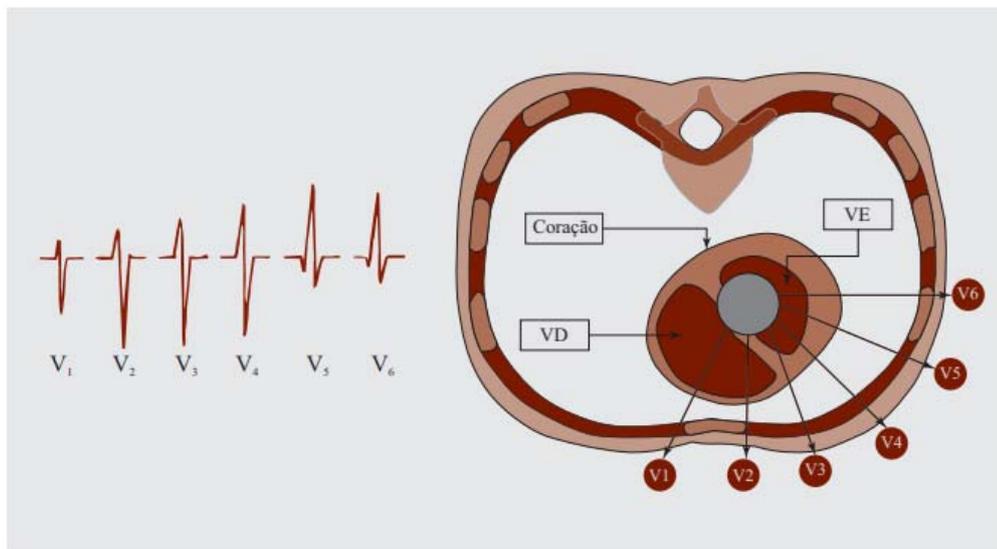


Figura 2.1.2.2 – Corte transversal do torax, correlações dos eletrodos e a região do coração.

Fonte: ECG-Manual prático de eletrocardiograma – Hcor, p. 21, 2021.

O coração possui 04 cavidades, sendo os átrios e ventrículos. Sua função é impulsionar o sangue que vem através das veias para os pulmões, e desses para todos os órgãos (sangue rico em oxigênio). Através de pulsos elétricos é determinado o ritmo cardíaco e sincronizado o batimento das 04 câmaras do coração. estes pulsos são detectados pelo ECG, e analisando o padrão destes, é possível avaliar a atividade do coração.

Um exame de eletrocardiografia, representa a atividade elétrica do coração de forma gráfica, composto pelas ondas P, Q, R, S e T. O intervalo é repetido em cada ciclo de batimento cardíaco, e cada onda tem relação com a despolarização e repolarização.

A onda P corresponde a despolarização (contração) dos átrios.

O intervalo PR corresponde ao intervalo da despolarização dos átrios e dos ventrículos.

O Complexo QRS corresponde a despolarização (contração) dos ventrículos.

O Segmento ST corresponde ao intervalo de tempo entre a despolarização e o início da repolarização dos ventrículos.

A onda T corresponde a repolarização dos ventrículos, que se tornam aptos a contrair novamente. A figura 2.1.2.3 demonstra como essas ondas são apresentadas graficamente:

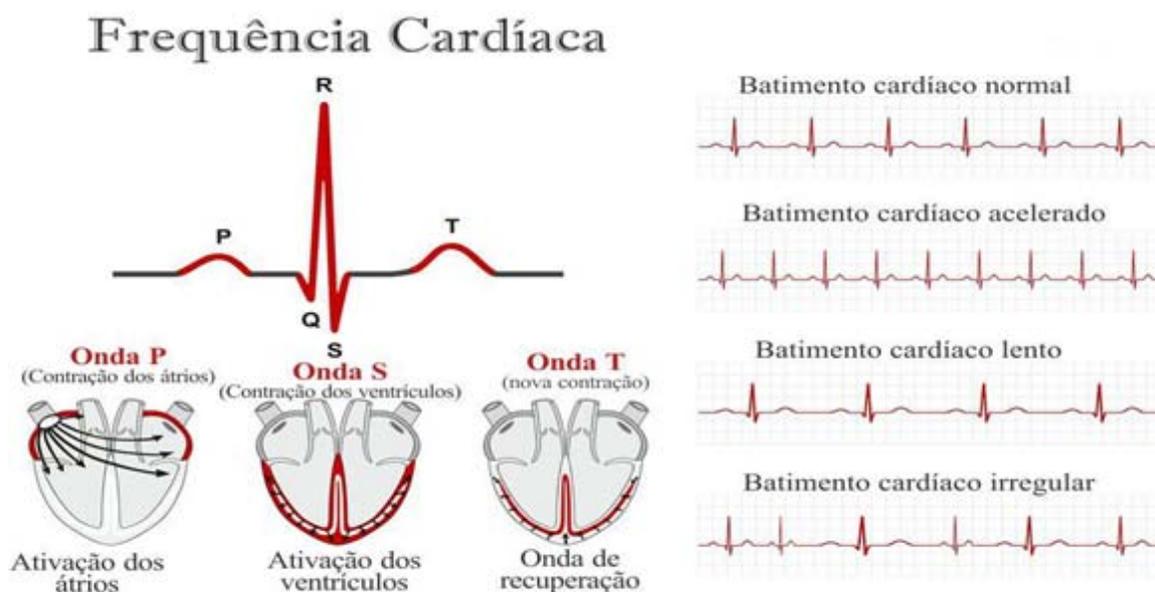


Figura 2.1.2.3 – Representação gráfica do batimento cardíaco no ECG.

Fonte: Arritmia Cardíaca e morte súbita, 2019.

A base de dados selecionada neste trabalho aborda 05 tipos de batimentos cardíacos:

- Normal:** Demonstrado na Figura 2.1.2.4, é o padrão de batimento de uma pessoa saudável.

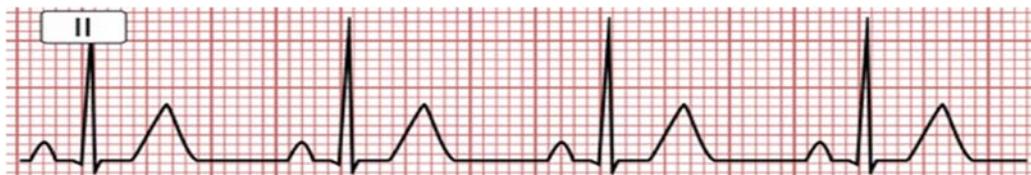


Figura 2.1.2.4 – Batimento cardíaco normal no ECG.

Fonte: ECGNOW, 2021.

- b) **Desconhecido:** Trata-se de padrões não abordados no trabalho, ou com alguma distorção ou ruído.
- c) **Batimento ventricular prematuro ou extrassístole ventricular:** É um batimento cardíaco adicional, ocasionado por uma ativação elétrica anormal, com início no ventrículos, antes do batimento normal. Elas se caracterizam pela aparição de um complex QRS largo, com morfologia aberrante (MY EKG – arritmias, 2023). A figura a seguir representa forma de onda característica, e pode-se observar a pausa pós extrassistólica em azul, na Figura 2.1.2.5.

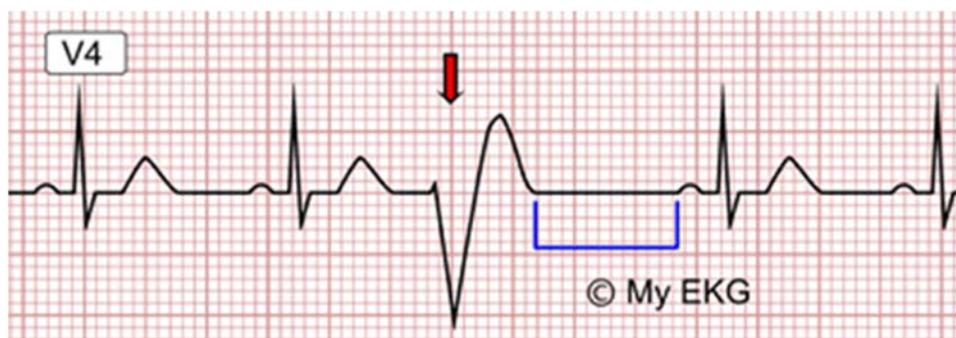


Figura 2.1.2.5 – Batimento cardíaco extrassístole ventricular no ECG.

Fonte: MY EKG – arritmias, 2023.

- d) **Batidas prematuras supraventriculares ou extrassístole atrial:** Tem origem nas câmaras superiores (átrios). Esse tipo de arritmia, na maioria das vezes é benigna, não necessitando de preocupações maiores. Trata-se de um batimento cardíaco extra, precoce em relação ao batimento precedente. A figura a seguir, representa o padrão desta arritmia no ECG. O pulso elétrico APB, representa a batida prematura que ocorre neste caso (ECGNOW, 2021), conforme demonstrado na figura 2.1.2.6.

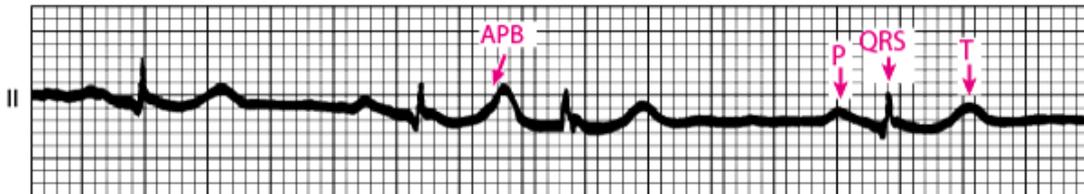


Figura 2.1.2.6 – Batimento cardíaco extrassístole atrial no ECG.

Fonte: MANUAL MSD, 2023.

- e) **Batida de fusão ventricular:** Demonstrado na figura 2.1.2.7, esse tipo de arritmia é caracterizado pela presença de batimentos tipicamente de origem ventricular. Normalmente esse fenômeno, representa a junção de duas frentes de onda, sendo uma do foco ventricular, e a outra do sistema normal de condução. Esse tipo de arritmia, é uma batida precoce em relação ao ciclo da taquicardia ventricular (SBC – Seção de Eletrocardiograma, 2014).

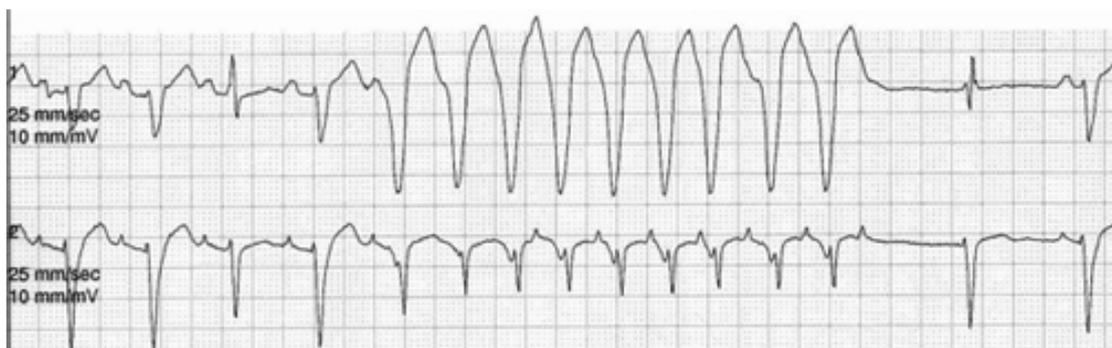


Figura 2.1.2.7 – Batimento de fusão ventricular no ECG.

Fonte: SBC - Seção de Eletrocardiograma, 2014.

2.2 – Big Data

O big data, refere-se a um enorme volume de dados, com elevada variedade e gerados em velocidade, tornando esse conjunto complexo, cujo técnicas tradicionais de processamento, não são eficientes.

Essa questão da Variedade, Volume e Velocidade são conhecidos com os três V's do big data. Um dos desafios, é referente ao armazenamento desses dados. Para se ter ideia, estima-se que esse volume dobre a cada 2 anos. As empresas precisam de maior elasticidade para armazenarem esses dados, e a solução para isso é a computação em nuvem. Outro desafio se refere a qualidade desses dados, onde a pluralidade de fontes acaba não permitindo uma

padronização, e esses dados para serem utilizados demandam processos de limpeza e estruturação.

O big data se tornou essencial, pois as empresas conseguem agregar valor em seus negócios através de análises de dados que respondem diversas perguntas sobre os negócios, principalmente referente ao comportamento dos consumidores. Os *insights* são gerados através dessas informações (ORACLE, 2022).

2.3 - DATASUS

Criado em 16 de abril de 1991, através do decreto nº 100, ele é o departamento de informática do SUS. Com o crescente volume de dados gerados na área da saúde, percebeu-se rapidamente a necessidade do emprego da tecnologia. Até hoje esse departamento foi responsável pela criação de mais de 200 sistemas para o ministério da saúde.

Seu objetivo principal é coletar, processar e disseminar dados sobre saúde no país. Uma característica importante é que ele permite o cruzamento com outras bases de dados como o IBGE, dessa maneira permite análises que direcionam o planejamento de políticas de saúde pública. É nesse contexto que se torna evidente a importância da CIENCIA DE DADOS.

O DATASUS⁶, oferece ferramentas de *Business Intelligence* (inteligência de negócios), onde os especialistas na área de saúde podem fazer pesquisas. Porém existe uma grande falta de informações por parte das prefeituras, e muitas informações são desatualizadas.

Além da desatualização, a falta de padronização também é um problema a ser superado. O poder público tem elaborado políticas voltadas a superar esses problemas através da unificação e padronização das bases de dados, além da melhor integração destas.

Através dessas bases, é possível obter indicadores para entender padrões no desenvolvimento de muitas doenças, para assim criar modelos capazes de fazer previsões cada vez mais precisas, e produzir conhecimento científico.

2.4 – Ciência de dados

Segundo definição da Amazon, é o estudo dos dados para a geração de *insights* relevantes para os negócios. É uma prática multidisciplinar, que utiliza matemática, estatística,

⁶ DATASUS. **Sobre o Datasus**. Disponível em: <https://datasus.saude.gov.br/sobre-o-datasus/>. Acesso em: 20 jul. 2022.

inteligência artificial e engenharia da computação. Embora seu conceito não seja novo, nos últimos anos tem ganhado popularidade, devido a sua importância nos negócios.

É um conceito presente em todas as áreas, e se faz necessário, com o aumento exponencial na geração de dados. Através do estudo das características e correlação de variáveis, esses dados dão pistas sobre o que se busca responder (AMAZON, 2022).

2.4.1 – Análise exploratória

Normalmente é a fase inicial do processo de estudo. Basicamente é aplicada métodos de estatística descritiva para organizar, resumir e descrever os principais aspectos dos dados, através de gráficos, tabelas ou narrativas revelando aspectos importantes dos dados (AMAZON, 2022).

2.4.2 – Análise diagnóstica

É uma análise mais profunda para entender o que houve. Utiliza técnicas como drill-down, data mining e correlações. Podem ser aplicados métodos de transformação de dados (AMAZON, 2022).

2.4.3 – Análise preditiva

Basicamente este tipo de análise visa fazer previsões, fundamentando-se em dados passados. Para isso utiliza-se de modelos de *machine learning*, correspondência de padrões e modelagem preditiva (AMAZON, 2022).

2.4.4 – Análise prescritiva

Trata-se de um patamar mais elevado do que os dados preditivos. Além das previsões, sugere uma resposta adequada para tal resultado. Pode analisar diferentes cenários e recomendar a melhor a ação a ser tomada (AMAZON, 2022).

2.4.5 – Google Colab

O Google Colab, ou Google Collaboratory, é um serviço de processamento em nuvem, oferecido pela Google, voltado à criação e execução de scripts em linguagens como Python ou R. Amplamente difundido, não necessita de instalação, pois roda direto no navegador.

2.4.6 – Python

Conforme definido pelos próprios criadores, é uma linguagem interpretada, interativa e orientada a objetos. Sintaxe clara e objetiva, tornam fácil sua compreensão. Ele incorpora módulos, exceções, tipagem dinâmica, tipos de dados dinâmicos de alto nível e classes. O Python suporta vários paradigmas de programação, como orientação a objetos, estruturada ou procedural e funcional (PYTHON, 2022). É muito difundida atualmente, e bastante utilizada no campo de ciência de dados. É uma das linguagens utilizada pelo SERPRO⁷ (maior empresa de TI pública do mundo).

2.4.7 – Pandas

É uma biblioteca para análise e manipulação de dados. Desenvolvido na linguagem Python, é uma ferramenta intuitiva, de fácil aprendizado e de código aberto (PANDAS, 2022).

2.4.8 – Numpy

É uma biblioteca para linguagem Python, que suporta o processamento da matrizes e arranjos multidimensionais. Possui variados conjunto de funções matemáticas de alto nível para operar sobre os arranjos e matrizes (NUMPY, 2022).

⁷ SERPRO. Disponível em <https://www.serpro.gov.br/>. Acesso em 04 mai. 2023

2.4.9 – Seaborn

É uma biblioteca utilizada para visualização de dados em Python. Baseado em matplotlib, é uma interface de alto nível que oferece recursos gráficos atraentes e informativos para visualização de informações estatísticas (SEABORN, 2022).

2.4.10 – Matplotlib

Matplotlib é uma biblioteca abrangente para criar visualizações estáticas, animadas e interativas em Python (MATPLOTLIB, 2022).

2.4.11 – Scikit-learn

É uma biblioteca destinada à análise preditiva de dados, sendo de código aberto para a linguagem de programação Python. Ela dispõe de ferramentas simples e eficientes, e neste trabalho é utilizada para o pré-processamento de dados (SCIKIT-LEARN, 2022).

2.4.12 – TensorFlow

É uma biblioteca de código aberto cuja finalidade é o aprendizado de máquina, computação numérica entre outros. Desenvolvido pelo Google, é muito utilizada, sendo um dos principais recursos de *machine learning* e *deep learning* (TENSORFLOW, 2022).

2.4.13 – Inteligência artificial

Segundo a Oracle, a inteligência artificial (IA), refere-se a sistemas ou máquinas que imitam a inteligência humana para realizar tarefas e podem se aprimorar iterativamente com base nas informações que coletam. Ela se manifesta de várias formas” (ORACLE, 2022).

Inteligência artificial é um conceito amplo, mas basicamente todos remetem a ideia de executar tarefas de forma independente, e capacidade de aprendizagem. Esse conceito se popularizou com o desenvolvimento e aprimoramento das técnicas de machine learning. Neste trabalho, demonstraremos como machine learning pode auxiliar no campo da saúde. A

helthtech brasileira Laura⁸, tem se destacado no mercado oferecendo soluções na área médica baseados em IA.

2.5 – Machine learning

É um ramo da IA e ciências da computação, que na sua essência utiliza os dados e algoritmos para efetuar tarefas, de maneira automatizada. Na medida que a base de dados utilizada no treinamento aumenta, o algoritmo vai melhorando sua precisão gradualmente, em alguns casos com recalibração do modelo de forma total, para assim revalidar a precisão de todo o modelo.

Através da utilização de métodos estatísticos, os algoritmos de *machine learning*, são treinados para fazerem previsões ou classificações. Trata-se de um poderoso recurso computacional, presente hoje em quase tudo que envolve IA.

Com base nos dados de entrada, rotulados ou não, o algoritmo produzirá uma estimativa sobre um padrão identificado. Uma função de erro oferece uma métrica para avaliar a precisão dos resultados. Na função do *Machine learning*, existem variáveis que são ajustados de modo a melhorar a precisão das análises. Estas podem ser classificadas quanto a “importância” ao modelo a ser aplicado sendo assim é de extrema importância que seja identificado o ciclo de vida destas.

Alguns casos de usos que podemos citar a utilização de *machine learning* são nas aplicações de reconhecimento da fala, *Bots*, mecanismos de recomendação, visão computacional, entre outros.

Basicamente, existem três categorias de classificadores, aprendizado supervisionado, aprendizado não supervisionado e aprendizado por reforço (IBM. Aprendizado de máquina, 2022).

2.5.1 – Aprendizado Supervisionado

Caracteriza-se pela utilização de dados rotulados que fazem a classificação ou predição com precisão.

⁸ LAURA. **Blog**. Disponível em <https://laura-br.com/blog/>. Acesso em:04 mai. 2023.

Na medida que a entrada de dados vai sendo alimentada com novos dados, esse método ajusta seus pesos buscando melhorar a precisão dos resultados na saída. Isso faz parte do processo de validação cruzada, cujo intuito é evitar o *underfitting* ou *overfitting*.

São exemplos desse modelo os algoritmos: máquinas de vetores de suporte (SVM), redes neurais, regressão linear, *naive bayes*, regressão logística, floresta aleatória e muitos outros (IBM, 2022).

2.5.2 – Aprendizado Não Supervisionado

A principal característica que difere do modelo anterior é a utilização de dados não rotulados. Esses algoritmos descobrem padrões ocultos ou agrupamentos nos dados considerados, sem a ajuda humana. A capacidade de descobrir diferenças e semelhanças nas informações, fazem com que esse modelo seja ideal para a análise exploratória de dados, segmentação, visão computacional, reconhecimento de padrões entre outros.

Também utilizado para redução de recursos, em um modelo através do processo de redução de dimensionalidade. Duas abordagens comuns para isso são a análise de componente principal e decomposição de valor singular.

Também são algoritmos desse modelo: redes neurais, armazenamento em cluster de k-médias, armazenamento em cluster probabilísticos, entre outros (IBM, 2022).

2.5.3 – Aprendizado semi-supervisionado

É um meio termo entre os dois modelos citados anteriormente. Para o treinamento é usado um conjunto de dados rotulado menor para orientar a classificação e extração de recursos de um conjunto maior, porém de dados não rotulados (IBM, 2022).

2.5.4 – Machine learning de reforço

É um modelo de *machine learning* comportamental, parecido com o modelo de aprendizado supervisionado, porém o algoritmo não é treinado usando dados de amostra.

Este é um modelo bastante interessante. Ele tem a capacidade de mapear uma série de entradas para saídas com dependências e não apenas uma entrada para uma saída (IBM, 2022).

Trata-se de um modelo bastante interessante, e um caso de uso conhecido é o Watson da IBM.

2.5.5 – Underfitting X Overfitting

Costuma-se dizer que o modelo está em *underfitting* (subajustando), quando este tem um desempenho insatisfatório nos dados de treinamento. Isso ocorre quando não é possível capturar o padrão dos dados na entrada e reproduzir na saída.

Quando o modelo está em *overfitting* (*superajustando*), é quando ocorre um excelente desempenho nos dados de treinamento, porém baixo nos dados de avaliação. Isso ocorre porque o modelo memorizou os dados, e não consegue generalizar em dados não vistos anteriormente (IBM. Aprendizado de máquina, 2022). Na figura 2.5.5 é demonstrado como isso ocorre.

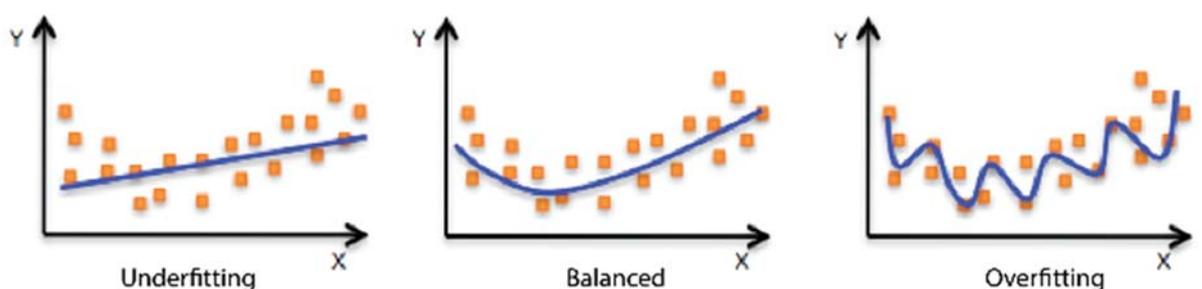


Figura 2.5.5.1 – Underfitting x Overfitting

Fonte: Amazon WorkDocs, 2022

2.5.6 – Rede Neural

É um processo de *machine learning*, também conhecido como *deep learning* (aprendizado profundo). Trata-se de um tipo de IA, inspirado no cérebro humano. Ela utiliza “nós” ou neurônios artificiais, interconectados em uma estrutura de camadas (AMAZON, 2022).

Um neurônio artificial ou matemático é inspirado no humano. Ele é um elemento da rede neural artificial e seu funcionamento ocorre da seguinte maneira:

- a) Calcula a soma ponderada de vários inputs.
- b) Aplica uma função e passa o resultado adiante.

A figura 2.5.6.1, ilustra o funcionamento de um neurônio artificial:

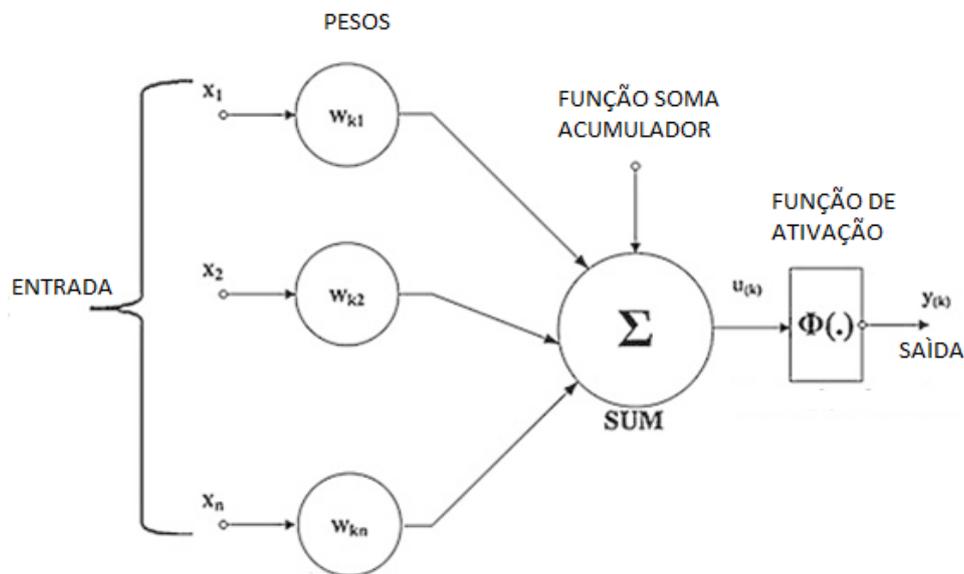


Figura 2.5.6.1 – Funcionamento básico do neurônio artificial

Fonte: Deep Learning Book – o neurônio biológico e matemático, 2022.

Uma rede neural cria um sistema adaptativo, o qual as máquinas utilizam para aprender com os erros. Uma rede neural básica, possui neurônios interconectados em 03 camadas:

1- Camada de entrada:

Por onde são inseridas as informações (dados) que entram na rede neural. Os nós de entrada processam os dados, analisam ou categorizam esses dados e os encaminham para a próxima camada.

2- Camada oculta:

As camadas ocultas utilizam as entradas da camada de entrada ou de outras camadas ocultas. As redes neurais artificiais podem ter muitas camadas ocultas. Cada camada analisa o resultado da camada anterior, processa-o mais um pouco e o encaminha para a próxima camada.

3 – Camada de saída:

Esta camada fornece o resultado de todos os dados processados na rede neural. Ela pode ter um ou vários nós. Se tiver um problema de classificação binária (sim ou não), a camada de saída terá um nó de saída, que fornece o resultado como 1 ou 0. Porém, se for um problema de classificação de várias classes, a camada de saída pode ter mais de um nó de saída.

Em sua arquitetura, uma rede neural profunda, ou rede de aprendizado profundo, têm várias camadas ocultas, com milhões de neurônios (nós) interligados. Uma variável, conhecida como peso, representa as conexões entre um nó e outro. O ‘peso’ será um número positivo se

um nó excitar o outro, ou negativo se um nó reprimir o outro. Os nós com “pesos” maiores têm mais influência nos outros nós.

Em tese, uma rede neural profunda pode direcionar qualquer tipo de entrada para qualquer tipo de saída. Entretanto, ela precisa de muito mais treinamento do que outros métodos de *machine learning*. Ela precisa de milhões de exemplos de dados de treinamento, enquanto uma rede simples talvez precise de apenas centenas ou milhares!

A figura 2.5.6.2, ilustra como essa arquitetura funciona:

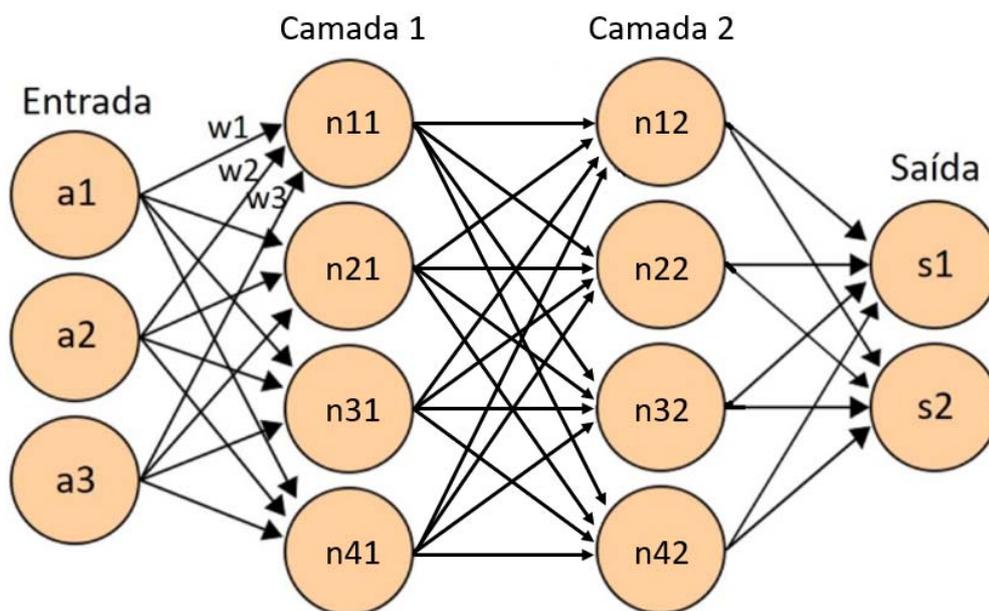


Figura 2.5.6.2 – Funcionamento básico das redes neurais

Fonte: Didática Tech, 2022.

2.5.7 – Rede neural convolucional (CNN)

É um tipo de rede neural bastante utilizada em visão computacional. É um algoritmo de *deep learning* que pode captar imagens na entrada, atribuir vieses que podem ser aprendidos a muitos aspectos diferentes, assim sendo possível diferenciar ou encontrar padrões ocultos nas imagens. Possui uma abordagem mais dimensionável para esse tipo de tarefa, pois possui desempenho superior no tratamento de imagem.

A figura 2.5.7.1, demonstra a arquitetura básica de uma rede neural convolucional.

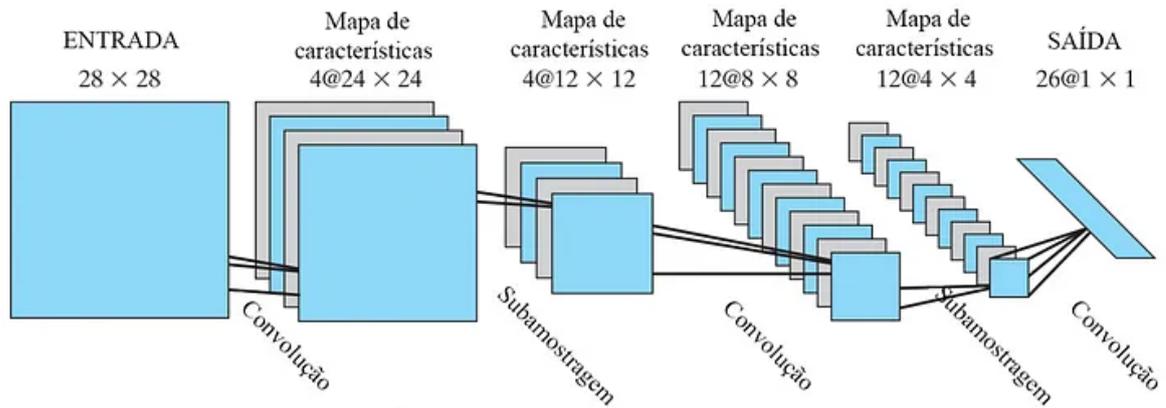


Figura 2.5.7.1 – Arquitetura de rede neural convolucional

Fonte: Medium, 2022.

A imagem é captada na camada de entrada e transformada numa matriz multidimensional, onde cada pixel é um valor e corresponde a um neurônio. Cada camada a seguir, é responsável por extrair determinadas informações, onde a informação flui através de cada camada da rede, com a camada anterior fornecendo a entrada para a camada seguinte. Essas são as camadas convolucionais, que captam as características mais marcantes e vai reduzindo a dimensionalidade a cada camada, até chegar na camada de saída, onde cada neurônio corresponderá a uma variável de saída.

CAPÍTULO 3

Propostas Tecnológicas

Neste trabalho é realizada análise breve das principais pesquisas em andamento sobre aplicações do *Big Data* e Inteligência Artificial na área da saúde.

Depois, demonstra-se a aplicação das redes neurais na classificação de arritmias cardíacas.

As informações apresentadas, são frutos de pesquisa referencial e bibliográfica em sites oficiais do governo, importantes institutos de pesquisas e gigantes da tecnologia. Os conjuntos de dados utilizados, foram obtidos na plataforma Kaggle. Esta é uma importante plataforma que permite a busca ou publicação de conjuntos de dados.

3.1 – Pesquisas sobre inteligência artificial na área da saúde

A faculdade de saúde pública da USP, possui o laboratório LABDAPS⁹, cujo objetivo é realizar pesquisas na área de *big data* e *machine learning* no campo da saúde. Atualmente possui pesquisas em andamento sobre:

- a) Algoritmos de *machine learning* para prever a expectativa de vida de municípios brasileiros.
- b) Predição de desfechos negativos de covid-19 em um hospital de São Paulo.
- c) Algoritmos de *machine learning* para prever a causa básica de óbito de uma amostra representativa de idosos do Município de São Paulo.

No Hospital São Lucas, pertencente a PUCRS, a equipe médica da radiologia, tem realizado estudos, com a aplicação da inteligência artificial na prática clínica, mais especificamente na área de tórax e neurorradiologia. Recentemente, foram realizadas pesquisas de revisão sistemática e metanálise sobre a performance do diagnóstico da inteligência artificial na detecção de câncer de pulmão e sobre o uso de inteligência artificial para prever o risco de ventilação mecânica no cenário de COVID-19¹⁰.

⁹ LABDAPS. Disponível em: <https://www.fsp.usp.br/labdaps/>. Acesso em 23 jul 2023.

¹⁰ PUCRS. **Inteligência artificial na medicina**. Disponível em: <https://www.pucrs.br/blog/inteligencia-artificial-na-medicina/>. Acesso em 23 jul 2023.

O uso do *big data* tem crescido em todas as áreas da ciência, e na saúde não poderia ser diferente. Existem três áreas importante da aplicação na saúde sendo: medicina de precisão, prontuários eletrônicos do paciente e internet das coisas (IOT).

Na medicina de precisão, o *big data* tem importante papel pois os conhecimentos científicos são baseados em grandes médias.

Pelo que temos presenciado nesse campo, tudo indica que a próxima fronteira da saúde será a análise do *big data*. O aumento nos estudos multicêntricos e cobrança na transparência dos gastos públicos tem fomentado esse campo (CHIAVEGATTO FILHO, 2015).

Alguns importantes institutos de pesquisa, como a Fiocruz, através do Instituto de Comunicação e Informação Científica e Tecnológica em saúde (ICICT), vem desenvolvendo pesquisas sobre aplicações do *big data* e ciência de dados em saúde nos campos (FIO CRUZ, 2022):

- a) Análise preditiva e algoritmos para mineração de dados e textos.
- b) Análise visual de dados para tomada de decisão em saúde.
- c) Infraestrutura, armazenamento e governança de dados em ecossistema Hadoop.

O prof. José Eduardo Krieger, cita em artigo publicado na revista FAPESP (2022), que o avanço da IA no campo da medicina, está no volume e qualidade dos dados.

No mesmo artigo, Krieger fala que os estudos do Instituto do Coração (InCor), no campo da IA, iniciaram no *big data*. O InCor é todo digital e utiliza um sistema de prontuário eletrônico, que alimenta um banco de dados com registro de 1.3 milhões de pacientes, ainda em 2020. Mais de 30 hospitais compartilham desse sistema, que permite o acesso a informações de mais de 10 milhões de registros, por parte de um grupo de pesquisadores, respeitando-se a privacidade dos pacientes, preservando a identidade destes, em conformidade com a lei geral de proteção de dados (LGPD). Foi através dessa base de dados, que foi possível alguns avanços nas pesquisas de IA.

O InCor é sede do Instituto Nacional de Ciências e Tecnologia em Medicina Assistida por Computação Científica (INCT-Macc). Articulam-se por meio desse instituto 31 laboratórios em 11 estados brasileiros e outros 17 com sede no exterior, distribuído em 7 países. No InCor, existe em andamento frentes de pesquisa nas áreas:

- a) Processamento de imagens
- b) Processamento de sinais e linguagem
- c) Integração de dados nas ciências ômicas (genômica, proteômica, metabolômica, entre outras).

Neste artigo também é citado que alguns modelos de redes neurais, são capazes de detectar anomalias que passam despercebidas ao ser humano. Atribui-se esse resultado a capacidade de criar correlações que ainda não foram feitas pela medicina e o desafio dos pesquisadores é tentar descobrir o caminho que as redes neurais fizeram para chegar a essas correlações.

Também é citado um projeto da Universidade Estadual de Campinas (UNICAMP), conduzido pelo Laboratório Aterolab da faculdade de Ciências Médicas (FCM), cujo objetivo é prever quais pacientes com doenças coronarianas podem vir a sofrer eventos clínicos, e são aplicados a eles algumas medidas de prevenção. Este estudo aponta que a utilização desses recursos tecnológicos pode prover uma economia de aproximadamente R\$50 milhões ao ano, para o SUS, somente nos gastos com pacientes cardíacos (FAPESP, 2022).

Em seu artigo “Inteligência Artificial e *Machine Learning* em Cardiologia – Uma mudança de Paradigma”, o prof. Claudio Tinoco Mesquita cita como a inteligência artificial tem auxiliado na análise de exames cardíacos. Através do reconhecimento de padrões eletrocardiográficos um algoritmo de machine learning, apresentou acurácia de 88% para classificação de arritmias. Esse percentual não foi maior devido a qualidade do sinal eletrocardiográfico. Cita também a importância de se construir “boas” bases de dados para o desenvolvimento de novos algoritmos. Também é colocado que a IA não visa substituir o médico, mas assistir estes de modo a diminuir erros cognitivos (MESQUITA, 2017).

3.2 – Base de dados

Diante da dificuldade de se obter dados disponíveis e de qualidade relacionados a eletrocardiogramas, recorreu-se a plataforma Kaggle, e foi utilizado a base de dados “ECG Heartbeat Categorization Dataset¹¹”.

Esta Base de dados é composta de duas coleções de exames de eletrocardiografia, ambos obtidos de dois conjuntos populares na classificação de batimentos cardíacos, o “MIT-BIH Arrhythmia Dataset” e “The PTB Diagnostic ECG Database”. Neste trabalho utilizou-se apenas o “MIT-BIH Arrhythmia Dataset”. Este conjunto tem quantidade de exames suficiente para treinar uma rede neural.

Os dados desta base correspondem aos seguintes tipos de batimentos cardíacos:

¹¹ KAGGLE. **Ecg heartbeat categorization dataset**. Disponível em: <https://www.kaggle.com/datasets/shayanfazeli/heartbeat>. Acesso em 26 jul 2023.

- a) Batidas normais
- b) Batidas desconhecidas
- c) Batidas prematuras supraventriculares ou extrassístoles atrial
- d) Batidas ventriculares prematuras ou extrassístoles ventricular
- e) Batidas de fusão ventricular

Esses dados são pré-processados e segmentados, tendo um rótulo, identificando a classe a qual estes batimentos pertencem. Os *scripts* foram desenvolvidos em linguagem Python no ambiente de processamento em nuvem google colab.

A figura 3.2.1, apresenta a forma que estão os dados armazenados. São 188 colunas, com os valores numéricos que correspondem aos pontos do gráfico do ECG, sendo a coluna 187, o rótulo referente a classificação. O total de amostras, consolidado os conjuntos de treino e teste, somam 109446 registros.

```
test_df.head()
```

	0	1	2	3	4	5	6	7	8	9	...	178	179	180	181	182	183	184	185	186	187
0	1.000000	0.758264	0.1111570	0.000000	0.080579	0.078512	0.066116	0.049587	0.047521	0.035124	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
1	0.908425	0.783883	0.531136	0.362637	0.366300	0.344322	0.333333	0.307692	0.296703	0.300366	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
2	0.730088	0.212389	0.000000	0.119469	0.101770	0.101770	0.110619	0.123894	0.115044	0.132743	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
3	1.000000	0.910417	0.681250	0.472917	0.229167	0.068750	0.000000	0.004167	0.014583	0.054167	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.570470	0.399329	0.238255	0.147651	0.000000	0.003356	0.040268	0.080537	0.070470	0.090604	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

5 rows x 188 columns

Figura 3.2.1 – Organização dos dados

Fonte: Elaborado pelo autor, 2023.

3.2.1 – Análise da base de dados

Esse conjunto de dados possui 188 colunas, com valores que correspondem aos pontos da onda no gráfico. No quadro 3.2.1.1, é informado os números de amostras e distribuição delas, no conjunto total de dados.

Quadro 3.2.1.1 – Distribuição do conjunto de dados”

		treino	teste
Tipo Batimento (arritmia)	variável	quantidade	quantidade
Batida Normal	0	72471	18118
Batida Desconhecida	4	6431	1608
Extrassístoles ventricular	2	5788	1448
Extrassístoles atrial	1	2223	556
Batida de fusão ventricular	3	641	162
Soma parcial de amostras		87554	21892
Total amostras			109446

Fonte: Elaborado pelo autor, 2023.

Para as análises e manipulações da base dados, foram utilizadas as bibliotecas Pandas e Numpy. Para as visualizações utilizou-se Matplotlib e Seaborn. Os gráficos 3.2.1.1, ao 3.2.1.4, demonstram a distribuição dos dados, entre as amostras de treino e teste.

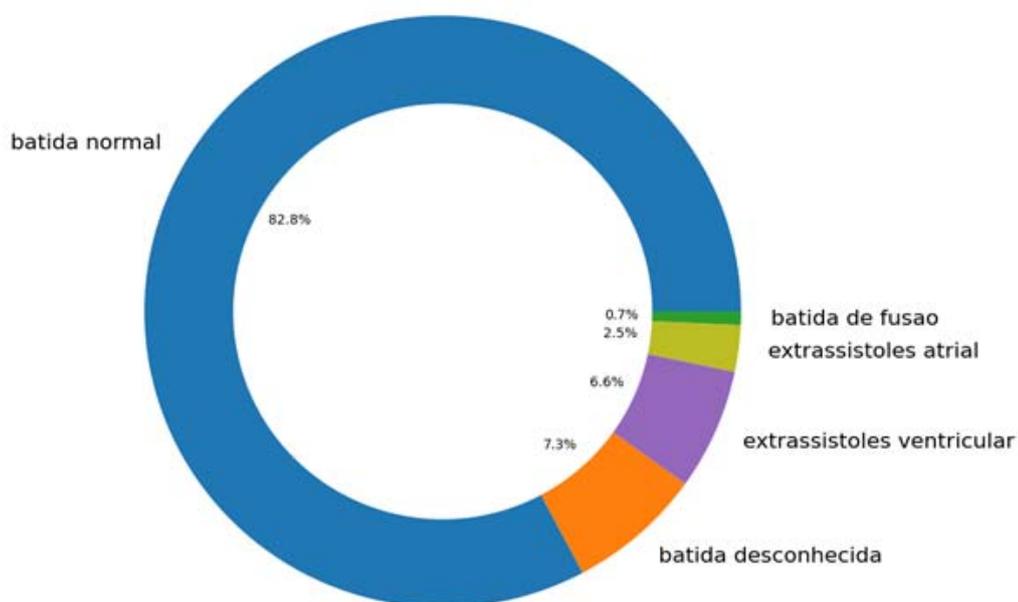
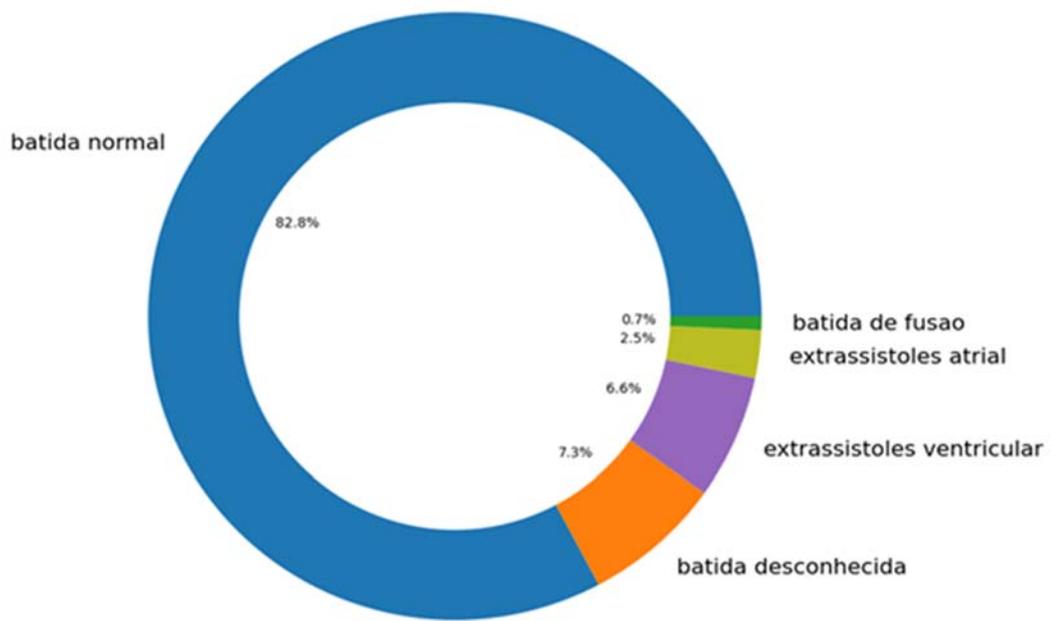


Figura 3.2.1.1 – Distribuição das amostras de treino

Fonte: elaborado pelo autor, 2023.



Conjunto de teste

Figura 3.2.1.2 – Distribuição das amostras de teste

Fonte: elaborado pelo autor, 2023.

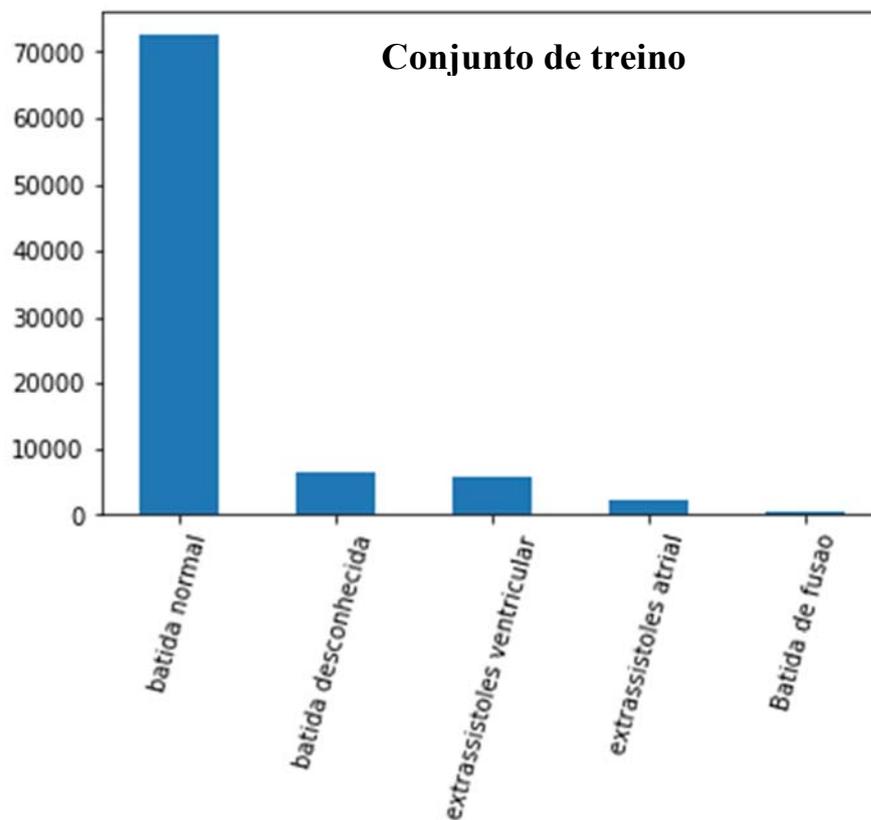


Figura 3.2.1.3 – Amostras de treino (barras)

Fonte: elaborado pelo autor, 2023.

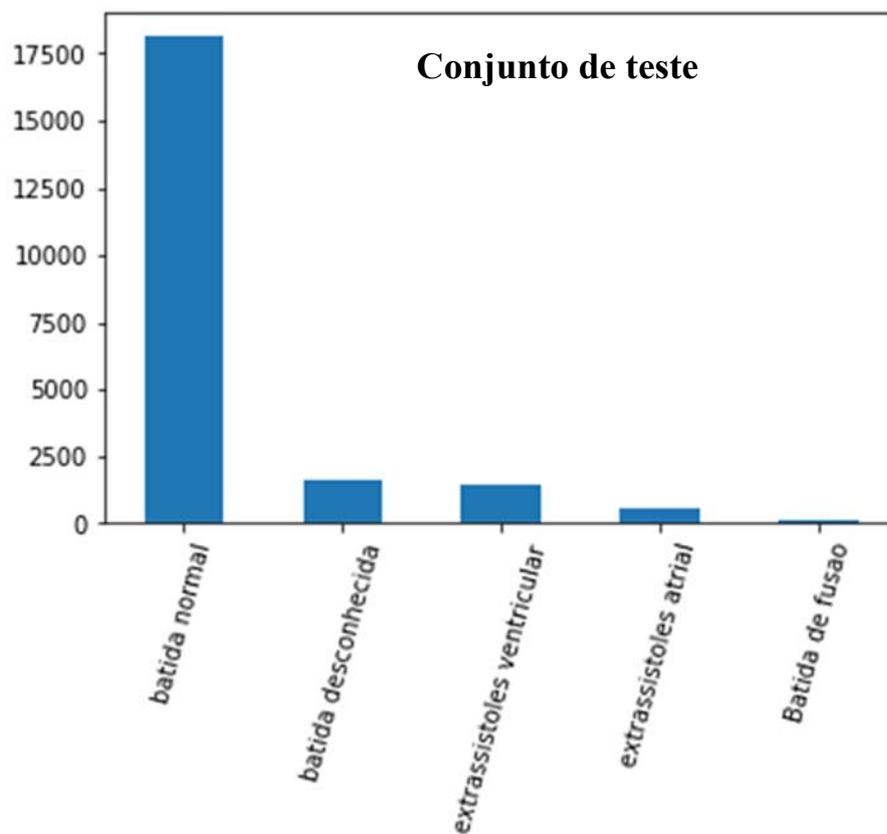


Figura 3.2.1.4 – Amostras de teste

Fonte: elaborado pelo autor, 2023.

O quadro 3.2.1.2, refere-se à verificação de dados faltantes. Essa verificação se faz necessária, pois valores ausentes, podem comprometer o registro.

Quadro 3.2.1.2 – Dados ausentes

Conjunto treino		Conjunto teste	
Coluna	Registro nulo	Coluna	Registro nulo
0	0	0	0
1	0	1	0
2	0	2	0
3	0	3	0
4	0	4	0
...
183	0	183	0
184	0	184	0
185	0	185	0
186	0	186	0
187	0	187	0

Fonte: Elaborado pelo autor, 2023.

Existe um desequilíbrio na quantidade de amostras, de acordo com o tipo de batimento, onde predomina a batida normal, conforme as figuras 3.2.1.2 e 3.2.1.4.

Por esse motivo foi feito o balanceamento, compondo um novo data set com 20 mil amostras de cada tipo de batida, sendo geradas a partir das amostras originais, mantendo as características predominantes, para treinamento do modelo. A figura a seguir demonstra como ficou a distribuição deste conjunto.

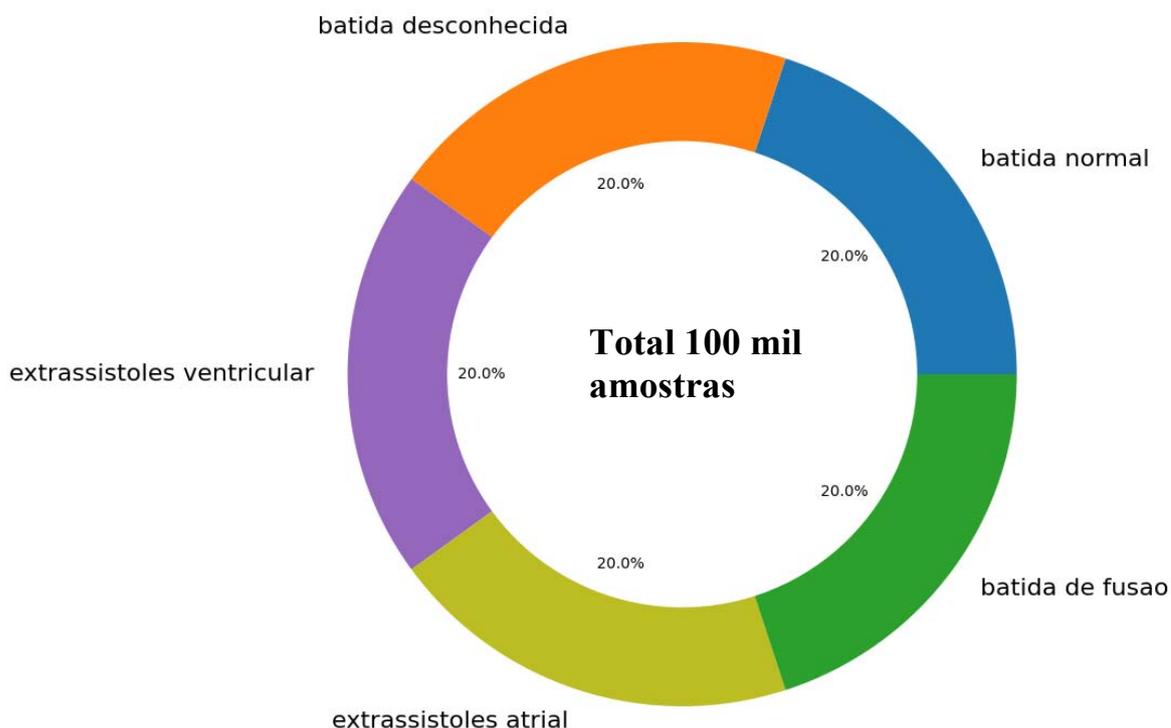


Figura 3.2.1.5 – Distribuição do conjunto de treino balanceado

Fonte: elaborado pelo autor, 2023.

3.3 – Modelagem rede neural Convolutacional (CNN)

As previsões foram realizadas através da análise de imagem dos exames ECG. Como os dados são numéricos, foi necessário convertê-los em imagem.

Para isso, as imagens foram plotadas e carregadas para os modelo.

A figura 3.3.1, representa as imagens obtidas como referência dos tipos de batimentos cardíacos presentes no conjunto de dados.

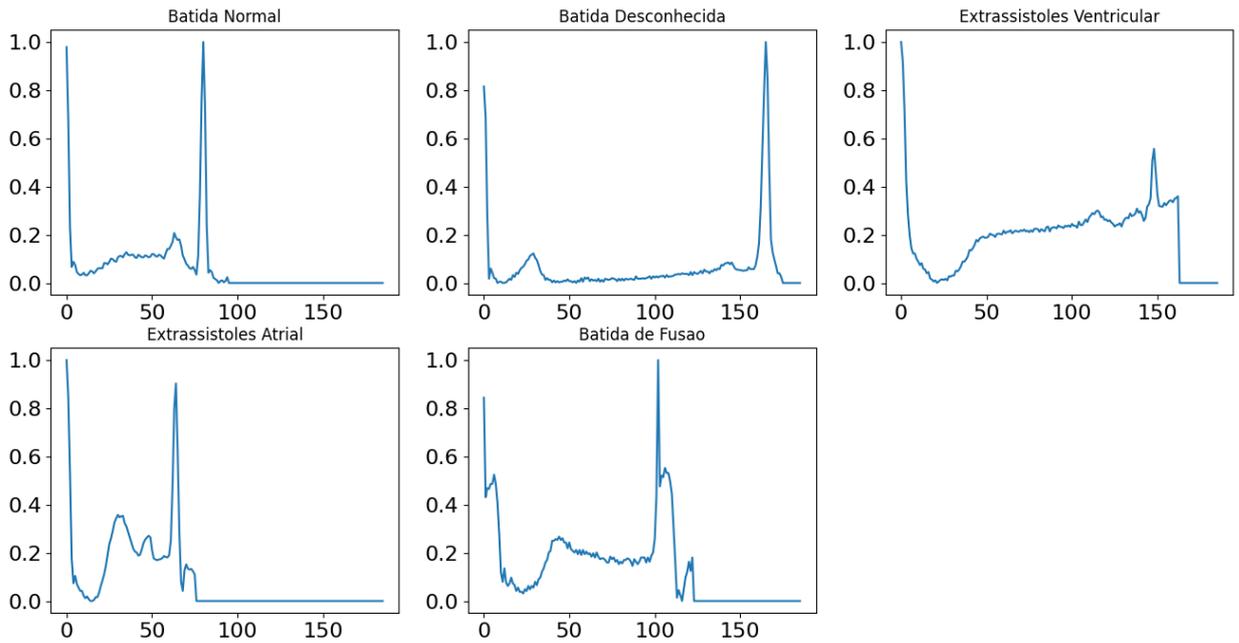


Figura 3.3.1 – Plotagem das batidas cardíacas

Fonte: elaborado pelo autor, 2023.

Buscando a melhor performance, foi testado um filtro de ruídos gaussiano, porém, o resultado não foi satisfatório, sendo este descartado. A figura 3.3.2, representa a atuação do filtro no modelo, onde a imagem inferior é sem tratamento, e a superior é tratada.

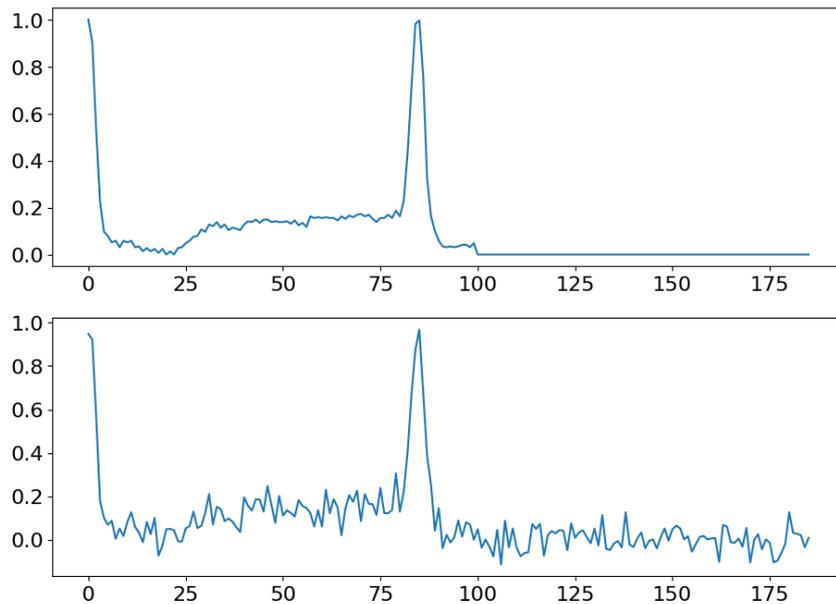


Figura 3.3.2 – Atuação do filtro gaussiano

Fonte: elaborado pelo autor, 2023.

A rede neural, foi treinada com dados desbalanceados e balanceados. As amostras de treino e teste, estão distribuídas em 75% e 25%, respectivamente. A tabela 3.3.1, apresenta os dados da melhor performance obtida, através dos dados originais.

Tabela 3.3.1 – Desempenho treinamento rede neural convolucional

Epoch 1/5
2737/2737 [=====] - 115s 41ms/step - loss: 0.1299 - accuracy: 0.9639 - val_loss: 0.1221 - val_accuracy: 0.9628
Epoch 2/5
2737/2737 [=====] - 110s 40ms/step - loss: 0.0771 - accuracy: 0.9777 - val_loss: 0.0986 - val_accuracy: 0.9727
Epoch 3/5
2737/2737 [=====] - 112s 41ms/step - loss: 0.0591 - accuracy: 0.9824 - val_loss: 0.0791 - val_accuracy: 0.9785
Epoch 4/5
2737/2737 [=====] - 111s 41ms/step - loss: 0.0466 - accuracy: 0.9856 - val_loss: 0.0692 - val_accuracy: 0.9820
Epoch 5/5
2737/2737 [=====] - 118s 43ms/step - loss: 0.0414 - accuracy: 0.9876 - val_loss: 0.0652 - val_accuracy: 0.9826
Accuracy: 98.26%

Fonte: elaborado pelo autor, 2023.

3.4 – Modelagem rede neural

Na modelagem desta rede, foram utilizados os dados na forma numérica, não sendo necessário converter em imagem.

Utilizou-se a biblioteca Sklearn no pré-tratamento dos dados e normalização deles. Também foi utilizado 75% para treino, e 25% para teste.

A figura 3.4.1, demonstra a distribuição total dos dados, consolidados os 02 conjuntos. Os conjuntos de dados assim como no modelo anterior, foram analisados e estão em conformidade para serem utilizados.

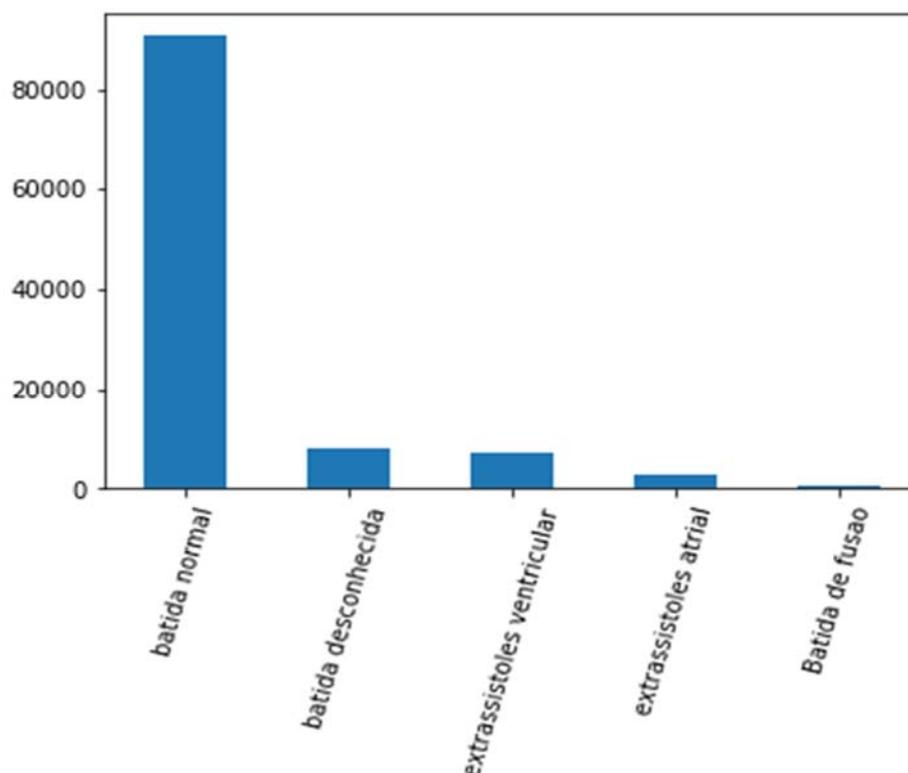


Figura 3.4.1 – Dados de treino e teste consolidados

Fonte: elaborado pelo autor, 2023.

Os dados foram normalizados, e a rede neural foi “parametrizada” com 05 camadas, sendo da primeira à quarta com 50 neurônios, e a camada de saída (softmax) com 05 neurônios. As distribuições dos conjuntos de treino e teste, foram de 75% e 25%, respectivamente.

A quantidade de EPOCHS (épocas), utilizada é 05.

A tabela 3.4.1, apresenta as estatísticas da performance do conjunto de teste.

Tabela 3.4.1 – Desempenho treinamento rede neural

[=====] - 4s 1ms/step - loss: 0.0897 - accuracy: 0.9763				
	precision	recall	f1-score	support
0	0.98	0.99	0.99	90589
1	0.92	0.64	0.76	2779
2	0.93	0.93	0.93	7236
3	0.93	0.54	0.68	803
4	0.98	0.98	0.98	8039
accuracy	0.98			109446
macro avg	0.95	0.82	0.87	109446
weighted avg	0.98	0.98	0.97	109446

Fonte: elaborado pelo autor, 2023.

CAPÍTULO 4

Resultados Obtidos ou Esperados

4.1 – Vantagens da aplicação da IA na área da saúde

O SUS atende aproximadamente 150 milhões de brasileiros. Se todos esses usuários do sistema tivessem um prontuário eletrônico, teríamos a possibilidade de fazer muitas campanhas, selecionando de maneira mais assertiva as pessoas a serem atendidas. Muitas vidas seriam poupadas, e muito dinheiro economizado.

A utilização de inteligência artificial na previsão de doenças já é bastante utilizada nos tempos atuais, com importantes centros de pesquisas atuando nessa direção.

Conforme estimado por pesquisadores da UNICAMP, a utilização desses recursos proporciona economia aos cofres públicos, pois a prevenção é mais benéfica e econômica, do que um tratamento tardio.

4.2 – Análise da base de dados

Os dados já se encontram divididos em conjunto de treino e teste, e não possui registros duplicados e nem dados faltantes. Os exames foram elaborados por importantes centros clínicos e laudados por médicos, respeitando a privacidade dos pacientes, de acordo com a LGPD.

A quantidade de registros é suficiente para treinar e testar redes neurais com boa confiabilidade. Todos os registros são numéricos, não possuindo variáveis qualitativas, conforme exposto na figura 3.2.1.

4.3 – Treinamento da rede neural

Esta rede neural foi testada com diferentes EPOCHS, e o melhor resultado foi 05. Os dados foram distribuídos em 75% no conjunto de treino, e 25% de teste. Também foi testado a base balanceada, porém a performance com a utilização de dados sintéticos não foi boa,

descartando este. A aplicação do filtro gaussiano para diminuir o ruído também piorou o resultado, sendo eliminado também. A função de ativação com melhor resultado é a RELU.

Com a função de ativação SIGMOIDE, a melhor acurácia não passou de 94%.

A tabela a 4.3.1, apresenta alguns dados da parametrização do modelo.

Tabela 4.3.1 – Desempenho conjunto de teste

Model: "sequential"
Layer (type) Output Shape Param #
dense (Dense) (None, 50) 9400
dense_1 (Dense) (None, 50) 2550
dense_2 (Dense) (None, 50) 2550
dense_3 (Dense) (None, 50) 2550
dense_4 (Dense) (None, 5) 255
Total params: 17,305
Trainable params: 17,305
Non-trainable params: 0
accuracy 0.98

Fonte: elaborado pelo autor, 2023.

4.4 – Treinamento da rede neural convolucional

No treinamento da CNN, os dados foram distribuídos em 75% no conjunto de treino, e 25% de teste. Também foi testado a base balanceada, porém a performance com a utilização de dados sintéticos não foi boa, descartando este, conforme também ocorreu com o modelo anterior. A aplicação do filtro gaussiano para diminuir o ruído também piorou o resultado, sendo eliminado também. A função de ativação com melhor resultado é a RELU. A quantidade de EPOCHS que melhor generalizou o modelo foi 05.

Este foi o modelo com melhor performance, atingido 98.26% de precisão, utilizando os dados originais, conforme apresentado na tabela 4.4.1.

Tabela 4.4.1 – Performance conjunto de teste

CNN	CNN
RELU	RELU
EPOCHS=5	EPOCHS=5
Base original	Base Balanceada
98.26%	97.40%

Fonte: elaborado pelo autor, 2023.

A figura 4.4.1, representa as métricas de desempenho obtido no treinamento do modelo.

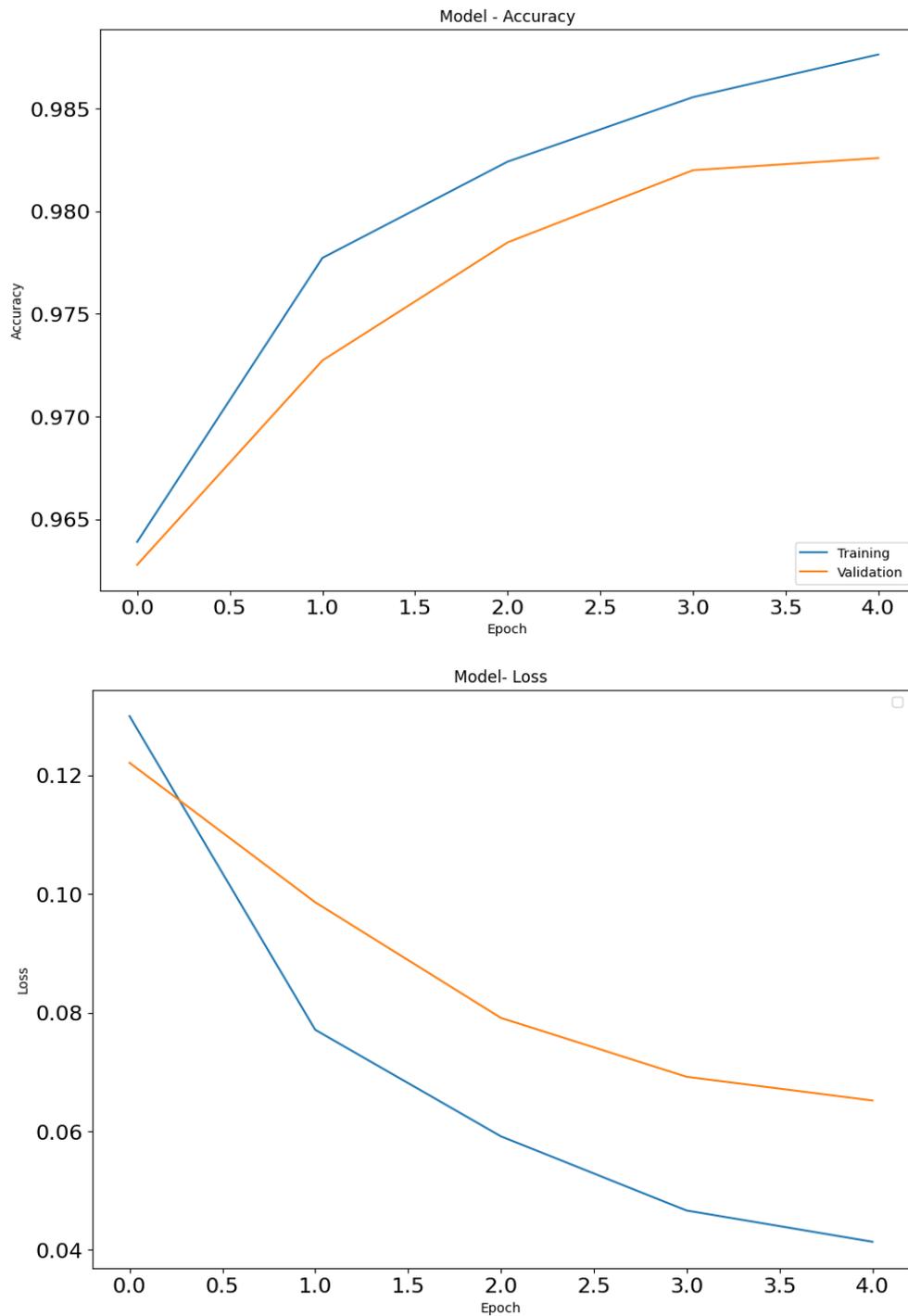


Figura 4.4.1 – Desempenho de treinamento CNN

Fonte: elaborado pelo autor, 2023.

A seguir, é apresentado na figura 4.4.2, os gráficos de performance deste modelo, que obteve a segunda melhor performance, com acurácia de 96.31%. Na curva de validação, é

possível notar uma queda de performance, evidenciado na variação negativa da acurácia, que acabou sendo corrigido pelo modelo.

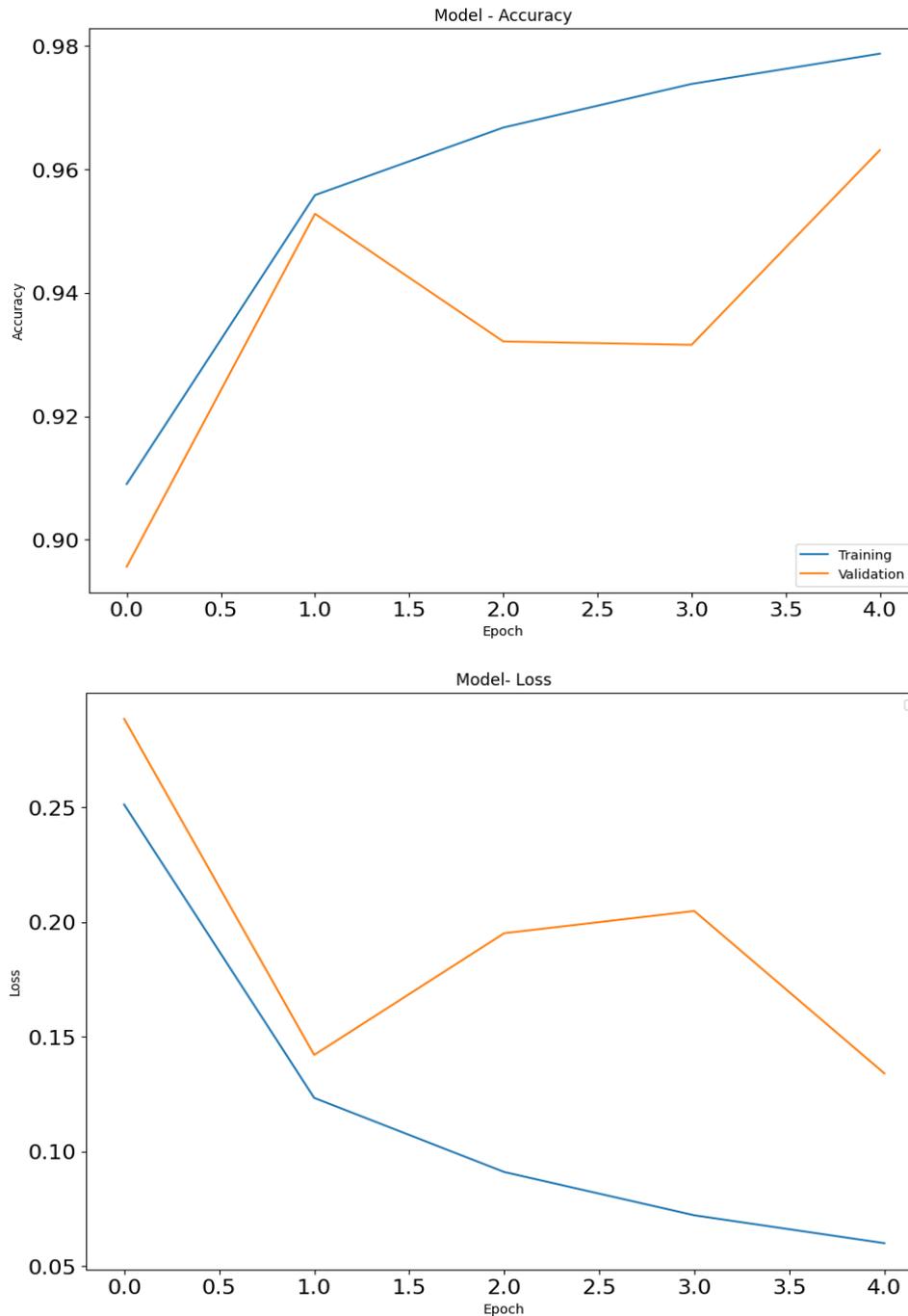


Figura 4.4.2 – Desempenho de treinamento CNN

Fonte: elaborado pelo autor, 2023.

No treinamento da CNN, os dados foram distribuídos em 75% no conjunto de treino, e 25% de teste. Também foi testado a base balanceada, porém a performance com a utilização de dados sintéticos não foi boa, descartando este, conforme também ocorreu com o modelo

anterior. A aplicação do filtro gaussiano para diminuir o ruído também piorou o resultado, sendo eliminado também. A função de ativação com melhor resultado é a RELU. A quantidade de EPOCHS que melhor generalizou o modelo foi 05.

Este foi o modelo com melhor performance, atingido 98.26% de precisão, utilizando os dados originais, conforme a tabela a seguir.

4.5 – Análise de erros na previsão

Foi criado e analisado um data frame com as previsões erradas. O percentual dos tipos de batidas neste conjunto é apresentado na figura 4.5.1.

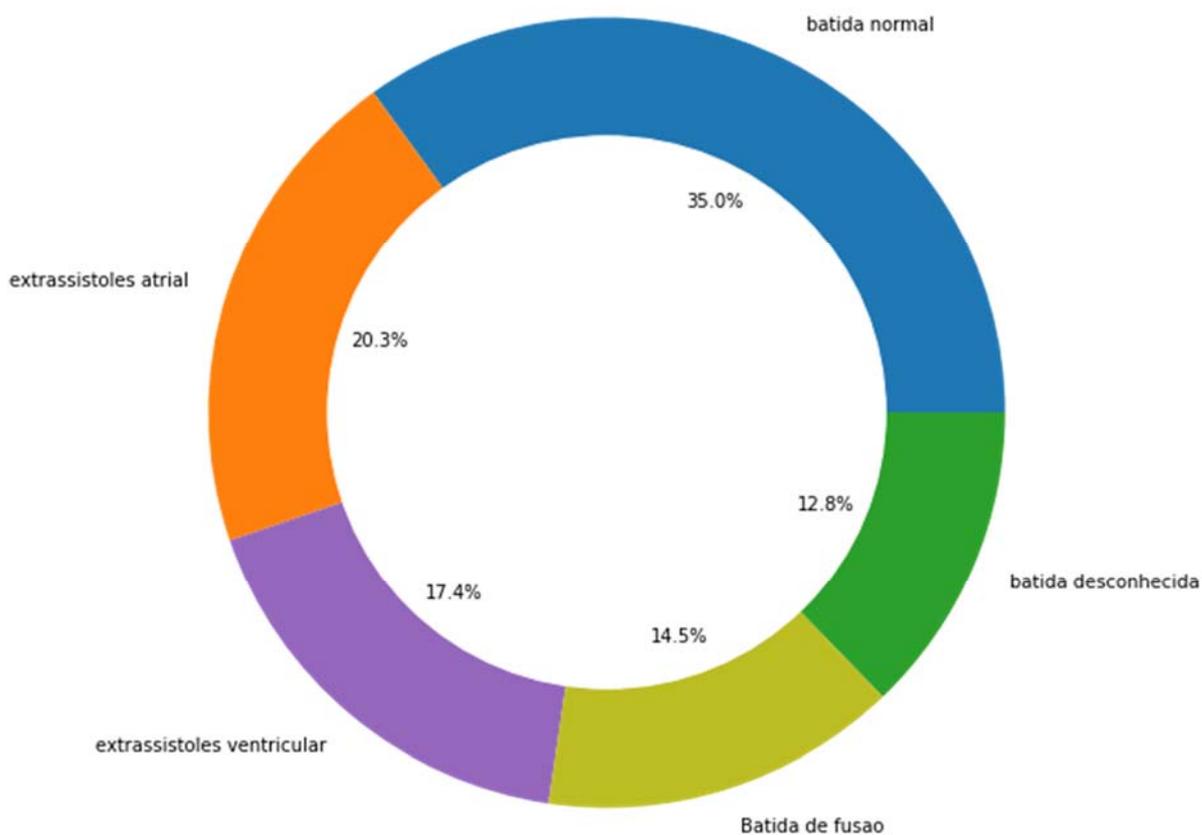


Figura 4.5.1 – Previsões erradas

Fonte: elaborado pelo autor, 2023.

A tabela 4.5.1, apresenta as distribuições do conjunto de treino e teste, conjunto de erro e acurácia distribuída por categoria de batida cardíaca.

Tabela 4.5.1 – Previsões erradas

	Tipos de Batida				
	normal	desconhecida	extrassístoles atrial	extrassístoles ventricular	Fusão
conjunto treino / teste	82.8%	7.3%	2.5%	6.6%	0.7%
conjunto de erro	35%	12.8%	20.3%	17.4%	14.5%
percentual	0.7%	0.26%	0.41%	0.35%	0.29%
Acuracia total	1.13%	6.28%	43.7%	8.35%	54.93%

Fonte: elaborado pelo autor, 2023.

Como pode ser observado na tabela acima, as acurácias precisam ser observadas por segmento de classificação, pois embora a acurácia final seja 98%, o desequilíbrio no conjunto de testes faz com que os segmentos com menos amostras tenham baixa acurácia, e isso seja mascarado no conjunto total, considerando todas as médias.

CAPÍTULO 5

Conclusão e Trabalhos Futuros

5.1 – Conclusão

A utilização de IA na área da saúde é realidade e muitos Laboratórios ligados a importantes institutos vem desenvolvendo modelos para a previsão de doenças, como o LABBDAPS, ligado a faculdade de saúde da USP, e o ATEROLAB, ligado a Faculdade de ciências médicas da UNICAMP.

Neste último, é citado em um artigo que em 2017 obtiveram 88% de precisão na análise gráfica de ECG com machine learning. Neste trabalho, foi obtido 98% de precisão total no mesmo tipo de análise. Porém, precisa avaliar os resultados segmentados por categoria.

Na análise de desempenho por categoria, fica evidente que a batida por fusão e extrassístoles atrial, o modelo não faz previsão confiável, muito evidentemente isso está relacionado a baixa quantidade de amostras no conjunto de treino, conforme apresentado na tabela 4.5.1.

Se o modelo for considerado para avaliar apenas batida normal ou não, a precisão é excelente, sendo necessário rever a avaliação caso classifique a primeira situação citada. A análise gráfica, através de CNN apresentou melhor performance, porém a rede neural utilizando tensorflow, numa rede neural convencional obteve um bom resultado.

Este trabalho é continuação da monografia defendida em julho/2023, da turma MB3B, nesta instituição. Nele foi avaliado a utilização de *machine learning*, utilizando KNN e Regressão Logística na previsão de doenças cardíacas e diabetes.

As performances não apresentaram números expressivos, sendo que na previsão de doenças do coração foram utilizadas variáveis categóricas, relacionadas a comportamentos de risco, e apresentou a pior performance, com acurácia de 71%. Para a previsão de diabetes, a acurácia foi de 84%, e neste, o conjunto de dados utilizado possui variáveis quantitativas, relacionados a métricas (exames) com maior correlação com a doença, como exemplo a glicose clicada. Nesse estudo também se observou, que a precisão oscila de acordo com o índice de probabilidades.

Os resultados obtidos demonstram o quanto a aplicação dessas tecnologias pode ser precisa, e um importante instrumento de apoio médico.

Certamente o futuro da medicina está no *BIG DATA* e IA, ferramentas essas que funcionam como vetores de aceleração em pesquisas, e apoio médico, e que irão tornar a medicina mais inclusiva e segura.

5.2 – Trabalhos Futuros

Como trabalho futuro, será analisado outros modelos de machine learning, aplicado em conjuntos de exames médicos, para maior aprofundamento neste campo.

Referências Bibliográficas

AMAZON AWS. **O que é uma rede neural?** Disponível em: <https://aws.amazon.com/pt/what-is/neural-network/>. Acesso em 22 jul 2022.

AMAZON WEB SERVICES. **Documentação do Amazon Machine Learning.** Disponível em: <https://docs.aws.amazon.com/machine-learning>. Acesso em: 22 jul. 2022.

BRASIL. Portal da Transparência. **Função 10 - Saúde. 2022.** Disponível em: <https://www.portaltransparencia.gov.br/funcoes/10-saude?ano=2022>. Acesso em: 18 jul. 2022.

CHIAVEGATTO FILHO, Alexandre Dias Porto. **Uso de big data em saúde no Brasil: perspectivas para um futuro próximo.** Epidemiol. Serv. Saúde, Brasília, v. 24, n. 2, p. 325-332, jun. 2015. Disponível em http://scielo.iec.gov.br/scielo.php?script=sci_arttext&pid=S1679-49742015000200015&lng=pt&nrm=iso>49742015000200015&lng=pt&nrm=iso>. Acessos em 23 jul. 2022.

DATASUS. **Sobre o Datasus.** Disponível em: <https://datasus.saude.gov.br/sobre-o-datasus/>. Acesso em: 20 jul. 2022.

DEEP LEARNING BOOK. **Neurônio biológico e matemático.** Disponível em: <https://www.deeplearningbook.com.br/o-neuronio-biologico-e-matematico/>. Acesso em 22 jul 2022.

DEEPLARNING BOOK. **Introdução as redes convolucionais.** Disponível em: <https://www.deeplearningbook.com.br/introducao-as-redes-neurais-convolucionais>. Acesso em 29 jul 2023

DIDATICA TECK. **O que são Redes Neurais e Deep Learning?** Disponível em: <https://didatica.tech/introducao-a-redes-neurais-e-deep-learning/>. Acesso em 22 jul 2023.

ECGNOW. **5 aspectos que você precisa saber sobre extrassístoles.** Disponível em: <https://www.ecgnow.com.br/blog/5-aspectos-que-voce-precisa-saber-sobre-extrassistoles/>. Acesso em 28 jul 2023.

FAPESP. **Inteligência artificial e doenças do coração.** Disponível em: https://revistapesquisa.fapesp.br/wp-content/uploads/2020/08/070-075_ia-e-doen%C3%A7as-cora%C3%A7%C3%A3o_294.pdf. Acesso em: 24 jul. 2022.

FIOCRUZ. **Big Data em Saúde.** Disponível em: <https://www.iciet.fiocruz.br/content/big-data-em-saude>. Acesso em: 24 jul. 2022.

IBGE. **Pesquisa Nacional de Saúde**. Disponível em: <https://www.ibge.gov.br/estatisticas/sociais/saude/9160-pesquisa-nacional-de-saude.html>. Acesso em: 16 jul. 2022.

IBM. **Aprendizado de Máquina**. Disponível em: <https://www.ibm.com/br-pt/cloud/learn/machine-learning>. Acesso em: 22 jul. 2022.

INSTITUTE FOR HEALTH METRICS AND EVALUATION. **Global Burden of Disease Study**. Disponível em: <http://www.healthdata.org/gbd>. Acesso em: 18 jul. 2022.

KAGGLE. **Ecg heartbeat categorization dataset**. Disponível em: <https://www.kaggle.com/datasets/shayanfazeli/heartbeat>. Acessado em 26 jul 2023.

MANUAL MSD. **Arritmias supraventriculares ectópicas**. Disponível em: <https://www.msmanuals.com/pt-br/profissional/doen%C3%A7as-cardiovasculares/arritmias-card%C3%ADacas-espec%C3%ADficas/arritmias-supraventriculares-ect%C3%B3picas>. Acesso em 20 jul 2023

MATPLOTLIB. Disponível em: <https://matplotlib.org/>. Acesso em: 28 mar. 2023.

MEDICINA MITOS E VERDADES. **Arritmia cardíaca e morte súbita**. Disponível em: <https://www.medicinamitoseverdades.com.br/blog/arritmia-cardiaca-e-morte-subita>. Acesso em 22 jul 2023.

MEDIUM. **Redes neurais convencionais**. Disponível em: <https://medium.com/itau-data/redes-neurais-convolucionais-2206a089c715>. Acesso em: 29 jul. 2023.

MESQUITA, Claudio Tinoco. **Inteligência Artificial e Machine Learning em Cardiologia – Uma Mudança de paradigma**. Int J. Cardiovasc Sci, v. 30, n. 3, p. 187-188, maio. 2017.

MY EKG. Disponível em: <https://pt.my-ekg.com/arritmias-cardiacas/extrassistolares-ventriculares.html>. Acessado em: 27 jul 2023.

NUMPY. Disponível em: <https://numpy.org/>. Acesso em: 05 ago. 2022.

ORACLE. **Big Data**. Disponível em: <https://www.oracle.com/br/big-data/>. Acesso em: 23 jul. 2022.

ORACLE. **Inteligência Artificial**. Disponível em: <https://www.oracle.com/br/artificial-intelligence/>. Acesso em: 20 jul. 2022.

PAHO. **Doenças Cardiovasculares**. Disponível em: <https://www.paho.org/pt/topicos/doencas-cardiovasculares>. Acesso em: 19 jul. 2022.

PANDAS. Disponível em: <https://pandas.pydata.org/>. Acesso em: 05 ago. 2022.

PYTHON. **Documentação técnica.** Disponível em: <https://www.Python.org/>. Acesso em: 05 ago. 2022.

RBEM. **Nova Metodologia de ensino do ECG:** Desmitificando a Teoria na Prática – Ensino Prático do ECG. Disponível em: <https://www.scielo.br/j/rbem/a/RXbsLmvxHH9jG7H9NFWRdJB/?lang=pt>. Acesso em 15 jul 2023

REIS, Helder José Lima, et al. **ECG manual prático de eletrocardiograma.** In: ECG manual prático de eletrocardiograma. 2013. p. 121-121.

SBC. **Seção de eletrocardiograma.** Disponível em: http://educacao.cardiol.br/2014/ecg/exibir_anterior.asp?cod=53. Acesso em: 28 jul 2023.

SCIKIT-LEARN. Disponível em: <https://scikit-learn.org/stable/index.html>. Acesso em: 19 abr. 2023.

SEABORN. Disponível em: <https://seaborn.pydata.org/>. Acesso em: 05 ago. 2022.