



UNIVERSIDADE FEDERAL DE GOIÁS
INSTITUTO DE INFORMÁTICA
PROGRAMA DE PÓS GRADUAÇÃO EM CIÊNCIA DA
COMPUTAÇÃO

LÍVIA MANCINE COELHO DE CAMPOS

**Mineração de dados de autópsia para
determinar as causas de morte na
depressão**

Goiânia
2021



UNIVERSIDADE FEDERAL DE GOIÁS
INSTITUTO DE INFORMÁTICA

TERMO DE CIÊNCIA E DE AUTORIZAÇÃO (TECA) PARA DISPONIBILIZAR VERSÕES ELETRÔNICAS DE TESES

E DISSERTAÇÕES NA BIBLIOTECA DIGITAL DA UFG

Na qualidade de titular dos direitos de autor, autorizo a Universidade Federal de Goiás (UFG) a disponibilizar, gratuitamente, por meio da Biblioteca Digital de Teses e Dissertações (BDTD/UFG), regulamentada pela Resolução CEPEC nº 832/2007, sem ressarcimento dos direitos autorais, de acordo com a [Lei 9.610/98](#), o documento conforme permissões assinaladas abaixo, para fins de leitura, impressão e/ou download, a título de divulgação da produção científica brasileira, a partir desta data.

O conteúdo das Teses e Dissertações disponibilizado na BDTD/UFG é de responsabilidade exclusiva do autor. Ao encaminhar o produto final, o autor(a) e o(a) orientador(a) firmam o compromisso de que o trabalho não contém nenhuma violação de quaisquer direitos autorais ou outro direito de terceiros.

1. Identificação do material bibliográfico

Dissertação Tese

2. Nome completo do autor

Livia Mancine Coelho de Campos

3. Título do trabalho

Mineração de dados de autopsia para determinar as causas de morte na depressão

4. Informações de acesso ao documento (este campo deve ser preenchido pelo orientador)

Concorda com a liberação total do documento SIM NÃO¹

[1] Neste caso o documento será embargado por até um ano a partir da data de defesa. Após esse período, a possível disponibilização ocorrerá apenas mediante:

- a) consulta ao(a) autor(a) e ao(a) orientador(a);
 - b) novo Termo de Ciência e de Autorização (TECA) assinado e inserido no arquivo da tese ou dissertação.
- O documento não será disponibilizado durante o período de embargo.

Casos de embargo:

- Solicitação de registro de patente;
- Submissão de artigo em revista científica;
- Publicação como capítulo de livro;
- Publicação da dissertação/tese em livro.

Obs. Este termo deverá ser assinado no SEI pelo orientador e pelo autor.



Documento assinado eletronicamente por **Rogério Lopes Salvini, Professor do Magistério Superior**, em 29/10/2021, às 10:37, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do [Decreto nº 10.543, de 13 de novembro de 2020](#).

Documento assinado eletronicamente por **LÍVIA MANCINE COELHO DE CAMPOS, Discente**, em 29/10/2021, às 11:27, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do



[Decreto nº 10.543, de 13 de novembro de 2020.](#)



A autenticidade deste documento pode ser conferida no site https://sei.ufg.br/sei/controlador_externo.php?acao=documento_conferir&id_orgao_acesso_externo=0, informando o código verificador **2454375** e o código CRC **DD4E74B3**.

LÍVIA MANCINE COELHO DE CAMPOS

Mineração de dados de autopsia para determinar as causas de morte na depressão

Dissertação apresentada ao Programa de Pós Graduação em Ciência da Computação do Instituto de Informática da Universidade Federal de Goiás (UFG), como requisito para obtenção do título de Mestre em Ciência da Computação.

Área de concentração: Ciência da Computação.

Linha de pesquisa: Sistemas Inteligentes e Aplicações.

Orientador: Prof. Dr. Rogerio Lopes Salvini

Goiânia
2021

Ficha de identificação da obra elaborada pelo autor, através do Programa de Geração Automática do Sistema de Bibliotecas da UFG.

Campos, Livia Mancine Coelho de
Mineração de dados de autópsia para determinar as causas de morte na depressão [manuscrito] / Livia Mancine Coelho de Campos. 2021.
XCV, 95 f.

Orientador: Prof. Dr. Rogerio Lopes Salvini.
Dissertação (Mestrado) - Universidade Federal de Goiás, Instituto de Informática (INF), Programa de Pós-Graduação em Ciência da Computação, Goiânia, 2021.
Bibliografia. Apêndice.
Inclui tabelas, lista de figuras, lista de tabelas.

1. depressão. 2. causa de morte. 3. autópsia. 4. mineração de dados. 5. aprendizado de máquina. I. Salvini, Rogerio Lopes, orient. II. Título.

CDU 004



UNIVERSIDADE FEDERAL DE GOIÁS

INSTITUTO DE INFORMÁTICA

ATA DE DEFESA DE DISSERTAÇÃO

Ata nº 24 da sessão de Defesa de Dissertação de **Livia Mancine Coelho de Campos**, que confere o título de Mestre em Ciência da Computação, na área de concentração em Ciência da Computação.

Aos oito dias do mês de outubro de dois mil e vinte e um, a partir das nove horas, via sistema de webconferência da RNP, realizou-se a sessão pública de Defesa de Dissertação intitulada “**Mineração de dados de autopsia para determinar as causas de morte na depressão**”. Os trabalhos foram instalados pelo Orientador, Professor Doutor Rogerio Lopes Salvini (INF/UFG) com a participação dos demais membros da Banca Examinadora: Professora Doutora Nilza Nascimento Guimarães (ICB/UFG), membra titular externa; Professor Doutor Fabrizzio Alphonsus Alves de Melo Nunes Soares (INF-UFG), membro titular interno. A realização da banca ocorreu por meio de videoconferência, em atendimento à recomendação de suspensão das atividades presenciais na UFG emitida pelo Comitê UFG para o Gerenciamento da Crise COVID-19, bem como à recomendação de isolamento social da Organização Mundial de Saúde e do Ministério da Saúde para enfrentamento da emergência de saúde pública decorrente do novo coronavírus. Durante a arguição os membros da banca não fizeram sugestão de alteração do título do trabalho. A Banca Examinadora reuniu-se em sessão secreta a fim de concluir o julgamento da Dissertação, tendo sido a candidata **aprovada** pelos seus membros. Proclamados os resultados pelo Professor Doutor Rogerio Lopes Salvini, Presidente da Banca Examinadora, foram encerrados os trabalhos e, para constar, lavrou-se a presente ata que é assinada pelos Membros da Banca Examinadora, aos oito dias do mês de outubro de dois mil e vinte e um.

TÍTULO SUGERIDO PELA BANCA



Documento assinado eletronicamente por **Rogerio Lopes Salvini, Professor do Magistério Superior**, em 08/10/2021, às 13:04, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do [Decreto nº 10.543, de 13 de novembro de 2020](#).



Documento assinado eletronicamente por **Fabrizzio Alphonsus Alves De Melo Nunes Soares, Professora do Magistério Superior**, em 08/10/2021, às 13:05, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do [Decreto nº 10.543, de 13 de novembro de 2020](#).



Documento assinado eletronicamente por **LÍVIA MANCINE COELHO DE CAMPOS, Discente**, em 08/10/2021, às 13:06, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do [Decreto nº 10.543, de 13 de novembro de 2020](#).



Documento assinado eletronicamente por **Nilza Nascimento Guimarães, Professor do Magistério Superior**, em 08/10/2021, às 14:27, conforme horário oficial de Brasília, com fundamento no § 3º do art. 4º do [Decreto nº 10.543, de 13 de novembro de 2020](#).

A autenticidade deste documento pode ser conferida no site
https://sei.ufg.br/sei/controlador_externo.php?



[acao=documento_conferir&id_orgao_acesso_externo=0](#), informando o código verificador **2361478** e o código CRC **30017812**.

Referência: Processo nº 23070.047800/2021-07

SEI nº 2361478

Dedico este trabalho à minha família, Leonardo, Júlia e Ester, pelo apoio incondicional em toda minha trajetória de estudos durante o mestrado.

Agradecimentos

A Deus pelo sustento em todos os momentos da minha vida.

Ao meu esposo e minhas filhas, vocês foram tão importantes que sem vocês eu não chegaria até aqui.

À minha rede de apoio, minha família e meus amigos e amigas, que me ajudaram nos momentos em que não pude estar presente com as minhas filhas.

Ao Prof. Rogerio Salvini, por todo aprendizado e paciência, pela oportunidade de trabalharmos juntos e por me orientar. Grata por sua excelência!

À Prof. Paula, que foi colaborada deste projeto. Obrigada por sua dedicação e por todo aprendizado. Você é incrível!

Aos Professores do INF, Prof. Dr. Fabrízio, Prof. Dr. Thierson, Prof. Dr. Flávio, Prof. Dr. Hugo que contribuíram para minha vida profissional. Em especial a Prof. Dra. Márcia Capelle por ser tão inspiradora!

Aos amigos Cleon, Marcos e Bruna, vocês foram essenciais! Muito obrigada pelo apoio (de sempre)!

Aos meus amigos do mestrado Pollyana, Divino, Naiane, João Lucas, Luíla e Jurandir, a caminhada com vocês foi mais serena. Obrigada pelo apoio e pelas conversas. Vocês são para toda a vida!

Resumo

Mancine, Lívia. **Mineração de dados de autópsia para determinar as causas de morte na depressão**. Goiânia, 2021. 95p. Dissertação de Mestrado. Instituto de Informática, Universidade Federal de Goiás.

A depressão está associada ao aumento da mortalidade, mas as causas de morte ligadas à depressão, bem como a expectativa de vida dos pacientes, ainda são pouco exploradas e ainda há controvérsias. Identificar possíveis doenças associadas à morte em pacientes com depressão pode auxiliar na tomada de decisão das políticas públicas de saúde e culminar em tratamentos mais específicos, estratégias de prevenção e em uma melhor expectativa de vida desses pacientes. Nos estudos sobre causas de morte onde exames de autópsia são analisados, é possível adquirir informações mais precisas sobre as doenças relacionadas à morte, uma vez que a autópsia determina a exatidão da causa da morte. Neste estudo, 1.136 indivíduos foram avaliados segundo laudos de autópsia provenientes do Serviço de Verificação de Óbitos da Capital (SVOC-USP) da região metropolitana de São Paulo, e a Entrevista Clínica Estruturada para o DSM-IV para o diagnóstico de depressão.

Foi realizada a mineração de dados baseada nas CIDs das causas relacionadas à morte, no qual onze algoritmos de Aprendizado de Máquina foram aplicados a fim de se buscar padrões para determinar as possíveis causas de morte relacionadas à depressão. Além da depressão maior, outros oito subgrupos de depressão foram analisados. Embora este estudo tenha feito uma ampla investigação na população em geral e em grupos específicos de pacientes, os resultados obtidos pelos modelos gerados não indicam diferenças de padrões nas causas de morte em indivíduos com e sem depressão. Este resultado corrobora com estudos anteriores da literatura onde as evidências das causas de morte por todas as causas e causas específicas e depressão não são significativas.

Palavras-chave

depressão, causa de morte, autópsia, mineração de dados, aprendizado de máquina

Abstract

Mancine, Lívia. **Full-body autopsy data mining to determine causes of death in major depression**. Goiânia, 2021. 95p. MSc. Dissertation. Instituto de Informática, Universidade Federal de Goiás.

Depression is associated with increased mortality, but the causes of death related to depression, as well as the life expectancy of patients, are still poorly explored and controversial. Identifying possible diseases associated with death in patients with depression can help in public health policy decision making and culminate in more specific treatments, prevention strategies, and improved life expectancy for these patients. In studies on causes of death which autopsy examinations are analyzed, it is possible to acquire more accurate information about the diseases related to death, since the autopsy determines the precise cause of death. In this study, we evaluated the causes of death of 1,136 subjects, according to autopsy reports from the Death Verification Service of the Capital (SVOC-USP) in the metropolitan region of São Paulo. The diagnosis of depression of these subjects was made according to the Structured Clinical Interview for DSM-IV (SCID). Data mining based on the ICDs of causes related to death was performed, in which eleven Machine Learning algorithms were applied in order to search for patterns to determine the possible causes of death related to depression. In addition to major depression, eight other subgroups of depression were analyzed. Although this study performed a broad investigation in the general population and in specific groups of patients, the results obtained by the generated models do not indicate differences in patterns in the causes of death in individuals with and without depression. This result corroborates with previous studies in the literature where the evidences for all-cause and cause-specific causes of death and depression are not significant.

Keywords

major depression, cause of death, autopsy, data mining, machine learning

Sumário

Lista de Figuras	13
Lista de Tabelas	14
1 Introdução	16
1.1 Contextualização	16
1.2 Justificativa	18
1.3 Motivação	18
1.4 Problema de pesquisa	19
1.5 Objetivos	19
1.6 Estrutura do Documento	20
2 Referencial Teórico	21
2.1 Classificação Internacional de Doenças (CID)	21
2.2 Declaração de Óbito	23
2.3 Autópsia	25
2.4 Mineração de Dados	27
2.5 Psiquiatria Computacional	31
2.6 Considerações Finais	32
3 Mapeamento Sistemático	33
3.1 Planejamento	33
3.2 Condução do processo de seleção dos estudos primários	36
3.3 Análise do Resultado sobre às Questões de Pesquisa	37
3.4 Considerações Finais	43
4 Materiais e Métodos	45
4.1 Conjuntos de dados	45
4.2 Método	48
4.2.1 Preparação do conjunto de dados	49
4.2.2 Mineração de Dados	58
4.2.3 Avaliação	58
4.3 Ferramentas	59
5 Resultados e Discussão	61
5.1 Experimentos	61
5.2 Causas de morte na depressão	62
5.3 Causas de morte em subgrupos da depressão	63
5.3.1 Depressão tardia	64

5.3.2	Depressão recorrente	65
5.3.3	Depressão primária	66
5.3.4	Depressão secundária	67
5.3.5	Depressão grave	68
5.3.6	Depressão por álcool e drogas	69
5.3.7	Depressão NPI_1	70
5.3.8	Depressão NPI 2	71
5.4	Discussão	72
6	Conclusões	75
6.1	Contribuições	75
6.2	Limitações	76
6.3	Trabalhos Futuros	77
	Referências Bibliográficas	78
A	Descrição dos dados de autópsia e dos dados da entrevista clínica	85
B	Algoritmos de Aprendizado de Máquina	91

Lista de Figuras

2.1	Bloco 5 da Declaração de Óbito com os respectivos campos: Óbitos de Mulher em Idade Fértil, Assistência Médica, Diagnóstico Confirmado por, Causas da Morte [42].	24
2.2	Encaminhamento dos óbitos no Brasil de acordo com a motivação	26
2.3	Aprendizado supervisionado	29
2.4	Aprendizado não supervisionado	29
2.5	Aprendizado semi supervisionado	30
2.6	Aprendizado por reforço	31
4.1	Formulário aplicado ao informante do falecido para a coleta de dados para o diagnóstico de episódio depressivo no passado	48
4.2	Metodologia referente ao processo de descoberta de conhecimento utilizada neste estudo.	49
4.3	Operações realizadas para o tratamento da qualidade da base de dados	51
4.4	Fluxograma do algoritmo de pareamento	53
5.1	Diagrama dos experimentos. Cada experimento utiliza um conjunto diferente de atributos como entrada para os 11 algoritmos de AM selecionados.	62
5.2	Subgrupos de depressão analisados.	64

Lista de Tabelas

2.1	Capítulos da CID-10 com os grupos correspondentes. Como exemplo, pode-se perceber que os grupos de doenças entre A00 a B99 pertencem a Categoria I Fonte: http://tabnet.datasus.gov.br/	22
3.1	Palavras-chave e seus sinônimos de acordo com parte da estrutura PICOC.	35
3.2	Resultado da questão de pesquisa referente ao diagnóstico de depressão de cada estudo selecionado.	38
3.3	Resultado da questão de pesquisa referente a média de idade da amostra estudada.	39
3.4	Resultado da questão de pesquisa referente ao tipo de pesquisa relacionada ao tempo.	40
4.1	Exemplos de casos do conjunto de dados de autópsia	47
4.2	Número de observações e de variáveis nos conjuntos de dados originais	49
4.3	Resumo quantitativo dos dados deste estudo após a integração do conjunto de dados e ao tratamento da qualidade dos dados.	51
4.4	Oito grupos da CID utilizados para análise da causa de morte na depressão. *Outras causas/Doenças infecciosas agrupa doenças como: <i>tuberculosis, sepsis, Chagas' disease, unknown cause of mortality, other disorders of brain</i> , que correspondem respectivamente aos CID's A15, A41, B57, R99, G93.	54
4.5	27 grupos da CID, que agrupam doenças similares, utilizados para análise da causa de morte na depressão.	55
4.6	Exemplo de um registro de autópsia de um indivíduo	56
4.7	Exemplo de um registro com os novos atributos preenchidos a partir dos CID's das causas do óbito.	57
5.1	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do grupo "Depressão". O número entre parênteses indica o desvio padrão.	63
5.2	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão tardia. O número entre parênteses indica o desvio padrão.	65
5.3	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão recorrente. O número entre parênteses indica o desvio padrão.	66
5.4	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão primária. O número entre parênteses indica o desvio padrão.	67

5.5	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão secundária. O número entre parênteses indica o desvio padrão.	68
5.6	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão grave. O número entre parênteses indica o desvio padrão.	69
5.7	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão por álcool e droga. O número entre parênteses indica o desvio padrão.	70
5.8	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão NPI_1. O número entre parênteses indica o desvio padrão.	71
5.9	Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão NPI_2. O número entre parênteses indica o desvio padrão.	72
A.1	Descrição dos dados do laudo de autópsia.	85
A.2	Descrição dos dados da entrevista clínica.	86

Introdução

Este capítulo apresenta os elementos essenciais para a definição do objeto de pesquisa, a saber: a Seção 1.1 contextualiza o assunto da pesquisa e o problema de pesquisa; a Seção 1.2 apresenta a justificativa da pesquisa; a Seção 1.3 aborda a motivação; a Seção 1.5 expõe o objetivo geral e os objetivos específicos. Por fim, a Seção 1.6 compreende a estrutura deste documento.

1.1 Contextualização

As causas de morte ligadas à depressão, bem como a expectativa de vida dos pacientes, ainda são pouco exploradas e ainda há controvérsias [15, 43]. No estudo de Damián et al. [15] foi percebida menor mortalidade da amostra com depressão em relação à amostra controle, provavelmente por conta de sua limitação da amostra, que era de idosos institucionalizados. Os resultado do estudo de Lausen et al. [43] sugerem que a mortalidade em pessoas com depressão é aproximadamente duas vezes maior em comparação com a amostra controle, e afirmam que a expectativa de vida em pessoas com depressão é ainda mais reduzida pela presença de doenças somáticas comórbidas ou abuso de substâncias.

O Transtorno Depressivo Maior conhecido também como Depressão Maior, Depressão Unipolar, ou apenas Depressão, é um distúrbio mental que afeta mais de 300 milhões de pessoas em todo o mundo, correspondendo a 4,4% da população mundial. Mais de 80% desta doença se concentra nos países onde a população possui renda média ou baixa [54]. Além disso, a depressão está associada a uma redução de expectativa de vida entre homens e mulheres de aproximadamente 10 (dez) e 7 (sete) anos, respectivamente [12]. O que demonstra que a redução da expectativa de vida é mais notável em homens do que em mulheres.

As causas de morte mais comuns associadas à depressão são os distúrbios cardiovasculares [63], derrame, síndromes metabólicas e câncer [36], sendo câncer de mama especificamente aumentado em mulheres [30]. Além disso, as causas naturais

são mais importantes do que as causas externas nos óbitos de pacientes com depressão unipolar [43].

As causas de morte de maneira geral são categorizadas em dois grupos, morte natural que é uma consequência devido a doenças ou mesmo por condições desconhecidas, e a morte não natural, que são mortes devido a causas externas [43], como mortes violentas, tais como, homicídios, acidentes, suicídios e mortes suspeitas.

Independente das causas de morte, as informações sobre suas causas devem constar na Declaração de Óbito (DO). A DO é um formulário padrão elaborado pela Organização Mundial da Saúde (OMS) e implementado praticamente em todos os países. O objetivo de se ter um formulário padrão para atestar óbito é para uniformizar a comparação de dados internacionalmente, além de ser suporte para produzir as estatística de mortalidade. A DO é um documento emitido pelo profissional da área médica e segue a legislação de cada país para atendimento de suas particularidades [42].

No Brasil, quando os óbitos são decorrentes de causa externa, ou seja, de morte violenta ou de morte suspeita, devem obrigatoriamente realizar exame de autópsia pelo Instituto Médico Legal (IML). No caso de morte natural, onde a morte aconteceu sem assistência médica ou mesmo por causas naturais desconhecidas, o exame de autópsia deve acontecer no Serviço de Verificação de Óbitos (SVO) [42].

A autópsia clínica é um importante procedimento realizado por médicos patologistas para determinar a causa da morte. A autópsia é baseada na combinação de avaliações macroscópicas e microscópicas e é importante como pré-condição para a descrição sistemática de doenças. Embora os avanços na ciência médica básica e nos exames de imagens médicas tenham sido grandes, a necessidade de autópsias continua sendo importante, principalmente para a garantia da qualidade e determinação da causa exata da morte [7], além de ser valiosa para manter estatísticas de mortalidade precisas, que permanecem essenciais para o planejamento de serviços públicos de saúde [21]. Quando a causa da morte é proveniente do exame da autópsia, são as informações realizadas pela autópsia que auxiliam no preenchimento da DO.

Trabalhos de revisão sistemática estabelecem uma relação entre causas de morte e depressão [14] identificam que há um risco aumentado de morte na depressão e que em muitos casos a depressão deve ser considerada um distúrbio com risco de vida. O trabalho de Wilson et al. [72] sugere um aumento de mortes relacionadas à depressão em algumas populações e associa à depressão a doença cardiovascular como causa da morte, porém, afirma que, para estimar esse efeito (seja um resultado direto da fisiopatologia da depressão ou resultado indireto da falta de autocuidado), são necessários estudos mais rigorosos.

1.2 Justificativa

Os dados relativos de quantas pessoas morrem, quais são suas características, como por exemplo, idade e sexo, mas principalmente, quais as causas responsáveis pelo óbito, são importantes para o conhecimento e o estabelecimento do perfil epidemiológico da população e de ações governamentais relacionadas a saúde pública [42].

A compreensão das causas de morte entre pessoas com transtornos mentais podem mostrar informações desconhecidas para abordar o problema e ampliar a discussão sobre o efeito dos transtornos mentais nestas mortes [67]. Identificar possíveis doenças associadas à morte em pacientes com depressão pode auxiliar na tomada de decisão das políticas públicas de saúde. Nesse sentido, o entendimento de possíveis doenças relacionadas à depressão, pode culminar em tratamento mais específico e em uma melhor expectativa de vida desses pacientes. Além disso, pode-se elaborar estratégias de prevenção, abrindo possibilidades de novos estudos fisiopatológicos.

Trabalhos prévios foram realizados para analisar de forma qualitativa e quantitativa, aliados principalmente a métodos estatísticos, a relação entre depressão e mortalidade [57, 69]. Contudo, existem poucos estudos nos quais a relação depressão e mortalidade são provenientes a partir de relatórios de autópsia. Assim, dada a relevância que a autópsia possui na exatidão da causa da morte, avaliar a depressão e as causas de morte utilizando dados da autópsia pode fornecer informações que serão importantes para se entender a relação entre depressão, mortalidade e, conseqüentemente, expectativa de vida.

1.3 Motivação

A realização deste trabalho tem como motivação o estudo sobre mortalidade e depressão em parceria com pesquisadores do Grupo de Estudos do Envelhecimento do Cérebro da Faculdade de Medicina da Universidade de São Paulo (FMUSP), que faz estudos em indivíduos com idade de 50 anos ou mais. Com o objetivo de entender o envelhecimento cerebral natural e patológico, este grupo de pesquisa criou o Banco de Cérebros, que possui dados coletados do Serviço de Verificação de Óbitos da Capital (SVOC-USP) - órgão que executa as autópsias obrigatórias para aqueles que morreram de causas naturais presumidas na região metropolitana de São Paulo.

Segundo Suemoto et al. [63], estes dados apresentam uma ampla gama de informações de uma grande amostra populacional, onde a causa imediata da morte pode ser acessada durante o exame de autópsia. Cabe ressaltar que grande parte dos estudos sobre depressão em indivíduos provenientes da comunidade (ou seja, não hospitalizados), e que avaliam as causas de morte, utilizam informações presentes no atestado de óbito e não por laudo de autópsia. Além disso, o diagnóstico de depressão foi baseado nos

critérios do Manual de Diagnóstico e Estatística dos Transtornos Mentais - 4ª edição (*Diagnostic and Statistical Manual of Mental Disorders - DSM-IV*), verificados através da aplicação de uma entrevista semiestruturada com um informante próximo do falecido, o que aumenta a especificidade do diagnóstico trazendo informação clínica mais precisa.

Por outro lado, não se encontrou na literatura o uso de técnicas de Mineração de Dados, ou mais especificamente, de métodos de Aprendizado de Máquina para extrair conhecimento de dados provenientes de exames de autópsia no estudo de mortalidade e depressão. Bzdok et al. [10] abordam que, por muitas décadas, o paradigma dominante de pesquisa com pacientes psiquiátricos tem encontrado dificuldades em sua realização e sugere o uso de algoritmos de Aprendizagem de Máquina como uma metodologia mais vantajosa que a análise que corresponde à estatística clássica. Librenza et al. [44] acrescentam que, no campo da psiquiatria, estudos sugerem as técnicas de Aprendizado de Máquina e Mineração de Dados como ferramentas importantes para o desenvolvimento de novas teorias e instrumentos, fornecendo métodos alternativos à estatística por inferência e auxiliando na criação de novas hipóteses.

1.4 Problema de pesquisa

O problema de pesquisa a ser abordado neste trabalho é um problema de classificação de indivíduos com e sem depressão. Deseja-se verificar se pessoas com depressão, comparadas às que não tem depressão, tem causas de morte diferentes, baseado em dados das doenças relacionadas à morte registradas em laudos de autópsia provenientes do SVOC-USP.

1.5 Objetivos

Este trabalho tem como objetivo geral analisar as doenças registradas em laudos de autópsias para verificar a associação e estabelecer as causas de morte em indivíduos com 50 anos ou mais que tinham depressão, buscando padrões de mortalidade destes indivíduos.

Como objetivos específicos, pode-se citar:

- Investigar a relação entre as causas múltiplas de morte declaradas nos laudos de autópsia e a depressão;
- Avaliar diferentes técnicas de aprendizado de máquina com o intuito de encontrar padrões que caracterizam as causas de morte em indivíduos com depressão;
- Comparar os resultados alcançados com a literatura médica.

1.6 Estrutura do Documento

O restante deste documento está organizado da seguinte maneira: o Capítulo 2 aborda os fundamentos teóricos que norteiam a pesquisa, tais como a Classificação Internacional de Doenças, Declaração de Óbito e Aprendizado de Máquina; o Capítulo 3 descreve o Mapeamento Sistemático relacionado a esta proposta. O Capítulo 4 descreve os Materiais e Métodos que foram utilizados nesta pesquisa; o Capítulo 5 apresenta os resultados obtidos dos experimentos, além da discussão em relação aos resultados dos experimentos. Por fim, no Capítulo 6 apresentam-se a conclusão acerca do trabalho, as contribuições, as limitações da pesquisa e os trabalhos futuros.

Referencial Teórico

Este capítulo apresenta o embasamento teórico necessário para o entendimento da pesquisa. A Seção 2.1 aborda conceitos relativos a Classificação Estatística Internacional de Doenças e Problemas Relacionados a Saúde (CID); a Seção 2.2 apresenta sobre Declaração de Óbito e sobre Causa de Morte; a Seção 2.3 compreende os conceitos sobre Autópsia. a Seção 2.4 apresenta conceitos sobre Aprendizado de Máquina e os tipos de aprendizado.

2.1 Classificação Internacional de Doenças (CID)

A Classificação Estatística Internacional de Doenças e Problemas Relacionados com a Saúde (CID) é o padrão mundial para gerar estatísticas de causa de morte [42]. Grande parte das despesas com saúde do mundo (70%) usam a CID para cálculo para reembolso e alocação de recursos; 110 países que coletivamente respondem 60% da população mundial usam dados de causas de morte da CID para planejar e monitorar de maneira sistemática as questões de saúde de seus países; e, a CID-10 é citada em mais de 20.000 artigos científicos [70].

O objetivo de ter uma CID padronizada para a coleta de dados de saúde é gerar condições de saúde comparáveis em nível internacional [70]. Em 1989, a Organização Mundial de Saúde (OMS) aprovou a sua décima edição, a CID-10 [42]. A estrutura da CID-10 é uma lista composta por um código alfanumérico com 4 posições, sendo uma letra na primeira posição e um número na segunda, terceira e quarta posição. A quarta posição é um caractere que segue após uma casa decimal. Portanto, os códigos possíveis, variam de A00.0 a Z99.9 [70].

A classificação da CID-10 é dividida em 22 capítulos. O primeiro código da CID-10 corresponde a um capítulo específico. De maneira hierárquica, para cada capítulo existem os grupos de doenças, dentro dos grupos de doença existem as categorias e para cada categoria as subcategorias. Na tabela 2.1, é apresentado os capítulos da CID-10 e os grupos de doenças correspondente a cada capítulo.

Capítulo	Descrição	Cód. grupos de doença
I	Algumas doenças infecciosas e parasitárias	A00 - B99
II	Neoplasias [tumores]	C00 - D48
III	Doenças do sangue e dos órgãos hematopoéticos e alguns transtornos imunitários	D50 - D89
IV	Doenças endócrinas, nutricionais e metabólicas	E00 - E90
V	Transtornos mentais e comportamentais	F00 - F99
VI	Doenças do sistema nervoso	G00 - G99
VII	Doenças do olho e anexos	H00 - H59
VIII	Doenças do ouvido e da apófise mastóide	H60 - H95
IX	Doenças do aparelho circulatório	I00 - I99
X	Doenças do aparelho respiratório	J00 - J99
XI	Doenças do aparelho digestivo	K00 - K93
XII	Doenças da pele e do tecido subcutâneo	L00 - L99
XIII	Doenças do sistema osteomuscular e do tecido conjuntivo	M00 - M99
XIV	Doenças do aparelho geniturinário	N00 - N99
XV	Gravidez, parto e puerpério	O00 - O99
XVI	Algumas afecções originadas no período perinatal	P00 - P96
XVII	Má formação congênita, deformidades e anomalias cromossômicas	Q00 - Q99
XVIII	Sintomas, sinais e achados anormais de exames clínicos e de laboratório, não classificados em outra parte	R00 - R99
XIX	Lesões, envenenamento e algumas outras consequências de causas externas	S00 - T98
XX	Causas externas de morbidade e de mortalidade	V01 - Y98
XXI	Fatores que influenciam o estado de saúde e o contato com os serviços de saúde	Z00 - Z99
XXII	Códigos para propósitos especiais	U04 - U99

Tabela 2.1: Capítulos da CID-10 com os grupos correspondentes. Como exemplo, pode-se perceber que os grupos de doenças entre A00 a B99 pertencem a Categoria I

Fonte: <http://tabnet.datasus.gov.br/>

Além dos capítulos e grupos, a CID ainda é composta de categorias e subcategorias onde são determinadas particularidades das doenças. Abaixo segue um exemplo da construção de um código da CID-10:

- Capítulo: Transtornos mentais e comportamentais (F00-F99);
- Grupo: Transtornos do humor (afetivos) (F30-F39);
- Categoria: Transtorno depressivo recorrente (F33);
- Subcategoria: Transtorno depressivo recorrente, atualmente em remissão (F33.4)

De acordo com o exemplo acima, a CID-10 que seria apresentada em um atestado médico, seria CID-10 F33.4.

2.2 Declaração de Óbito

A declaração de óbito (DO) é o documento que confirma o óbito de uma pessoa e é fornecido pelo médico ao qual vinha prestando assistência (exceto em situações específicas) ao paciente. Ele possui a finalidade de definir uma *causa mortis* e responde aos interesses de ordem legal, ética e médico-sanitária [55]. O modelo da DO em vigor é composto por 9 blocos divididos em 59 campos [42]. Dentre os blocos da DO o bloco 5 é o mais relevante para este estudo, uma vez que são os dados referentes a este bloco que serão utilizados nesta pesquisa. No Brasil, assim como na maioria dos países, os códigos da CID-10 são a referência para o preenchimento dos campos da DO relacionados às causas da morte.

O bloco 5 da DO (Figura 2.1) se refere às condições e causas do óbito, e sua importância decorre do fato de ser a fonte da causa básica do óbito e dos agravos que contribuíram para óbito. Além disso trás informações como o fato do indivíduo ter recebido assistência médica durante a doença que levou à morte e ter sido ou não realizada a autópsia. Basicamente, o bloco 5 é formado por 4 campos, descritos abaixo [42]:

- *Óbito de Mulher em Idade Fértil*: informação em caso de óbito de mulher em idade fértil (no Brasil, considerada de 10 a 49 anos). Trata-se de uma informação importante e que representa subsídio para melhor conhecimento das mortes maternas.
- *Assistência Médica*: refere-se ao atendimento médico continuado que o paciente recebeu, ou não, durante a enfermidade que ocasionou o óbito.
- *Diagnóstico confirmado por*: corresponde à execução ou não da autópsia.
- *Causas da Morte*: preenchido pelo médico que atestou o óbito, deve apresentar o diagnóstico mais preciso da causa e das circunstâncias da morte, é dividido em duas partes, como pode ser observado na seção 2.2.

ÓBITO DE MULHER EM IDADE FÉRTIL		ASSISTÊNCIA MÉDICA		DIAGNÓSTICO CONFIRMADO POR:	
<input checked="" type="checkbox"/> A morte ocorreu <input type="checkbox"/> Na gravidez <input type="checkbox"/> No aborto <input type="checkbox"/> De 43 dias a 1 ano após o parto <input type="checkbox"/> Ignorado <input type="checkbox"/> No parto <input type="checkbox"/> Até 42 dias após o parto <input type="checkbox"/> Não ocorreu nestes períodos <input type="checkbox"/> Ignorado		<input checked="" type="checkbox"/> Recebeu assist. médica durante a doença que ocasionou a morte? <input checked="" type="checkbox"/> Sim <input type="checkbox"/> Não <input type="checkbox"/> Ignorado		<input type="checkbox"/> Necrópsia? <input type="checkbox"/> Sim <input checked="" type="checkbox"/> Não <input type="checkbox"/> Ignorado	
V Condições e causas do óbito	CAUSAS DA MORTE PARTE I Doença ou estado mórbido que causou diretamente a morte.		ANOTE SOMENTE UM DIAGNÓSTICO POR LINHA		
	CAUSAS ANTERCEDENTES Estados mórbidos, se existirem, que produziram a causa acima registrada, mencionando-se em último lugar a causa básica.				
	PARTE II Outras condições significativas que contribuíram para a morte, e que não entraram, porém, na cadeia acima.				
	a	Infarto agudo miocárdio	Tempo aproximado entre o início da doença e a morte	1 dia	CID
	b	Isquêmica miocárdica	anos		
c	Aterosclerose generalizada	anos			
d	Diabetes mellitus	20 anos			
	Hipertensão arterial	-			

Figura 2.1: Bloco 5 da Declaração de Óbito com os respectivos campos: Óbitos de Mulher em Idade Fértil, Assistência Médica, Diagnóstico Confirmado por, Causas da Morte [42].

Causas da Morte

O campo causas da morte é composta de 2 partes, como se segue [42]:

- Parte 1: referente à causa que provocou a morte e as possíveis causas que contribuíram para a morte.
 - Linhas ‘a’, ‘b’ e ‘c’: complicações da causa básica. São possíveis causas que podem ter contribuído para a causa terminal ou imediata; a linha ‘a’ refere-se à causa imediata ou terminal, e as linhas ‘b’ e ‘c’, às causas intermediárias.
 - Linha ‘d’: causa básica da morte, informação mais importante sob a ótica epidemiológica, evento que iniciou o processo da morte, caracterizado como a doença ou lesão que inicia a sequência de estados mórbidos, ou as circunstâncias do acidente ou da violência, que levaram diretamente à morte [55].
- Parte 2: devem constar as doenças que o paciente era portador, mas que não estejam diretamente relacionadas à causa terminal da morte [55].

Conforme o exemplo da Figura 2.1, é possível notar que a causa imediata do óbito foi proveniente de um “infarto agudo do miocárdio”, as causas intermediárias foram as doenças “isquêmica miocárdica” e “aterosclerose generalizada”, e a causa básica da morte foi “diabetes mellitus”. Além disso, existe uma doença que está diretamente relacionada à causa imediata da morte que é a “hipertensão arterial”. Para fins de estatísticas oficiais de mortalidade a causa selecionada é sempre a causa básica da morte. Do ponto de vista de prevenção da morte, a causa básica da morte é a causa precipitante que inicia a sucessão de eventos mórbidos que provocam a morte [42].

De acordo com Laurenti et al. [41] as estatísticas de mortalidade refletem muito pouco sobre a frequência dos transtornos mentais na população, visto que, diferentemente de outros agravos à saúde, como por exemplo, as doenças infecciosas e parasitárias,

as doenças dos aparelhos circulatório, respiratório, digestivo e as causas externas, os transtornos mentais são pouco citados como causas básica de morte.

Ishitani e França [34] afirmam que dificilmente uma morte tem uma causa única, uma vez que além da doença ou doenças que determinam a morte, existe uma grande frequência de doenças presentes ou mesmo associadas ou não entre si, que contribuem direta ou indiretamente no processo que leva à morte.

Os dados de múltiplas causas de morte têm o potencial de ajudar a apontar possíveis preocupações quanto à precisão, além de fornecer uma imagem mais completa da mortalidade por causas que frequentemente não são registradas como a causa subjacente (básica) da morte [68]. Nesse contexto, as estatísticas de mortalidade sob o enfoque das causas múltiplas de morte podem retratar melhor o perfil de mortalidade quando se aproveitam as informações sobre todas as causas de morte mencionadas na DO. Além disso, os estudos com causas múltiplas de morte podem apresentar um melhor conhecimento da morbidade da população e da frequência de óbitos por aquelas doenças que não são selecionadas como causa básica, bem como da possibilidade do estudo da associação de causas [34].

Com o avanço da computação na análise de dados médicos, é possível ter análises importantes sobre as múltiplas causas de morte. Laurenti e Jorge [42] apresentam que a possibilidade de pesquisas que envolvam a computação na análise de dados possibilita estudos de associação de causas, além de permitir que diagnósticos com pouca probabilidade de serem selecionadas como causa básica possam ser conhecidos e mensurados.

2.3 Autópsia

Na literatura, desde os primórdios, pode-se observar que a autópsia desempenha um papel importante na história da humanidade, principalmente no que se refere ao entendimento do corpo humano e conseqüentemente na elucidação das doenças [37]. Doenças como a Síndrome da Imunodeficiência Adquirida (AIDS) e Alzheimer são exemplos de doenças que foram melhor entendidas por meio da autópsia [28].

O exame de autópsia, necropsia ou do latim *post-mortem*, independente do termo utilizado, é um exame do corpo após a morte. Autópsia significa “ver por si próprio” e é geralmente considerada sinônimo do termo necropsia e *post-mortem* (após a morte) [28]. Por ser um termo mais convencional, para esse estudo o termo utilizado será autópsia.

Os objetivos da autópsia é determinar a causa da morte, determinar como ocorreu a morte bem como os processos patológicos envolvidos, além de adquirir informações mais confiáveis sobre a natureza e a causa da doença [27, 28]. Por meio da autópsia é

possível avaliar os efeitos de novas doenças e novos tipos de terapia medicamentosa e outras intervenções terapêuticas [28].

Além disso, em estudos sobre autópsia, há evidências da falta de precisão na declaração da causa de morte quando declarações de óbito originais foram comparadas com os dados obtidos em relatórios de autópsias. Por exemplo, certas causas de morte teriam um maior número de declarações, tais como, tuberculose (16,7% a mais), doença reumática crônica do coração (22,0% a mais), hiperplasia da próstata (33,3% a mais) e as anomalias congênitas (14,5% a mais) [35].

No Brasil, a autópsia é realizada pelo SVO ou pelo IML. No caso de morte natural sem assistência médica, ou com assistência médica, mas sem causa estabelecida, a morte será analisada pelo SVO. A morte natural é a consequência de um processo esperado e previsível que podem estar atribuídas a uma doença ou mesmo ao envelhecimento natural. A morte não natural são as ditas violentas, que podem ser decorrentes, por exemplo de homicídio, suicídio, intoxicação e acidentes de trânsito. A Figura 2.2 apresenta um resumo da realização da autópsia no Brasil de acordo com a motivação do óbito.

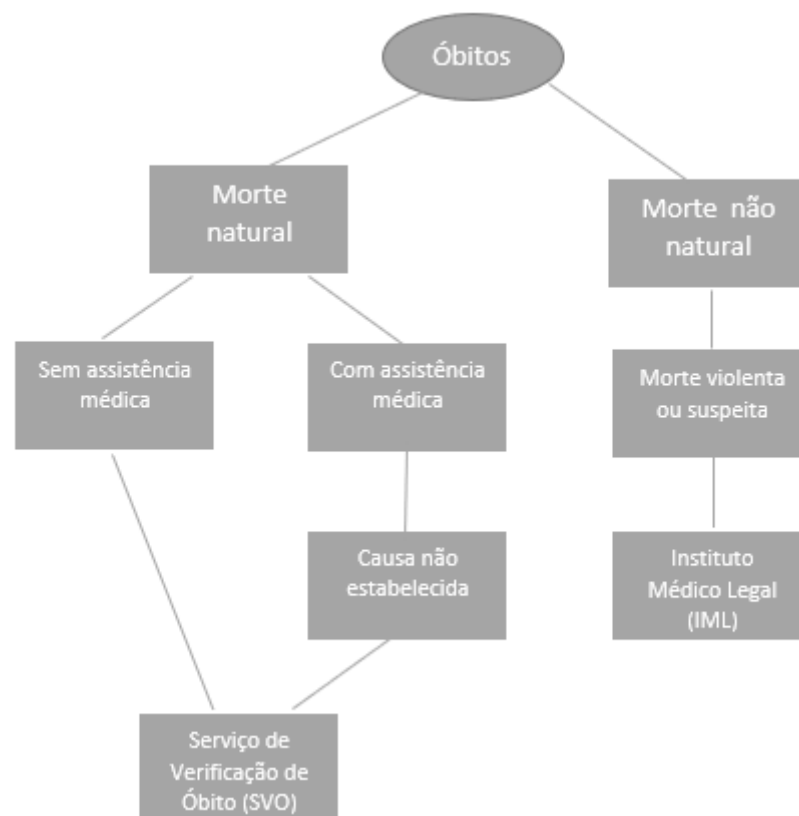


Figura 2.2: Encaminhamento dos óbitos no Brasil de acordo com a motivação

Os SVOs possibilitam a detecção das emergências epidemiológicas, o diagnóstico isolado ou surtos de doenças emergentes e reemergentes e ainda agravos fora do comum. Nesse contexto, é possível ter informações para orientar a tomada de decisão para o controle das doenças e permitir o aprimoramento da qualidade da informação de mortalidade para subsidiar as políticas de saúde [19]. Algumas diferenças metodológicas em pesquisas sobre mortalidades estão relacionadas a origem da população estudada: alguns estudos são de amostra de comunidade (não hospitalizado), outros de serviços especializados ou de pacientes internados (hospitalizado). Ao realizar estudos onde a amostra é de uma comunidade, é possível analisar a realidade dos casos menos graves de depressão, que acontecem com a maioria da população. Muitas formas de depressão não são reconhecidas ou são tratadas na rede primária de saúde, não refletindo nos estudos sobre depressão e causas de mortes.

A autópsia fornece meios de avaliar a prestação de cuidados de saúde e assim, desempenha um importante serviço social. Além disso, a autópsia confirma a precisão da causa da morte. No entanto, ao validar a certificação de morte, torna as estatísticas vitais mais precisas e significativas. [28]

Nesse contexto, dada a relevância que a autópsia possui na exatidão da causa morte, avaliar a depressão e as causas de morte utilizando dados da autópsia, pode fornecer informações que serão importantes para se entender a relação entre depressão, mortalidade e, conseqüentemente expectativa de vida.

2.4 Mineração de Dados

A mineração de dados refere-se à aplicação de algoritmos específicos para descoberta e extração de padrões em base de dados e se baseia em técnicas conhecidas de aprendizagem de máquinas, reconhecimento de padrões e estatísticas para encontrar tais padrões. Neste contexto, o objetivo da mineração de dados é encontrar padrões em base de dados até então desconhecidos, e uma vez encontrados esses padrões, eles podem ser usados para tomada de decisões.

Para [4, 31], a mineração de dados é sinônimo do processo de Descoberta de Conhecimento em Bases de Dados (KDD, do inglês, *Knowledge Discovery in Databases*). O processo KDD, pode ser definido como um processo não trivial de identificação de novos padrões válidos, úteis e compreensíveis [20]. No entanto, ainda que não há consenso na definição dos termos KDD e mineração de dados. De acordo com [20] o processo KDD refere-se a todo o processo de descoberta de conhecimento, e a mineração de dados é uma das atividades deste processo. Porém, há uma concordância entre os autores, em que o processo de mineração deve ser iterativo, interativo e dividido em fases.

A mineração de dados enquanto processo, envolve a coleta e seleção de dados, o pré-processamento de dados, análise de dados, incluindo a visualização de resultados, interpretação de descobertas e a aplicação do conhecimento. Para pré-processar e analisar dados, algoritmos de aprendizado de máquina e métodos estatísticos são utilizados nesta etapa [61]. Além disso, as análises referente ao processo de mineração de dados podem ser distinguidas em descritivas, onde o conhecimento é representado na forma de modelos que retratam padrões e relações em dados e preditivos, onde o conhecimento é representado em previsões sobre eventos futuros.

Aprendizado de Máquina

O Aprendizado de Máquina (AM) é um campo de pesquisa da Inteligência Artificial que estuda algoritmos capazes de fazer descobertas automáticas de padrões a partir de informações em conjuntos de dados. Tais padrões podem ser usados para fazer previsões, além de dar suporte no processo de tomada de decisão [2].

No AM, os algoritmos são programados para aprender com a experiência passada e empregam o princípio de inferência denominado indução [11]. Os algoritmos de AM induzem uma hipótese ou função a partir de um conjunto de dados de treinamento que deve ser capaz de ser *generalizada* para novos exemplos de um mesmo domínio, ou seja, deve ser válida para um conjunto de dados diferente daquele utilizado no seu aprendizado.

Quando a hipótese apresenta uma baixa capacidade de generalização, a razão pode ser que o modelo gerado (hipótese) pode estar superajustado (*overfitting*) ou subajustado (*underfitting*) aos dados de treinamento [11]. Segundo Tan et al [66], *overfitting* e *underfitting* são problemas que estão relacionados à complexidade do modelo e sugere que as causas desses problemas podem ser devido à presença de ruído ou mesmo à falta de amostras representativas [20].

As formas de aprendizado de máquina podem ser classificadas como: aprendizado supervisionado, aprendizado não-supervisionado, aprendizado semi-supervisionado e aprendizado por reforço [2].

No aprendizado supervisionado, a indução da hipótese é feita em função de uma variável de interesse (variável dependente). Algoritmos de aprendizado supervisionado recebem um conjunto de dados rotulados, ou seja, dados onde os valores da variável de interesse (chamada também de variável de saída, classe ou atributo alvo) são conhecidos. O conjunto de dados é dividido em um conjunto de treinamento, usado para a indução do modelo por meio de um algoritmo de aprendizado; e um conjunto de teste, usado para a aplicação do modelo gerado para predizer as classes dos exemplos deste conjunto. Com isso, o algoritmo mapeia uma função a partir dos valores dos dados de entrada e dos rótulos do atributo alvo na etapa conhecida como *treinamento*. Em seguida, essa

função, conhecida também como *modelo de classificação* [66], pode ser utilizada em um novo conjunto de dados para prever valores de saída, na etapa conhecida como *teste*. Os valores preditos no conjunto de teste são comparados com os rótulos dos dados para a avaliação do modelo. A Figura 2.3 representa o aprendizado supervisionado. Em geral, os algoritmos de aprendizado supervisionado são usados para gerar modelos preditivos para tarefas de *classificação* (quando a variável dependente possui valores discretos) ou *regressão* (quando a variável dependente possui valores contínuos).

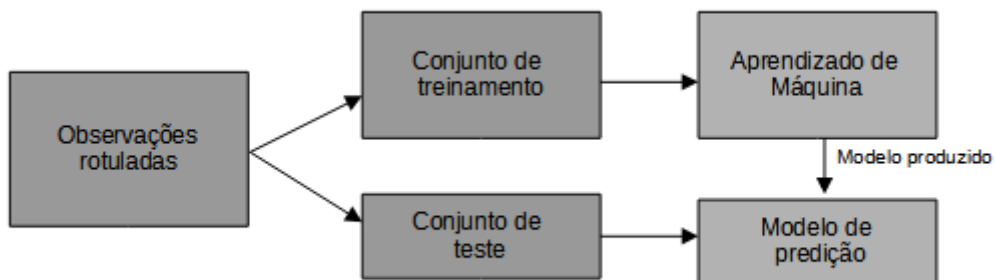


Figura 2.3: Aprendizado supervisionado

No aprendizado não-supervisionado, a indução da hipótese é feita em função da estrutura intrínseca dos dados, não havendo um variável dependente para guiar o aprendizado, ou seja, os dados não são rotulados. Assim, os algoritmos não supervisionados, se referem a identificar informações relevantes nos dados sem a presença de um elemento externo para guiar o aprendizado [11], isto é, seu objetivo é descrever um conjunto de dados, como por exemplo, encontrar um grupo de objetos semelhantes ou determinar associações entre os atributos de um conjunto de dados. A Figura 2.4 representa o modelo de aprendizado não supervisionado. Os algoritmos que utilizam este tipo de aprendizado geram modelos descritivos para tarefas tais como segmentação (*clustering*), associação ou sumarização.

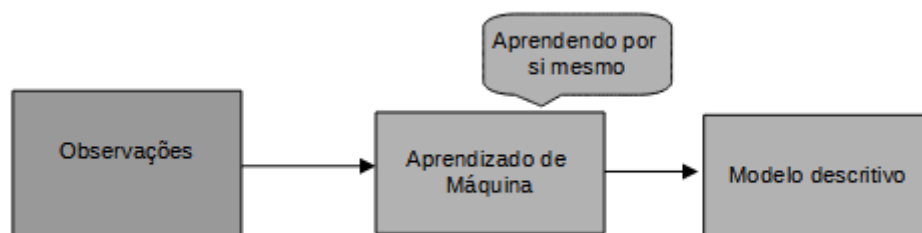


Figura 2.4: Aprendizado não supervisionado

No aprendizado semi-supervisionado, dados não rotulados são acrescentados ao conjunto de treinamento com intuito de aumentar a eficiência do classificador. Inicialmente, o algoritmo é treinado com um pequeno grupo de dados rotulados. Com o classificador treinado, são então classificados os dados não rotulados. Os dados que conseguirem um alta confiança na classificação são adicionados ao conjunto de treinamento. Dessa maneira o algoritmo é novamente treinado com o novo conjunto. Esse processo é repetido até que os dados classificados tem em sua maioria uma alta confiança. Assim, o classificador ensina a si mesmo, com suas próprias previsões [48]. Algoritmos de classificação semi-supervisionada são especialmente úteis em situações onde o conjunto de dados não é representativo o suficiente para o aprendizado indutivo de um classificador. A Figura 2.5 demonstra a técnica de aprendizado semi supervisionada.

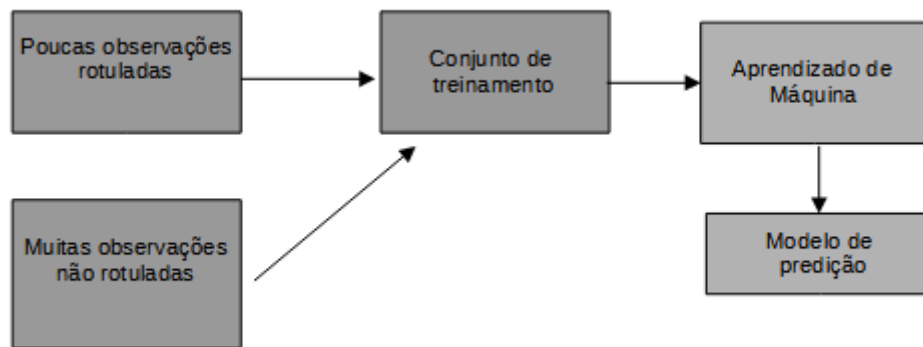


Figura 2.5: Aprendizado semi supervisionado

O aprendizado por reforço é diferente dos aprendizados anteriores porque não há conjunto de treinamento, rotulado ou não. O aprendizado por reforço trata de problemas onde um determinado agente de software deve agir em determinados ambientes de modo a maximizar alguma noção de recompensa (ou reforço) cumulativa. Isso quer dizer que um agente precisa saber que algo bom aconteceu quando ganhar ou que algo ruim aconteceu quando perder. O que define uma recompensa é a tarefa que o agente deve desempenhar em encontrar uma política π , que mapeia estados em ações, que maximiza a medida de reforço. O aprendizado por reforço é representado pela Figura 2.6 Os tipos de aprendizagem por reforço são aprendizagem passiva e ativa. A aprendizagem passiva utiliza uma representação baseada em estados em um ambiente completamente observável. Além disso, a política π do agente é fixa e a tarefa consiste em aprender o estado-ação. Na aprendizagem ativa o agente pode escolher as ações e deve aprender uma política ótima, pois não há uma política π fixa, nesse caso o agente deve experimentar tanto quanto possível do seu ambiente a fim de aprender como se comportar nele [58].

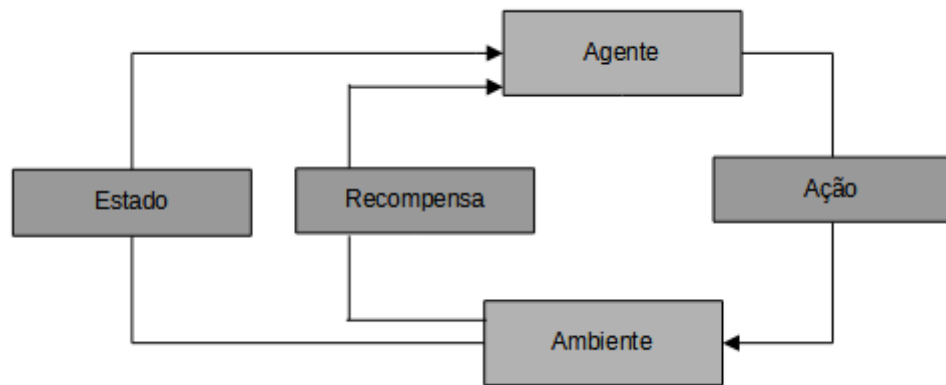


Figura 2.6: Aprendizado por reforço

2.5 Psiquiatria Computacional

A psiquiatria computacional é um campo heterogêneo, que incorpora métodos da psiquiatria, psicologia, neurociência, economia comportamental e aprendizado de máquina [33]. As variáveis que envolvem as doenças psiquiátricas são complexas e a psiquiatria computacional fornece duas maneiras de abordar essa complexidade: análise de dados guiados por dados teoricamente agnósticos, que nesse caso utiliza-se AM; e, análise de dados orientados por teoria, onde os modelos resumem uma compreensão teórica, muitas vezes mecanicista, dos fenômenos em questão [33, 59].

Huys et al. [33] apresentam que a análise de dados guiada por dados, pode ser limitada em sua capacidade de capturar as complexidades das variáveis e as suas interações e sugere que o campo da psiquiatria faça combinação das abordagens orientadas por teoria e orientada por dados, por essa combinação ser especialmente poderosa e promissora nos resultados. Nesse sentido, Rutledge et al. [59], afirmam que as abordagens que combinam grandes quantidade de dados, aprendizado de máquina e desenvolvimento orientado por teoria têm um enorme potencial que podem trazer informações úteis que darão suporte para melhorar o tratamento de doenças mentais.

A geração de conhecimento em neurociência básica e tomada de decisão clínica em psiquiatria foi fundamentada em estatísticas clássicas com testes formais para diferenças em amostras frequentemente pequenas [10]. Nesse contexto, Rutledge et al. [59] apresentam que abordagens de AM têm contribuído para a descoberta de informações que conduzem a um melhor entendimento sobre os distúrbios psiquiátricos e, consequentemente tem colaborado para aliviar o sofrimento resultante desses distúrbios.

2.6 Considerações Finais

Neste capítulo foram descritos os principais conceitos concernentes a este trabalho. Inicialmente, foram apresentadas uma breve explicação sobre a formação da estrutura da CID e da DO. A importância de se entender a estrutura da CID é que ela codifica o que provocou a morte e as possíveis causas que contribuíram para a morte, como consta na DO, e que também são relacionados às causas de morte no laudo de autópsia. Além disso, os dados das causas de morte provenientes do exame de autópsia são considerados precisos, uma vez que a investigação destas causas é realizado por um médico patologista que faz o exame macroscópico e microscópico do cadáver.

Um outro conceito apresentado neste capítulo, é sobre o uso de AM em pesquisas relacionadas às doenças mentais. No estudo de Librenza et al. [44], os autores apresentam que a medicina baseada em evidências usando métodos estatísticos clássicos contribuíram para entender fatores de risco e bons tratamentos para os transtornos mentais. No entanto, é possível perceber que as técnicas que visam o desenvolvimento de cuidados psiquiátricos, como aprendizado de máquina, vêm ganhando espaço na pesquisa psiquiátrica [59, 44].

Mapeamento Sistemático

O Mapeamento Sistemático (MS) é um estudo secundário que possibilita ter uma visão geral dos estudos relevantes de uma determinada pesquisa a partir de um protocolo definido. O processo de MS envolve as fases de planejamento, condução e publicação dos resultados, executadas de forma iterativa [53]. Para a condução deste estudo foram utilizadas a ferramenta Parsifal, que permitiu executar a fase de planejamento, a ferramenta Mendeley, onde foi possível gerenciador as referências bibliográficas selecionadas e o suporte de uma planilha eletrônica, onde foi realizado uma tabela dos estudos primários finais selecionados, no qual foram realizadas diversas anotações importantes destes estudos. A Seção 3.1 apresenta o Planejamento do MS, onde é descrito o protocolo de pesquisa. A Seção 3.2 aborda a condução desta pesquisa, apresenta o comportamento da *string* de busca na base de dados selecionada e a quantidade de estudos selecionados. Em seguida na Seção 3.3 é apresentado os resultados em relação às questões de pesquisa. Por último, na Seção 3.4, as considerações finais do capítulo.

3.1 Planejamento

A fase de planejamento é a primeira etapa do MS e consiste em criar um protocolo que descreve o processo da pesquisa. Permite identificar a avaliação de qualidade que são questões que definem o escopo dos artigos aceitos. Além disso, é indicado as palavras-chave e seus sinônimos, a *string* de busca genérica, as bases de dados de pesquisa e os critérios de inclusão e exclusão. Os parâmetros utilizados nessa fase estão descritos a seguir.

Protocolo

- Objetivo

Identificar na literatura estudos primários que relacionem doenças que motivaram a causa de morte em indivíduos com depressão.

- *Population, Intervention, Comparison, Outcome, Context (PICOC)*

O PICOC do acrônimo População, Intervenção, Comparação, Resultado e Contexto, auxilia na identificação das palavras-chave e conseqüentemente na construção da *string* de busca. São informações que ajudam na investigação da pesquisa que deseja realizar. Segue a estrutura PICOC deste estudo:

- a. População: Grupo de pessoas diagnosticada com depressão, no qual foi possível detectar a motivação da morte;
- b. Intervenção: Métodos computacionais e estatísticos que apoiaram os estudos em relação a causa de morte e a depressão;
- c. Comparação: Literatura médica que relaciona depressão e causa de morte;
- d. Resultado: informações sobre técnicas e métodos utilizadas para relacionar depressão e causa de morte;
- e. Contexto: População efetiva do estudo compreende em indivíduos com e sem depressão onde foi possível verificar qual a causa de morte.

- Questões de pesquisas

As questões de pesquisa de um estudo de Mapeamento Sistemático são genéricas, uma vez que este tipo de estudo fornece uma visão mais ampla e geral, normalmente de caráter exploratório, de um tópico de pesquisa [53]. No entanto, a questão de pesquisa que norteou este MS foi a seguinte:

Indivíduos com depressão morrem de causas diferentes quando comparados a indivíduos sem depressão?

Adicionalmente, outras questões de pesquisas foram usadas para direcionar a busca de estudos em relação as publicações existentes na literatura, como se segue:

- a. O paciente foi diagnóstico com depressão? Qual método utilizado para este diagnóstico.
- b. Qual a idade média dos pacientes que foram avaliados no estudo?
- c. Os pacientes avaliados eram de uma comunidade não hospitalizada? Quais eram as condições da população em relação ao estudo no que tange à interferência e ao tempo
- d. Quais métodos (estatístico ou computacionais) utilizados para relacionar a depressão e a causa de mortis?
- e. Os dados analisados são provenientes de exame de autópsia?

- Identificação de Palavras-Chave e Sinônimos

As palavras-chave foram identificadas a partir das questões de pesquisa e conforme a estrutura: PICOC, porém foram consideradas relevantes apenas parte da estrutura

PICOC, como: População, Intervenção e Contexto. A tabela 3.1 apresenta as palavras-chaves e seus sinônimos e o relacionamento entre a parte da estrutura PICOC.

Palavra-chave	Sinônimo	Relacionado com
autopsy	<i>dissection, necropsy, postmortem</i>	Intervenção
causa mortis	<i>death, mortality</i>	Contexto
depression	<i>depressive disorder, mental disease</i>	População

Tabela 3.1: Palavras-chave e seus sinônimos de acordo com parte da estrutura PICOC.

- Geração de *strings* de busca

Uma *string* de busca é um agrupamento de palavras e termos referente ao tema de pesquisa conectados por operadores lógicos, tais como AND, OR e NOT e que tenha relação com as palavras-chave. Com o intuito de validar a *string* de busca, inicialmente foi realizado buscas preliminares na base de dados selecionada (PubMed) juntamente com a colaboradora desta pesquisa e especialista em psiquiatria. A busca preliminar, ou busca piloto é essencial para refinar e/ou estender a *string* de busca, além de aumentar a chance de viabilizar uma busca eficaz e eficiente [53]. Assim, alguns métodos de filtragem foram adotados, como o título (*title*) e o resumo (*abstract*) das publicações. Após a busca piloto, percebeu-se que a *string* de busca adequada foi a seguinte:

(depression[Title] AND mortality[Title]) NOT (suicide[Title/Abstract])

- Base de busca

A base de busca selecionada foi a PubMed, que é um motor de busca para a *Medical Literature Analysis and Retrieval System Online* (Medline). A Medline é uma base de dados online de acesso gratuito a referências e resumos de revistas científicas da área Biomédica. A revisão e o acompanhamento da especialista foi essencial para identificar os estudos primários relevantes e verificar que a PubMed seria o suficiente para o MS. A PubMed pode ser acessada pelo link: <https://www.ncbi.nlm.nih.gov/pubmed/>

- Seleção de estudos primários

A seleção de estudos primários foi dividida em: critérios de inclusão, critérios de exclusão e o processo de seleção dos estudos. Os critérios de inclusão e exclusão determina o rigor da pesquisa e impossibilita o viés do pesquisador no momento da seleção [50]. A seguir segue os critérios principais de inclusão e exclusão.

Critério de inclusão:

- a. Estudos que associam os motivos de mortalidade quando comparados a indivíduos com e sem depressão.
- b. Estudos que abordam a causa de mortis de indivíduos com depressão por meio de autópsia ou por outro tipo de registro.

Critério de exclusão:

- a. Estudos que estudam comorbidades muito específicas e associam os motivos de mortalidade quando comparados a indivíduos com e sem depressão.
- b. Estudos que estudam populações específicas e associam os motivos de mortalidade quando comparados a indivíduos com e sem depressão.
- c. Estudos cujo o objetivo é estudar diferenças entre características de pessoas com depressão, tais como sexo e idade.
- d. Estudos que comparam grupos de indivíduos com e sem depressão para avaliar previsão de mortalidade.
- e. Estudos específicos (como falta de vitaminas? - como descrever?, ou de comunidades muito específicas) que fazem relação entre depressão e mortalidade.

3.2 Condução do processo de seleção dos estudos primários

Para validar documentos e procedimentos do Mapeamento, foi realizado reuniões com a colaborada e especialista deste projeto e o orientador deste trabalho.

Durante a condução deste estudo, os estudos primários foram identificados conforme o processo descrito abaixo:

- a. Refinamento da *string* de busca, por meio de testes piloto;
- b. Após a busca dos estudos por meio da *string* de busca, os artigos foram catalogados na ferramenta Mendeley Desktop.
- c. A partir da leitura do título e do resumo, os artigos foram avaliados quanto aos critérios de inclusão e exclusão e o resultado foi registrado em uma planilha eletrônica;
- d. Em seguida, os estudos incluídos, foram avaliados individualmente por mim e pela especialista. Nesta etapa, quando houve dúvida pela inclusão ou exclusão do estudo, foi possível recorrer a especialista e decidir pela inclusão ou exclusão;
- e. Foram realizadas reuniões para validar os estudos. Nesta etapa, a especialista auxiliou no entendimento dos artigos quando não houve uma unanimidade na inclusão de um artigo;

- f. Em seguida, os estudos selecionados foram avaliados novamente em reuniões conjuntas com a especialista.

O MS foi realizado entre o meses de fevereiro e março de 2021. Foi definido que a busca seria realizada em relação aos últimos 10 anos, ou seja, entre 2001 a 2021. Com a aplicação da *string* de busca na base da PubMed foram encontrados 637 artigos. Após a análise dos critério de inclusão e exclusão, foi realizado a leitura do título e resumo para os resultados retornados, obtendo-se 61 artigos. Por fim, foi realizada uma leitura dos artigos por completo resultando em 11 artigos relevantes. Na Seção 3.3 é apresentado quais estudos foram selecionados, respondendo às questões de pesquisa e uma breve explicação sobre o achado principal de cada estudo é apresentado.

3.3 Análise do Resultado sobre às Questões de Pesquisa

A seguir os resultados do estudo são apresentados como respostas às questões de pesquisa apresentadas na Seção 3.1 item Questões de Pesquisa.

A Tabela 3.2 apresenta quais testes foram utilizados para o diagnóstico da depressão nos estudos selecionados e responde à questão de pesquisa: “O paciente foi diagnóstico com depressão? Qual teste utilizados para o diagnóstico”. Pode se observar que em todos os estudos os pacientes foram diagnósticos com depressão utilizando diversos critérios, dos quais dois estudos [62, 17] usaram o teste *General Health Questionnaire (GHQ)*, outros dois estudos [51, 60] usaram *World Health Organization Composite International e Diagnostic Interview Short-Form (CIDI-SF)*. Os demais estudos usaram critérios diferentes entre si.

Autor (Referência)	Método de avaliação da depressão
STEWART, Ralph AH <i>et al.</i> , 2003 [62]	General Health Questionnaire (GHQ)
DINIZ, Breno S. <i>et al.</i> , 2014 [17]	General Health Questionnaire (GHQ-12)
MENG, Ruiwei, <i>et al.</i> , 2020 [51]	Chinese version of the World Health Organization Composite International Diagnostic Interview–Short Form and a 7-item symptoms questionnaire modified from the Composite International Diagnostic Interview–Short Form (CIDI-SF)
SAINT ONGE, Jarron M.; KRUEGER, Patrick M.; ROGERS, Richard G., 2014 [60]	World Health Organization Composite International, Diagnostic Interview Short-Form (CIDI-SF)
WU, C.S.; HSU, L.Y.; WANG, S.H., 2020 [71]	Psychiatric Central Research Register Danish
PENNINX, Brenda WJH <i>et al.</i> , 2001 [56]	Diagnostic and Statistical Manual of Mental Disorders (DSM-III) – major depression; Center for Epidemiologic Studies-Depression Scale scores of 16 or higher – minor depression
DAS-MUNSHI, Jayati <i>et al.</i> , 2019 [16]	International Classification of Disorders-10 (ICD-10)
GALLO, Joseph J. <i>et al.</i> , 2005 [24]	Structured Clinical Interview for Axis I DSM-IV Disorders (SCID)
SULLIVAN, Mark D. <i>et al.</i> , 2007 [65]	Patient Health Questionnaire (PHQ-9)
BUTNORIENE, Jurate <i>et al.</i> , 2015 [8]	MiniNeuropsychiatric Interview (MINI)
HO, Cyrus SH <i>et al.</i> , 2016 [32]	Geriatric Mental State Examination (GMS)

Tabela 3.2: Resultado da questão de pesquisa referente ao diagnóstico de depressão de cada estudo selecionado.

Em relação a questão de pesquisa “Qual a idade média dos indivíduos que foram avaliados no estudo?”, foi observado que 7 estudos [56, 24, 65, 8, 32, 17, 60], a amostra da população tinha 50 anos ou mais de idade. Os outros estudos [71, 51, 62, 16], avaliaram amostras onde a idade da população apresentava 15 anos ou mais. A Tabela 3.3 apresenta o resumo sobre esta pergunta.

Autor (Referência)	Idade média da amostra
PENNINX, Brenda WJH <i>et al.</i> , 2001 [56]	50 anos ou mais
GALLO, Joseph J. <i>et al.</i> , 2005 [24]	50 anos ou mais
SULLIVAN, Mark D. <i>et al.</i> , 2007 [65]	50 anos ou mais
BUTNORIE, Jurate <i>et al.</i> , 2015 [8]	50 anos ou mais
HO, Cyrus SH <i>et al.</i> , 2016 [32]	50 anos ou mais
SAINT ONGE, Jarron M.; KRUEGER, Patrick M.; ROGERS, Richard G., 2014 [60]	50 anos ou mais
DINIZ, Breno S. <i>et al.</i> , 2014 [17]	50 anos ou mais
WU, C.S.; HSU, L.Y.; WANG, S.H., 2020 [71]	15 anos ou mais
MENG, Ruiwei, <i>et al.</i> , 2020 [51]	15 anos ou mais
STEWART, Ralph AH <i>et al.</i> , 2003 [62]	15 anos ou mais
DAS-MUNSHI, Jayati <i>et al.</i> , 2019 [16]	15 anos ou mais

Tabela 3.3: Resultado da questão de pesquisa referente a média de idade da amostra estudada.

Sobre a questão de pesquisa: “Os pacientes avaliados eram de uma comunidade não hospitalizada? Quais eram as condições da população do estudo em relação à interferência e ao tempo?”. Procurou observar em quais condições a amostra da população em estudo foi avaliada. Em relação à interferência um estudo pode ser observacional ou um ensaio clínico. No caso de um estudo observacional, o pesquisador não impõe um tratamento para os grupos de indivíduos que estão sendo avaliados, porém usa informações disponíveis sobre o indivíduo. Já um ensaio clínico, é um estudo experimental em que o pesquisador aloca indivíduos para um tratamento específico.

Sobre a população, os estudos longitudinais tinham aspectos semelhantes, pois os indivíduos eram acompanhados por um espaço de tempo. Nos estudos transversais, os pacientes estavam hospitalizados ou em fase de tratamento de alguma doença específica. Considerando o tipo de interferência de cada estudo, apenas um estudo apresentou como característica o ensaio clínico [62], os demais estudos caracterizaram-se como observacionais. Quanto ao tipos de pesquisas relacionado ao tempo, os estudos [71, 51, 65] eram estudos transversais (realizado em um determinado instante de tempo) e os demais estudos [56, 62, 16, 24, 8, 32, 17, 60] caracterizaram-se como estudos longitudinais (realizado ao longo do tempo). A Tabela 3.4 apresenta os resultado sobre os tipos de estudo. Todos os estudos apresentaram casos com depressão e casos sem depressão, ou seja, os grupos foram comparados durante os experimentos.

Autor (Referência)	Tipo de pesquisa relacionado ao tempo
WU, C.S.; HSU, L.Y.; WANG, S.H.,2020 [71]	Estudo transversal
MENG, Ruiwei, <i>et al.</i> , 2020 [51]	Estudo transversal
SULLIVAN, Mark D. <i>et al.</i> , 2007 [65]	Estudo transversal
PENNINX, Brenda WJH <i>et al.</i> , 2001 [56]	Estudo longitudinal
STEWART, Ralph AH <i>et al.</i> , 2003 [62]	Estudo longitudinal
DAS-MUNSHI, Jayati <i>et al.</i> , 2019 [16]	Estudo longitudinal
GALLO, Joseph J. <i>et al.</i> , 2005 [24]	Estudo longitudinal
BUTNORIE, Jurate <i>et al.</i> , 2015 [8]	Estudo longitudinal
HO, Cyrus SH <i>et al.</i> , 2016 [32]	Estudo longitudinal
DINIZ, Breno S. <i>et al.</i> , 2014 [17]	Estudo longitudinal
SAINT ONGE, Jarron M.; KRUEGER, Patrick M.; ROGERS, Richard G., 2014 [60]	Estudo longitudinal

Tabela 3.4: Resultado da questão de pesquisa referente ao tipo de pesquisa relacionada ao tempo.

Ao avaliar a pergunta “Quais métodos (estatístico ou computacionais) utilizados para relacionar a depressão e a causa de morte?”, verificou que apenas um estudo [24], usou como método estatístico o *Relative Odds* para identificar uma possível associação causal entre a depressão e a causa de morte. Os demais estudos utilizaram *Cox proportional hazards regression model*.

Sobre a pergunta “Os dados analisados são provenientes de exame de autópsia?”, não foi encontrado durante este Mapeamento estudos que faziam a relação de depressão e a causa de morte por meio de análises de autópsias.

Ao avaliar a pergunta principal: “Indivíduos com depressão morrem de causas diferentes quando comparados a indivíduos sem depressão?”, foi observado que as evidências dos estudos não permite estabelecer inferências causais ao associar a depressão e a causa da morte por todas as causas e por causa específica. Ainda, estudos revelam algumas controvérsias na literatura existente. Por exemplo, embora evidências anteriores tenham sugerido que a depressão pós-infarto agudo do miocárdio possa estar associada à redução da sobrevida, de acordo com o estudo de revisão de Machado et.al. [49], nenhuma evidência conclusiva indicou que o tratamento da depressão se traduz em um aumento da sobrevida nesta população específica. Em seguida é apresentado uma descrição como resposta à pergunta principal.

Na pesquisa de WU, C.S.; HSU, L.Y.; WANG, S.H.,2020 [71], analisaram dados da *Taiwan’s National Health Insurance Research Database* que correlacionaram

pacientes com diabetes que tinham transtornos depressivos e pacientes com diabetes sem depressão, dentre esses foram selecionados pacientes que desenvolveram complicações macrovasculares e microvasculares, mortalidade por todas as causas e mortalidade por causa específica. No entanto, não houve associação de depressão com complicações microvasculares, mortalidade por doenças cardiovasculares ou mortalidade por diabetes mellitus. O efeito da depressão nas complicações e mortalidade do diabetes foi mais proeminente entre os adultos jovens do que entre os adultos de meia-idade e idosos.

No estudo de MENG, Ruiwei *et al.*, 2020 [51], buscou investigar se a depressão é um fator de risco para a mortalidade por todas as causas e mortalidade por doenças cardiovasculares em adultos na China. Os resultados sugerem que a depressão é um fator de risco independente de mortalidade por doenças cardiovasculares e por todas as causas em adultos na China, especialmente em homens.

A pesquisa de PENNINX, Brenda WJHet *et al.*, 2001 [56], analisou uma amostra aleatória de idosos residentes na comunidade, estratificada por idade e sexo, que foi retirada dos registros populacionais de 11 municípios da Holanda. O estudo avaliou os efeitos da depressão maior e menor na mortalidade cardíaca. Foram examinados indivíduos com diagnóstico de doença cardíaca e indivíduos sem doença cardíaca. Não foram encontradas evidências de um efeito cardiovascular adverso mais forte da depressão em pacientes cardíacos quando comparados a indivíduos sem doença cardíaca. Em ambos os subgrupos, o risco excessivo de mortalidade cardíaca associado à depressão estava presente e muito semelhante. No entanto, os riscos de mortalidade cardíaca diferiam de acordo com o nível de depressão: os riscos eram cerca de duas vezes maiores para indivíduos com maior grau de depressão quando comparados a depressão mais leve.

Em STEWART, Ralph AH *et al.*, 2003 [62], foi realizado um estudo em pacientes com doença coronária estável, porém os sintomas depressivos eram mais comuns em pacientes com diagnóstico considerado médio a ruim. Foi observado que após o ajuste para fatores de risco cardiovascular, variáveis socioeconômicas e sintomas de doença cardiovascular, não houve associação entre sintomas depressivos e eventos cardiovasculares no caso de indivíduos fatais ou não fatais durante o seguimento de longo prazo.

No estudo de DAS-MUNSHI, Jayati *et al.*, 2019 [16], foram avaliados indivíduos com depressão em uma localidade etnicamente diversa no sudeste de Londres, acompanhados por 8 anos (2007-2014) vinculados a certidões de óbito. As taxas de mortalidades foram comparadas com a população geral (não deprimida) da Inglaterra e País de Gales. Os indivíduos foram padronizados por idade e sexo. Como resultado, apresentaram que pessoas com depressão experimentaram um risco aumentado de mortalidade por todas as causas, causas naturais e não naturais, em relação à população geral. Ao realizarem análises detalhadas por causa de morte, indicaram um grau de heterogeneidade por etnia e, em alguns casos, um menor risco de mortalidade por depressão para alguns grupos

étnicos, em comparação com a população de referência não deprimida. Estas análises indicaram que, indivíduos de grupos de minorias étnicas com depressão, apresentaram um risco reduzido de mortalidade por todas as causas em relação ao grupo branco-britânico com depressão. Isso foi especialmente marcado para indivíduos com depressão da África Negra, do Caribe Negro e do sul da Ásia, com tendências semelhantes observadas para mortes por causas naturais. Para mortes por causas não naturais, os negros caribenhos com depressão tiveram um risco de mortalidade mais baixo em relação aos brancos britânicos com depressão.

Na pesquisa de GALLO, Joseph J. *et al.*, 2005 [24], analisaram pacientes de 20 clínicas da Atenção Primária à Saúde das cidades de Nova York, Filadélfia e Pittsburgh. Os participantes foram identificados por meio de uma triagem para identificar a depressão de uma amostra aleatória. A amostra incluiu pacientes com triagem positiva para depressão e pacientes com triagem negativa para depressão. Neste estudo, a influência da depressão no risco de morte entre os pacientes mais velhos, após um intervalo de acompanhamento de 2 anos, contribuiu tanto para a mortalidade quanto para o infarto do miocárdio ou diabetes. Estimaram que a fração populacional atribuível de morte devido à depressão foi de 13%.

Neste estudo, SULLIVAN, Mark D. *et al.*, 2007 [65], avaliaram pacientes do *Action to Control Cardiovascular Risk in Diabetes (ACCORD)*. Examinaram a relação entre depressão, mortalidade e eventos cardiovasculares em uma amostra que recebia tratamento padronizado para diabetes. O estudo também incluiu critérios rigorosos para definir complicações macrovasculares e microvasculares e a causa da morte. Além disso, examinaram o efeito da depressão na mortalidade entre aqueles com e sem doença cardiovascular prévia. Como resultado, indicaram que a depressão aumenta o risco de mortalidade por todas as causas e pode aumentar o risco de eventos macrovasculares entre adultos com diabetes tipo 2 com alto risco de eventos cardiovasculares.

Em BUTNORIENE, Jurate *et al.*, 2015 [8], pacientes foram selecionados aleatoriamente do Centro de Atenção Primária à Saúde e avaliados prospectivamente pela presença de Síndrome Metabólica (SM), fatores de riscos para doenças cardiovasculares, depressão e ansiedade. Os achados deste estudo demonstraram que a SM e a ansiedade previram maior mortalidade cardiovascular em mulheres de meia-idade, assim a SM e ansiedade foram considerados preditores de sobrevida neste grupo. Este efeito, quando analisados em homens, não obteve o mesmo resultado, sendo que, ao avaliar depressão, ansiedade e SM, não houve associação à mortalidade.

No estudo de HO, Cyrus SH *et al.*, 2016 [32], avaliaram dados coletados do *National Pesquisa de Saúde Mental de Idosos (NMHS-E)* de Cingapura. A amostra tinha pessoas com e sem depressão. Os resultados desta pesquisa indicaram que depressões principais e subliminares estão associadas ao aumento da mortalidade, em grande parte

devido ao estilo de vida e comportamentos de saúde, comorbidade médica e incapacidade funcional. No entanto, ao avaliarem uma subamostra de 50 indivíduos, a depressão maior foi associada a um impacto independente sobre a mortalidade, especialmente mortes resultantes de doenças cardiovasculares e acidente vascular cerebral.

Nesta pesquisa DINIZ, Breno S.*et al.*, 2014 [17], avaliaram especificamente a depressão tardia (indivíduos com 60 anos ou mais). A análise utilizou dados do Estudo de Coorte de Envelhecimento de Bambuí, Brasil. Este estudo concluiu que a depressão tardia está associada ao aumento do risco de mortalidade por todas as causas em idosos, em particular naqueles com sintomas depressivos mais graves. Além disso, o gênero pode moderar o risco de mortalidade, sendo mais elevado em homens do que em mulheres.

Os dados da pesquisa de SAINT ONGE, Jarron M.; KRUEGER, Patrick M.; ROGERS, Richard G., 2014 [60], são provenientes do *National Health Interview Survey*, no qual a amostra é referente à adultos não institucionalizados nos Estados Unidos. O estudo examinou uma amostra com mortalidade específica por causa (cardiovascular e câncer). A depressão maior entre adultos mais velhos foi associada a um risco aumentado de mortalidade não suicida durante um acompanhamento de 6 anos. Como resultado, a depressão foi associada a 2,68 vezes ao risco de mortalidade por doença cardiovascular entre aqueles que não tinham doença cardiovascular preexistente, em comparação com 1,82 vezes para aqueles com doença cardiovascular.

3.4 Considerações Finais

Este capítulo descreveu o processo de Mapeamento Sistemático, cuja as fases foram planejamento, condução e publicação dos resultados. Não foi encontrado nenhum estudo que fizesse uma relação entre mortalidade e a depressão buscando associação entre todas as causas de morte por meio da autópsia.

De acordo com os resultados encontrados, parte das pesquisas se concentraram em buscar associações entre a depressão e doenças específicas, como por exemplo, doenças cardiovasculares e diabetes. Ainda, outros estudos mostraram a associação de depressão com risco aumentado de morte por todas as causas e causas específicas em populações em geral e em grupos específicos de pacientes. No entanto, a relação causal entre depressão e as causas de mortes permanecem inconclusivas. [14, 49, 52]. Embora, a relação entre e depressão e causa de morte não esteja clara, é possível perceber na literatura que a depressão pode afetar consideravelmente a qualidade de vida e o bem estar. No entanto, é aceito que a maioria das pesquisas apontam que a depressão está relacionada ao aumento da mortalidade [14], dessa maneira estudos prospectivos, em que a aferição da causa da morte é mais precisa (por meio de autópsia) e ainda que a amostra

da população seja da comunidade, devem ser realizados para entender se há uma relação entre depressão e causas de morte.

Materiais e Métodos

Este capítulo apresenta os materiais e métodos referente a esta pesquisa. A Seção 4.1 apresenta os dados utilizados neste trabalho. A Seção 4.2 aborda os métodos para a execução da proposta e, a Seção 4.3, as ferramentas necessárias para o desenvolvimento desta pesquisa.

4.1 Conjuntos de dados

Os conjuntos de dados usados neste estudo são provenientes de uma parceira com pesquisadores do Grupo de Estudos do Envelhecimento do Cérebro da Faculdade de Medicina da Universidade de São Paulo (FMUSP). Todos os protocolos da FMUSP, o termo de consentimento e os procedimentos seguem os regulamentos internacionais e brasileiros para pesquisas envolvendo seres humanos e foram aprovados pelos comitês de pesquisas locais e federais. Os indivíduos foram incluídos após os procedimentos do estudo terem sido explicados aos informantes e eles concordaram em participar assinando o termo de consentimento. O informante era um familiar próximo ou cuidador que teve pelo menos um contato semanal com o falecido nos últimos 6 meses antes da morte e foi capaz de recontar e fornecer detalhes sobre a saúde do falecido. Neste trabalho, não foi realizada a coleta de dados. Os conjuntos de dados desta pesquisa já estavam tabelados e foram usados em estudos anteriores do Grupo de Estudos do Envelhecimento do Cérebro da FMUSP. Além disso, os procedimentos metodológicos foram descritos em outros estudos que exploraram estes conjuntos de dados [29, 63, 23, 64]. Estes dados estão separados em dois arquivos, um contendo as informações dos laudos de autópsia e outro com informações de uma entrevista clínica.

Conjunto de dados da autópsia

Os dados de autópsia são da região metropolitana da cidade de São Paulo, consiste em uma amostra de indivíduos não hospitalizados (comunidade) e foram coletados pelo Serviço de Verificação de Óbitos da Capital da Universidade de São Paulo (SVOC-

USP), entre os anos de 2004 a 2016. Estes dados consistem em uma ampliação dos dados utilizados em um estudo anterior realizado por Suemoto et al. [63]. Todos os indivíduos desta amostra tinham 50 anos ou mais.

O conjunto de dados de autópsia contém informações relacionadas às causas de morte de um indivíduo. Este conjunto é composto pelos atributos “nsvo”, “Causa do óbito”, “Doença principal”, “Causa básica1”, “Causa básica2”, “Outras doenças” e “Diagnósticos”, sendo que “Outras doenças” e “Diagnósticos” possuem diversos outros campos que não foram utilizados neste trabalho. O atributo “nsvo” é a identificação do indivíduo que faleceu, e é composta por 3 a 5 dígitos, uma barra e mais 2 dígitos que são referentes ao ano do óbito, como por exemplo: “3941/09”¹. A “Causa do óbito” se refere à causa imediata ou terminal da morte; a “Doença principal”, é a doença que contribuiu diretamente para o óbito, e “Causa básica 1” e “Causa básica 2” são as causas intermediárias que contribuíram para o óbito. Estes atributos estão codificados com o código da CID-10 das respectivas doenças.

Na Tabela 4.1 são apresentados três exemplos de registros do conjunto de dados de autópsia. A variável *cid_obito*, que representa a “Causa do óbito”, é a única que deve estar preenchida em todas as observações. No entanto, as demais variáveis relativas às outras doenças associadas à morte podem ter valores faltantes. No primeiro registro desta tabela de exemplo, temos que o indivíduo cujo identificador *nsvo* é “9287/08” teve como causa do óbito um “insuficiência cardíaca congestiva” (CID I50), no entanto, apresentou como doença principal “insuficiência ventricular esquerda” (CID I51). Além disso, este indivíduo possuía mais outras duas doenças que foram associadas ao óbito: “doença cardíaca hipertensiva” (CID I11) e “hipertensão arterial sistêmica” (CID I10). Pode-se notar que uma mesma doença pode ser a causa do óbito, doença principal ou doença básica em diferentes indivíduos. Por exemplo, “hipertensão arterial sistêmica” (CID I10) em negrito na Tabela 4.1 é a “Causa básica2” para o primeiro indivíduo, mas é a “Causa do óbito” (causa da morte) para o segundo indivíduo (*nsvo* “8340/16”), e “Doença principal” para o terceiro indivíduo (*nsvo* “0901/07”).

Em relação as doenças relacionadas a morte, considerou-se que não há erros no preenchimento das CIDs neste conjunto de dados, uma vez que o registro das CIDs no laudo de autópsia é realizado por um médico patologista (que faz o exame macroscópico e microscópico do cadáver), e que apenas casos onde não se sabe a causa da morte passam por autópsia no SVOC-USP.

¹Os dados referentes ao atributo “nsvo” usados neste trabalho como exemplos são fictícios para manter a confidencialidade dos indivíduos

nsvo	cid_obito	cid_principal	cid_basica1	cid_basica2
9287/08	I50	I51	I11	I10
8340/16	I10			
0901/07	J11	I10		

Tabela 4.1: Exemplos de casos do conjunto de dados de autópsia

Conjunto de dados da entrevista clínica

Os dados clínicos são provenientes do *Biobanco para Estudos em Envelhecimento* da FMUSP, e foram coletados por meio de uma entrevista completa com um familiar próximo (enquanto aguardava a autópsia). Foram coletados dados sócio-econômicos, demográficos, e clínicos, identificados pela variável *nsvo*, tais como idade ao morrer, sexo, anos de escolaridade, etnia (branca, preta, parda e asiática), frequência de contato entre o falecido e o informante, diagnóstico médico prévio de hipertensão, diabetes, doença arterial coronariana, insuficiência cardíaca, dislipidemia, acidente vascular cerebral, tabagismo e uso de álcool e avaliação cognitiva por meio do *Clinical Dementia Rating (CDR)* [64].

Além dos dados socio-demográficos, este conjunto também possui informações da Entrevista Clínica Estruturada para o DSM-IV (SCID, do inglês, *Structured Clinical Interview for DSM-IV*). O SCID é uma entrevista semiestruturada criada para fazer diagnósticos psiquiátricos confiáveis em adultos de acordo com o DSM-IV [40]. O questionário aplicado ao informante do falecido para obter estes dados é apresentado na Figura 4.1.

<p>HUMOR DEPRIMIDO: 1) Já houve, pelo menos uma vez na vida dele(a), algum momento em que se sentiu deprimido(a), ou na pior na maior parte do dia, quase todos os dias? (Como é que foi isso?) Quando foi isso? _____</p> <p>ANEDONIA: Naquele tempo ou em outro momento, ele(a) perdeu o interesse ou o prazer pelas coisas que costumava gostar? (Como é que foi isso?) Isso acontecia quase todos os dias? Quando foi isso? _____</p>
<p>3) Naquele período ele(a) perdeu ou ganhou peso? Ele(a) estava tentando perder peso? SE NÃO: Como estava o apetite dele(a)? (Tinha que forçar para comer? Comia menos/mais que o habitual? Isso acontecia quase todos os dias?)</p>
<p>4) Como estava o sono dele(a)? (Dificuldade para dormir, acordando frequentemente, dificuldades em ficar acordado, acordando muito cedo, ou, dormindo demais? Quantas horas por noite comparado com o habitual? Isso acontecia quase todas as noites?)</p>
<p>5) Ele(a) estava tão inquieto(a) ou agitado(a) que não era capaz de ficar parado(a)? (Isso acontecia todos os dias?) SE NÃO: E o contrário: falando e se movendo mais devagar do que o normal dele(a)? (Isso acontecia todos os dias?)</p>
<p>6) Como estava a energia dele(a)? (Cansado(a) o tempo inteiro? Quase todos os dias?)</p>
<p>7) Como ele se sentia em relação a si mesmo? (Inútil, sem valor?) Quase todos os dias? SE NÃO: Ele se sentia culpado por coisas que fazia ou deixava de fazer? Quase todos os dias?</p>
<p>8) Ele(a) tinha dificuldade para se concentrar ou pensar? (Com que tipo de coisas isso interferia?) Quase todos os dias? SE NÃO: Era difícil para ele(a) tomar decisões sobre coisas do dia a dia? Quase todos os dias?</p>
<p>9) As coisas estavam tão ruins que ele(a) pensava muito em morte ou que estaria em melhor situação se estivesse morto(a)? E quanto a se ferir?</p>
<p>10) Durante aquele tempo ficou difícil para ele(a) trabalhar, cuidar das coisas em casa, ou se relacionar com as pessoas?</p>
<p>11) Pouco antes de tudo isso começar ele estava bebendo em excesso ou usando drogas?</p>
<p>12) Pouco antes disso começar ele estava fisicamente doente?</p>
<p>13) Isso começou logo depois de alguém próximo a ele(a) morrer? Durou mais que 2 meses?</p>

Figura 4.1: Formulário aplicado ao informante do falecido para a coleta de dados para o diagnóstico de episódio depressivo no passado

Para cada pergunta do questionário, é possível uma das seguintes respostas: 1 - ausente ou falso, 2 - subliminar (duvidoso), 3 - limiar ou verdadeiro, e “?” - informação inadequada. O diagnóstico da depressão foi realizado conforme a avaliação destas respostas. Os critérios para classificar como depressão é quando cinco (ou mais) respostas são verdadeiras (valor=3), sendo que dentre as cinco, a primeira (humor deprimido) e/ou a segunda (anedonia - perda do interesse ou prazer) deve ser verdadeira. Além das 13 questões do formulário apresentado na Figura 4.1, existe também uma questão sobre a ocorrência de outros episódios depressivos do indivíduo.

4.2 Método

Para a realização desta pesquisa, a metodologia aplicada foi a do processo de Descoberta de Conhecimento em Bases de Dados (KDD, do inglês, *Knowledge Discovery in Databases*) [20]. O KDD é um processo não trivial de identificação de padrões que é interativo e iterativo no qual envolve diversas fases com várias decisões sendo feitas por um analista, que é um especialista do domínio dos dados. As fases do KDD podem

ser generalizadas em: compreensão do domínio, preparação do conjunto de dados (que envolve as etapas de seleção, limpeza ou pré-processamento, e transformação dos dados), mineração de dados, interpretação e avaliação de resultados, e finalmente a obtenção do conhecimento. A adaptação das fases do KDD para o presente projeto é apresentada na Figura 4.2, e a descrição dos processos relacionados ao KDD são apresentados nas próximas seções deste capítulo. A análise dos resultados é apresentada no Capítulo 5.

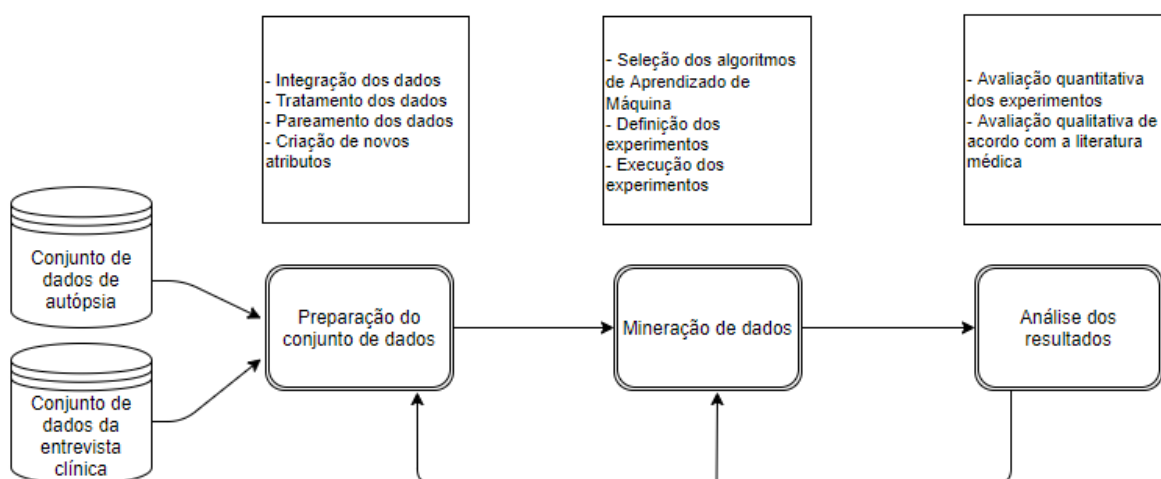


Figura 4.2: Metodologia referente ao processo de descoberta de conhecimento utilizada neste estudo.

4.2.1 Preparação do conjunto de dados

Os dados originais estavam em arquivos separados conforme foram coletados, mas é possível relacionar estas informações através do código de identificação (“nsvo”) do indivíduo. Vale ressaltar que, não se teve acesso a nenhuma outra informação do indivíduo que não fosse os dados clínicos e os dados de morte. Todos os indivíduos foram anonimizados, sendo a única referência de identificação destes indivíduos a variável “nsvo”. A Tabela 4.2 apresenta o número de variáveis e de observações dos conjuntos de dados de autópsia e da entrevista clínica. A descrição das variáveis utilizadas neste estudo estão no Apêndice A. Existe uma diferença entre o número de observações destes conjuntos, o motivo desta diferença é que os familiares entrevistados de alguns indivíduos não souberam responder às perguntas da entrevista clínica.

Conjunto de dados	Número de observações	Número de variáveis
Laudo de autópsia	2189	11
Entrevista clínica	1136	54

Tabela 4.2: Número de observações e de variáveis nos conjuntos de dados originais

Primeiramente, foi realizada a integração dos conjuntos de dados de autópsia e entrevista clínica em um único conjunto por meio do identificador “nsvo”. Em seguida, foram realizadas tarefas relacionadas ao tratamento da qualidade dos dados. As operações realizadas foram:

- verificação de registros repetidos;
- a inclusão das CID's faltantes nos casos onde havia a descrição da doença, mas não o seu código;
- exclusão de registros:
 - 17 registros foram excluídos pois não tinham o registro de autópsia;
 - 9 registros foram excluídos devido a causa do óbito ser por doença inespecífica;
 - 3 registros foram excluídos pois com os dados da entrevista clínica não foi possível verificar a depressão ou não;
 - 1 registro foi excluído por divergência de idade entre o conjunto de dados da entrevista clínica e de autópsia;
 - 1 registro foi excluído, pois a causa do óbito apresentava como motivo de morte uma informação inconsistente, não sendo possível verificar qual a causa correta da morte.
 - 3 registros foram excluídos, pois houve divergência da informação do sexo entre os conjuntos de dados de autópsia e da entrevista clínica e não foi possível avaliar qual informação estava correta.

A Figura 4.3 é apresentado um resumo sobre a exclusão de registros.

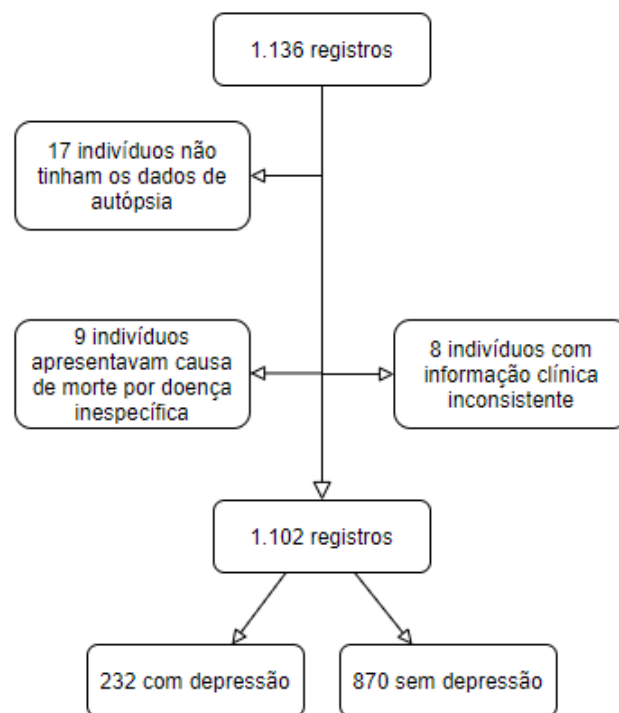


Figura 4.3: Operações realizadas para o tratamento da qualidade da base de dados

Na Tabela 4.3 é mostrado o quantitativo dos dados após a integração e tratamento dos dados.

Total observações	Número atributos	Classe “Depressão”	Classe “Sem depressão”
1102	65	232 (21,05%)	870 (78,95%)

Tabela 4.3: Resumo quantitativo dos dados deste estudo após a integração do conjunto de dados e ao tratamento da qualidade dos dados.

Pareamento dos dados

O pareamento refere-se à seleção de indivíduos de maneira que as distribuições dos fatores considerados entre casos e controles sejam similares [47]. Neste contexto, o pareamento pondera a amostra do grupo “Depressão” a fim de aumentar a semelhança com os indivíduos do grupo “Sem depressão”, visando balancear as características observadas de ambos os grupos da amostra. Além disso o processo de pareamento gera um elevado grau de balanceamento das variáveis entre a amostra do grupo “Depressão” e do grupo “Sem depressão” que foram utilizadas no processo de classificação ao utilizar os algoritmos de AM.

Dessa maneira, a composição de casos e controles da base de dados procedeu da comparação de variáveis, cujos valores foram comparados por pares de registros.

O pareamento levou a um número de amostras de casos igual ao número de controle, onde cada caso de depressão corresponde a um caso sem depressão com as mesmas características. As variáveis usadas para fazer o pareamento foram:

1. idade
2. idade \pm 4 anos (ao comparar utilizando esta variável, a idade correspondente poderia ser de até 4 anos de diferença)
3. sexo biológico
4. demência
5. escolaridade
6. etnia

Foi desenvolvido um algoritmo para realizar o pareamento entre as classes “Depressão” e “Sem depressão”. Inicialmente, o algoritmo buscou na base de dados os registros da classe “Depressão” e os comparou aos registros da classe “Sem depressão”, onde as variáveis idade, sexo biológico, demência, escolaridade e etnia eram iguais. Como não foi possível obter todos os registros pareados entre as classes quanto às variáveis citadas, uma nova seleção foi realizada. Porém, os registros já pareados foram excluídos antes de iniciar outra seleção. Dessa maneira, iniciou-se uma nova busca, neste caso considerando as variáveis: idade \pm 4 anos, sexo biológico, demência, escolaridade e etnia. E então os registros pareados com estas variáveis foram excluídos, pois uma nova busca foi necessária, uma vez que ainda não havia pareado todos os registros da classe “Depressão”. O algoritmo de pareamento realizou esse processo de comparação entre as variáveis e a exclusão dos registros já pareados por mais 3 vezes, seguindo os critérios de comparação entre idade \pm 4 anos, sexo biológico, demência e escolaridade. Em seguida comparou idade \pm 4 anos, sexo biológico e demência e por fim, idade \pm 4 anos e sexo biológico. Dessa maneira, não foi descartado nenhum registro da classe “Depressão”, pois todos os registros foram pareados de acordo com as buscas citadas anteriormente. A Figura 4.4 mostra o fluxograma deste algoritmo. Após a etapa de pareamento, obteve-se uma base de dados balanceada entre as classes “Depressão” e “Sem depressão”. Estes registros pareados foram utilizados em todos os experimentos deste estudo.

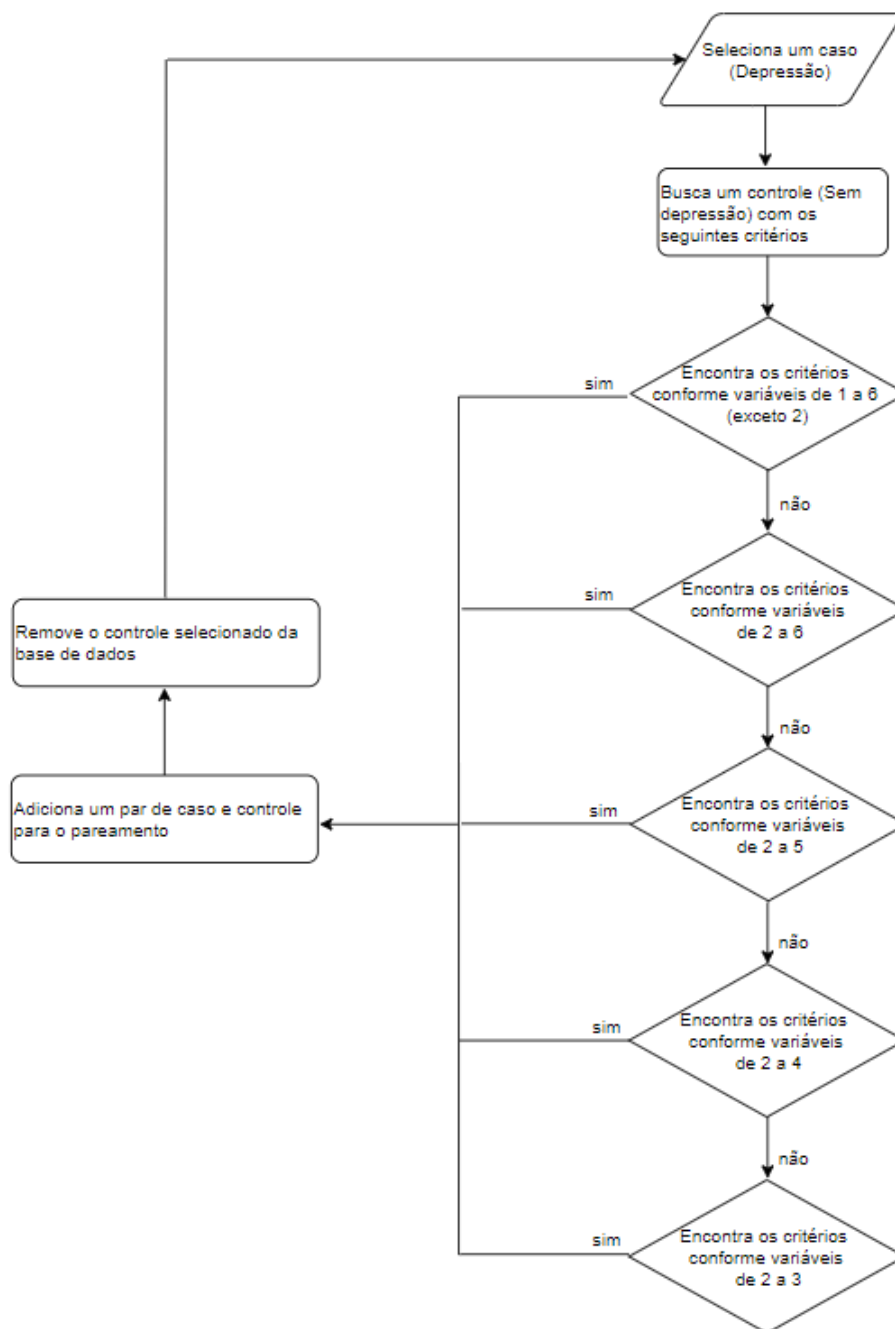


Figura 4.4: Fluxograma do algoritmo de pareamento

Criação de atributos

A criação de atributos é uma das tarefas da fase de transformação do KDD que permite o enriquecimento semântico das informações. Neste estudo, foram criados novos atributos a partir dos capítulos e do agrupamento de categorias da CID-10. As categorias de três dígitos da CID são excessivamente numerosas para que se possa incluí-las como variável categórica na análise, assim, formas alternativas de agrupamento de categorias de

três dígitos se fazem necessárias para melhorar o tratamento da informação diagnóstica da CID. Desta forma, foram criados dois conjuntos de agrupamentos da CID-10 que foram denominados: “Capítulos da CID-10” e “Categorias da CID-10”, descritos abaixo:

1. “Capítulos da CID-10” (Tabela 4.4) corresponde ao agrupamento dos grandes capítulos da CID. Esse agrupamento tem como meta principal manter a coerência clínica dentro de um mesmo capítulo da CID.

Grupos de doenças	Código CID
Neoplasmas (tumores)	C
Doenças do sangue e dos órgãos hematopoéticos e alguns transtornos imunitários	D
Doenças endócrinas, nutricionais e metabólicas	E
Doenças do aparelho circulatório	I
Doenças do aparelho respiratório	J
Doenças do aparelho digestivo	K
Doenças do aparelho geniturinário	N
Outras causas/Doenças infecciosas*	O

Tabela 4.4: Oito grupos da CID utilizados para análise da causa de morte na depressão.

*Outras causas/Doenças infecciosas agrupa doenças como: *tuberculosis, sepsis, Chagas' disease, unknown cause of mortality, other disorders of brain*, que correspondem respectivamente aos CID's A15, A41, B57, R99, G93.

2. “Categorias da CID-10” (Tabela 4.5) corresponde ao agrupamento de categorias de três dígitos, respeitando a numeração dos capítulos da CID-10. Esse agrupamento foi feito de forma empírica tendo como meta principal manter a coerência clínica e diminuir o número de variáveis.

Grupos de doenças	Grupos de categorias de CIDs
câncer do trato respiratório	C33, C34
tumor primário desconhecido	C80
câncer pancreático, hepático e colestático	C22, C24, C25
câncer geniturinário	C50, C53, C61, C67, R19
metástase	C26, C38, C41, C56, C64, C76, C77, C78, C79
câncer de cabeça e pescoço	C00, C01, C02, C10, C12, C14, C32
câncer do trato digestivo	C15, C16, C17, C18, C20
anemia	D50, D53, D57, D64, D68
diabetes mellitus	E14
desnutrição protéico-calórica, caquexia	E43, E44, E46, R64
doenças relacionadas com hipertensão arterial	I10, I11, I12
doença isquêmica do coração	I20, I21, I22, I23, I24, I25
aterosclerose	I70
aneurisma aórtico	I71
pericardite	I30, I31
doenças pulmonares isquêmicas relacionadas	I26, I27
Outros cardioI: hemorragia das vias respiratórias, hemorragia, não classificada em outra parte	R04, R58
Outros cardioII: cardiomiopatia, choque cardiogênico, insuficiência cardíaca	I42, I50, R57
infecções do trato respiratório inferior, Pitórax	J06, J18, J20, J21, J34, J41, J69, J85, J86
inalação de conteúdo gástrico	W78
doença pulmonar obstrutiva crônica	J42, J43, J44
doenças hepáticas secundárias ao uso de álcool	K70, K71, K72, K73, K74, K76
doenças colestáticas	K80, K81, K83
doença peritoneal, peritonite	K66, K65
outras doenças digestivas: úlcera duodenal, outras doenças do aparelho digestivo, distúrbio vascular do intestino, hemorragia gastrointestinal, íleo paralítico	K26, K92, K55, K56, K85, K86
doença renal crônica e aguda	N03, N08, N12, N15, N18, N25, N35, N17
infecções do trato urinário	N10, N11, N30, N39

Tabela 4.5: 27 grupos da CID, que agrupam doenças similares, utilizados para análise da causa de morte na depressão.

Cada agrupamento de CIDs apresentado acima se tornou um atributo binário. Estes novos atributos foram preenchidos com valores *sim* ou *não* conforme a codificação da causa do óbito e das doenças relacionadas ao óbito, ou seja, se o código CID da causa do óbito ou de alguma das doenças relacionadas (principal, básica 1 ou básica 2), se enquadrar em algum dos grupos das tabelas 4.4 ou 4.5, o atributo correspondente recebe o valor “sim”, caso contrário, recebe o valor “não”. Por exemplo, o indivíduo representado na Tabela 4.6 teria os atributos referentes aos capítulos “Doenças do aparelho digestivo” (CID K) e “doenças do aparelho circulatório” (CID I), e os grupos de categorias “doença peritoneal, peritonite” (CID K66), “aneurisma aórtico” (CID I71) e “aterosclerose”(CID 70), marcados com *sim*, enquanto que para todos os outros atributos seria atribuído o valor *não*. O registro completo deste indivíduo pode ser visualizado na Tabela 4.7.

nsvo	cid_obito	cid_principal	cid_basica1	cid_basica2
6047/04	K66	I71	I70	

Tabela 4.6: Exemplo de um registro de autópsia de um indivíduo

Atributo	Valor	Atributo	Valor
nsvo	6047/04	metástase	não
depressão	sim	câncer de cabeça e pescoço	não
sexo biológico	masculino	câncer do trato digestivo	não
idade	65	anemia	não
demência	sim	diabetes mellitus	não
escolaridade	fundamental	proteína calórica maln, caquexia	não
etnia	caucasiano	doenças relacionadas a hipertensão	não
cid_obito	K66	doenças isquêmica do coração	não
cid_principal	I71	sistema aterosclerose	sim
cid_basica1	I70	aneurisma aórtico	sim
cid_basica2		pericardite	não
cancer	não	embolia pulmonar	não
tumor benigno e células sanguíneas	não	outros cardioI	não
nutrição do sistema endócrino	não	outros cardioII	não
cardiovascular	sim	infecção do trato respiratório baixo	não
doença do trato respiratório	não	inalação do conteúdo gástrico	não
doença do sistema digestivo	sim	crônica e fibrose	não
sistema geniturinário	não	doenças hepáticas secundárias ao uso de álcool	não
outras causas	não	doenças colestáticas	não
câncer trato respiratório	não	peritoneal e peritonite	não
tumor primário desconhecido	não	outras doenças digestiva	não
câncer de pâncreas, fígado e colestático	não	rim crônico e agudo	não
câncer geniturinário	não	doença do trato urinário	não

Tabela 4.7: Exemplo de um registro com os novos atributos preenchidos a partir dos CID's das causas do óbito.

4.2.2 Mineração de Dados

Nesta etapa do KDD são selecionados os métodos utilizados para localizar padrões nos dados. Foram selecionados e aplicados algoritmos específicos de aprendizado de máquina supervisionados, ou seja, algoritmos cujo propósito é gerar modelos de classificação dos indivíduos nas classes “Depressão” ou “Sem depressão”. Para atender ao objetivo de se explorar diferentes classificadores e avaliar o desempenho preditivo, 11 algoritmos de AM foram usados nos experimentos computacionais. Estes algoritmos foram escolhidos baseados em estudos recentes em periódicos com alto fator de impacto que utilizaram dados médicos [26, 45, 39, 13, 25, 33]. Assim, foram incluídos nos experimentos os seguintes algoritmos: *Logistic Regression (LR)*, *Decision Tree (DT)*, *Naïve Bayes (NB)*, *Support Vector Machine (SVM)*, *K-Nearest-Neighbor (KNN)*, *Random Forest (RF)*, *Multilayer Perceptron (MLP)*, *AdaBoost (AD)*, *Gradient Boosting (GB)*, *Extreme Gradient Boosting (XGBoost)*, *Light Gradient Boosting Machine (LGBM)*. Estes algoritmos pertencem a diferentes paradigmas de AM e conferem uma diversidade de abordagens para tentar encontrar padrões nos dados. No Apêndice B estes algoritmos são descritos resumidamente de acordo com o viés de aprendizagem que assumem.

4.2.3 Avaliação

A tarefa dos algoritmos de aprendizado de máquina deste trabalho é a classificação binária dos indivíduos da base de dados. A classificação binária se dá quanto os casos são rotulados em 2 classes, usualmente denominadas “positiva” e “negativa”. No presente estudo, a classe “Depressão” foi considerada como a classe positiva enquanto que “Sem depressão” foi considerada como classe negativa. Assim, um atributo da base de dados que representa a variável de classe possui os valores “sim” ou “não” conforme o caso.

Para aferir o desempenho dos modelos gerados, foi considerada a acurácia, dada pela Equação 4-1. A acurácia avalia a taxa de acerto do classificador em ambas as classes, ou seja, ela avalia a proporção entre as instâncias previstas corretamente e todas as instâncias no conjunto de dados. No presente estudo, os dados foram pareados e conseqüentemente ficaram balanceados, isto é, com a mesma quantidade de casos positivos e negativos. Logo, a acurácia é considerada uma métrica confiável para a avaliação [6, 13].

$$Acuracia = \frac{VP + VN}{VP + VN + FP + FN} \quad (4-1)$$

onde:

- VP são verdadeiros positivos, quantidade de casos positivos corretamente classificados como positivos;

- *VN* são verdadeiros negativos, quantidade de casos negativos corretamente classificados como negativos;
- *FP* são falsos positivos, quantidade de casos negativos classificados incorretamente como positivos;
- *FN* são falsos negativos, quantidade de casos positivos classificados incorretamente como negativos.

Para avaliar a capacidade de generalização (ou seja, a capacidade de classificar corretamente novos casos) dos modelos gerados pelos algoritmos de AM, foi utilizada a técnica de validação cruzada estratificada com *k folds* (*stratified k-fold cross validation*) [38], na qual o conjunto total de dados é dividido em *k* subconjuntos de mesmo tamanho e com a mesma proporção de exemplos de cada classe. A partir disso, um subconjunto é usado para avaliação (teste), e o restante para a geração dos modelos (treinamento). Esse processo é repetido *k* vezes e, em cada iteração, um subconjunto diferente é usado para teste enquanto que os outros são unidos para treinamento. No presente trabalho, o valor de *k* utilizado foi 10 conforme recomendado em [38]. Além dos subconjuntos da validação cruzada estarem estratificados, os exemplos destes subconjuntos também se encontravam pareados conforme o procedimento descrito na Seção 4.2.1.

4.3 Ferramentas

Para o desenvolvimento dos experimentos, foi utilizado um computador com as seguintes configurações: processador Intel core i5, 6GB de memória RAM, sistema operacional Windows 10 de 64 bits. A linguagem de programação utilizada foi a linguagem Python 3.6. Foi utilizado o Google *Collaboratory*, conhecido como Colab, que é um serviço em nuvem gratuito hospedado pelo Google. Por ser executado nos servidores do Google não é utilizado todo o recurso computacional local, uma vez que recursos computacionais são reservados na nuvem em tempo de execução. O Colab permite aos desenvolvedores criar e executar códigos e escrever textos em um único documento interativo que são agrupados por células, chamado de notebooks. Estes notebooks são salvos na conta do Google Drive do executor. Algumas vantagens do Colab são aceleração de GPU sem custo e bibliotecas pré-instaladas, como Scikit-learn e Matplotlib, entre muitas outras.

As bibliotecas que foram utilizadas para a execução dos experimentos realizados neste trabalho foram:

- Bibliotecas para análise de dados

- *NumPy*: responsável pelo processamento de grandes matrizes e matrizes multidimensionais, além disso apresenta uma extensa coleção de funções matemáticas de alto nível e a execução de várias operações com esses objetos;
- *Pandas*: fornece ferramentas de análise de dados e estruturas de dados de alta performance;
- Bibliotecas para visualização:
 - *Matplotlib*: é uma biblioteca de baixo nível para criar diagramas e gráficos bidimensionais;
 - *Seaborn*: API de alto nível baseada na biblioteca matplotlib, contém configurações padrão mais adequadas para o processamento de gráficos.
- Bibliotecas para algoritmos de aprendizado de máquina
 - *Scikit-learn*: fornece algoritmos para muitas tarefas padrão de aprendizado de máquina e mineração de dados.

Resultados e Discussão

Neste capítulo são apresentados os resultados dos experimentos realizados e a discussão de acordo com a literatura médica. Na Seção 5.1 são apresentados os resultados na Seção 5.4 a discussão.

5.1 Experimentos

Os experimentos foram elaborados para analisar diferentes perspectivas das causas de morte. Foram realizados 3 experimentos distintos, nos quais os dados de entrada variaram conforme os agrupamentos das doenças apresentados na Seção 4.2.1. Inicialmente, foram analisadas as doenças de forma mais ampla, agrupadas conforme os capítulos da CID (Tabela 4.4), em seguida grupos mais específicos de doenças (Tabela 4.5) e por último os dois agrupamentos juntos (Tabela 4.4 e 4.5). Em todos os experimentos, as análises foram feitas com os 11 algoritmos de AM selecionados. Os parâmetros destes algoritmos foram ajustados com os valores padrões previamente definidos pelo ambiente de programação utilizado. A descrição dos algoritmos de classificação encontra-se no Apêndice B. Adicionalmente, foram realizados experimentos com ajustes automáticos de parâmetros. A técnica utilizada foi a *Random Search*, que testa combinações aleatórias de hiperparâmetros. Como os resultados foram muito semelhantes quando comparados aos valores dos parâmetros padrão, optou-se por considerar os experimentos com este último, uma vez que o tempo de execução dos experimentos é muito menor.

Para cada experimento, foi selecionado um conjunto de atributos relacionados aos agrupamentos da CID-10, como apresentado na Subseção 4.2.1. Dessa maneira, os experimentos tiveram o seguinte formato:

- **Experimento 1:** corresponde aos atributos baseados nos capítulos da CID-10 (Tabela 4.4).
- **Experimento 2:** corresponde os atributos baseados nas categorias da CID-10 (Tabela 4.5).

- **Experimento 3:** considera os atributos baseados tanto nos capítulos quanto nas categorias da CID-10 (Tabelas 4.4 e 4.5).

Na Figura 5.1 é apresentado um esquema de organização destes experimentos. A numeração dos experimentos corresponde aos atributos selecionados para cada análise.

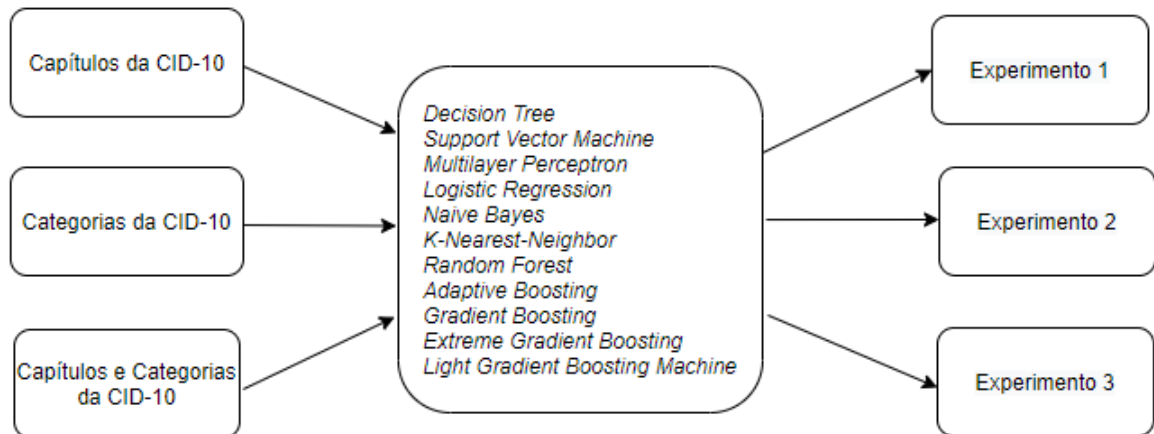


Figura 5.1: Diagrama dos experimentos. Cada experimento utiliza um conjunto diferente de atributos como entrada para os 11 algoritmos de AM selecionados.

5.2 Causas de morte na depressão

Na Tabela 5.1 a seguir, estão listados todos os resultados das acurácias individuais de cada algoritmo para os 3 experimentos. No geral, embora este estudo tenha explorado diversos algoritmos de AM com diferentes vieses de aprendizado, todos os resultados tiveram acurácias semelhantes, variando de 43 % a 52 %, independentemente do conjunto de variáveis utilizadas como entrada.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,51 (0,04)	0,47 (0,06)	0,47 (0,07)
<i>Support Vector Machine</i>	0,50 (0,04)	0,46 (0,04)	0,45 (0,04)
<i>Multilayer Perceptron</i>	0,50 (0,03)	0,49 (0,07)	0,47 (0,05)
<i>Logistic Regression</i>	0,51 (0,03)	0,48 (0,03)	0,48 (0,04)
<i>Naïve Bayes</i>	0,51 (0,03)	0,50 (0,03)	0,49 (0,04)
<i>K-Nearest-Neighbor</i>	0,52 (0,05)	0,52 (0,08)	0,48 (0,05)
<i>Random Forest</i>	0,50 (0,04)	0,49 (0,07)	0,48 (0,06)
<i>Adaptive Boosting</i>	0,50 (0,03)	0,48 (0,03)	0,48 (0,05)
<i>Gradient Boosting</i>	0,50 (0,03)	0,44 (0,05)	0,43 (0,05)
<i>Extreme Gradient Boosting</i>	0,51 (0,04)	0,47 (0,03)	0,45 (0,05)
<i>Light Gradient Boosting Machine</i>	0,50 (0,03)	0,46 (0,02)	0,47 (0,06)

Tabela 5.1: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do grupo “Depressão”. O número entre parênteses indica o desvio padrão.

5.3 Causas de morte em subgrupos da depressão

Considerando os resultados apresentados anteriormente na Seção 5.2, onde nenhum algoritmo gerou um modelo com acurácia significativa para distinguir os indivíduos com depressão e sem depressão, foi realizado um estudo com subgrupos específicos de depressão, bem como com o Inventário Neuropsiquiátrico (NPI) dos indivíduos. Os subgrupos de depressão foram formados a partir das informações do SCID (Figura 4.1), presentes nos dados clínicos dos indivíduos. Cada subgrupo de depressão corresponde a um subconjunto dos dados originais do grupo de “Depressão”, conforme apresentado na Figura 5.2. Além disso, foram analisados 2 grupos relacionados ao NPI. O NPI é um instrumento de medida global de sintomas neuropsiquiátricos ou dos seus domínios sintomáticos específicos, como foi o caso deste estudo ao avaliar sintomas de depressão. Por meio do NPI, foi avaliado a intensidade da depressão nos últimos 3 meses que antecederam a morte de um determinado indivíduo. Dessa maneira, foram criadas novas variáveis que remetem a cada subgrupo de depressão e ao NPI.

Foram então realizados experimentos adicionais com 6 subgrupos de depressão e 2 grupos relacionados ao NPI, cujos resultados são apresentados nas seções a seguir. Para fins de simplificação os grupos do NPI serão tratados também como subgrupos de depressão. Portanto, os 3 experimentos descritos anteriormente (Seção 5.1) foram realizados nestes 8 subgrupos de depressão onde, para cada subgrupo, foi criada a sua variável de classe respectiva na base de dados.

A seleção dos subconjuntos de depressão manteve o pareamento das amostras. A Figura 5.2 ilustra os subgrupos de depressão com suas respectivas quantidades de casos positivos (indivíduos com depressão). Vale ressaltar que, a análise dos subgrupos de depressão seguiu todos os aspectos metodológicos da análise do grupo original “Depressão”. Apenas o subgrupo “Depressão álcool e droga” teve o número de *folds* reduzido (para $K = 6$) devido a sua quantidade pequena de amostras ($n = 29$). Nas seções seguintes são apresentados os resultados obtidos nos 3 experimentos em cada subgrupo.

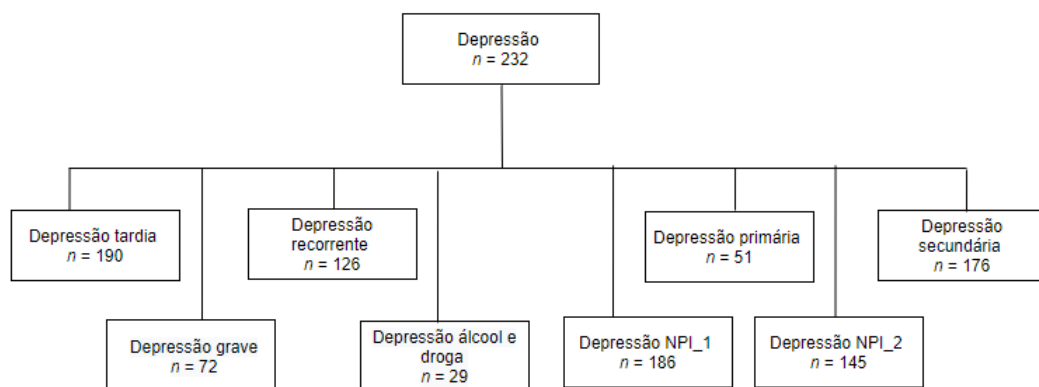


Figura 5.2: Subgrupos de depressão analisados.

5.3.1 Depressão tardia

O subgrupo de depressão tardia refere-se aos indivíduos que tiveram depressão com 60 anos ou mais. Neste grupo, a acurácia dos modelos gerados variou de 46% a 52%, independentemente do conjunto de variáveis utilizadas como entrada para os algoritmos, conforme Tabela 5.2.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,51 (0,09)	0,48 (0,07)	0,47 (0,07)
<i>Support Vector Machine</i>	0,50 (0,09)	0,48 (0,08)	0,47 (0,10)
<i>Multilayer Perceptron</i>	0,50 (0,07)	0,50 (0,08)	0,46 (0,05)
<i>Logistic Regression</i>	0,51 (0,08)	0,48 (0,08)	0,48 (0,10)
<i>Naïve Bayes</i>	0,52 (0,06)	0,47 (0,03)	0,47 (0,04)
<i>K-Nearest-Neighbor</i>	0,52 (0,08)	0,51 (0,06)	0,47 (0,07)
<i>Random Forest</i>	0,49 (0,08)	0,48 (0,07)	0,49 (0,07)
<i>Adaptive Boosting</i>	0,51 (0,08)	0,50 (0,07)	0,47 (0,09)
<i>Gradient Boosting</i>	0,50 (0,08)	0,48 (0,07)	0,47 (0,09)
<i>Extreme Gradient Boosting</i>	0,50 (0,07)	0,50 (0,07)	0,49 (0,07)
<i>Light Gradient Boosting Machine</i>	0,49 (0,07)	0,52 (0,09)	0,51 (0,09)

Tabela 5.2: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão tardia. O número entre parênteses indica o desvio padrão.

5.3.2 Depressão recorrente

A depressão recorrente é quando episódios depressivos são repetitivos e constantes. Os indivíduos com o indicador de haver outros episódios de depressão no questionário de avaliação do SCID fizeram parte deste subgrupo. Como pode ser observado na Tabela 5.3, os resultados da acurácia dos modelos gerados correspondente a este subgrupo variou de 41% a 53%, independentemente do conjunto de variáveis utilizadas como entrada.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,52 (0,08)	0,51(0,08)	0,49 (0,08)
<i>Support Vector Machine</i>	0,51 (0,08)	0,44 (0,06)	0,47 (0,06)
<i>Multilayer Perceptron</i>	0,52 (0,06)	0,46 (0,05)	0,48 (0,06)
<i>Logistic Regression</i>	0,53 (0,07)	0,49 (0,07)	0,47 (0,05)
<i>Naïve Bayes</i>	0,50 (0,03)	0,51 (0,05)	0,52 (0,06)
<i>K-Nearest-Neighbor</i>	0,47 (0,05)	0,47 (0,08)	0,49 (0,07)
<i>Random Forest</i>	0,51 (0,08)	0,50 (0,09)	0,51 (0,08)
<i>Adaptive Boosting</i>	0,52 (0,08)	0,49 (0,08)	0,48 (0,08)
<i>Gradient Boosting</i>	0,50 (0,08)	0,46 (0,08)	0,44 (0,07)
<i>Extreme Gradient Boosting</i>	0,50 (0,07)	0,41 (0,05)	0,41 (0,07)
<i>Light Gradient Boosting Machine</i>	0,51 (0,06)	0,47 (0,05)	0,49 (0,08)

Tabela 5.3: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão recorrente. O número entre parênteses indica o desvio padrão.

5.3.3 Depressão primária

A depressão primária aparece gradualmente e sem motivo aparente, não estando relacionada a qualquer doença (seja ela orgânica ou psiquiátrica) ou ao uso abusivo de álcool ou drogas. Os indivíduos incluídos neste subgrupo não apresentaram sintomas de depressão relacionados ao uso de substâncias, condição médica geral ou luto, avaliados como negativos no questionário de avaliação do SCID. A acurácia dos modelos gerados neste subgrupo variou de 44% a 60%, independentemente do conjunto de variáveis utilizadas como entrada. Os resultados do subgrupo depressão primária pode ser observado na Tabela 5.4.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,55 (0,18)	0,56 (0,20)	0,58 (0,22)
<i>Support Vector Machine</i>	0,53 (0,21)	0,55 (0,15)	0,47 (0,14)
<i>Multilayer Perceptron</i>	0,55 (0,15)	0,57 (0,14)	0,54 (0,18)
<i>Logistic Regression</i>	0,53 (0,15)	0,47 (0,13)	0,49 (0,13)
<i>Naïve Bayes</i>	0,51 (0,05)	0,53 (0,09)	0,53 (0,09)
<i>K-Nearest-Neighbor</i>	0,56 (0,15)	0,54 (0,17)	0,55 (0,18)
<i>Random Forest</i>	0,55 (0,16)	0,60 (0,16)	0,60 (0,18)
<i>Adaptive Boosting</i>	0,54 (0,11)	0,47 (0,10)	0,49 (0,14)
<i>Gradient Boosting</i>	0,54 (0,18)	0,57 (0,13)	0,57 (0,13)
<i>Extreme Gradient Boosting</i>	0,53 (0,14)	0,53 (0,16)	0,44 (0,12)
<i>Light Gradient Boosting Machine</i>	0,55 (0,20)	0,45 (0,15)	0,55 (0,12)

Tabela 5.4: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão primária. O número entre parênteses indica o desvio padrão.

5.3.4 Depressão secundária

A depressão secundária tem relação causal com algum tipo de doença (seja ela orgânica ou psiquiátrica) ou com uso de álcool e/ou drogas. Os indivíduos incluídos neste subgrupo apresentaram sintomas de depressão relacionados ao uso de substâncias, condição médica geral ou luto, com pelo menos uma resposta avaliada como positiva nestes itens no questionário de avaliação do SCID. Os resultados relacionados ao subgrupo depressão secundária variaram entre 42% e 51%, entre os Experimento 1, Experimento 2 ou Experimento 3, conforme Tabela 5.5.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,44 (0,08)	0,44 (0,04)	0,43 (0,05)
<i>Support Vector Machine</i>	0,47 (0,10)	0,43 (0,04)	0,42 (0,04)
<i>Multilayer Perceptron</i>	0,43 (0,10)	0,47 (0,06)	0,44 (0,04)
<i>Logistic Regression</i>	0,51 (0,14)	0,46 (0,05)	0,47 (0,06)
<i>Naïve Bayes</i>	0,44 (0,07)	0,45 (0,04)	0,45 (0,04)
<i>K-Nearest-Neighbor</i>	0,46 (0,11)	0,51 (0,07)	0,51 (0,07)
<i>Random Forest</i>	0,46 (0,10)	0,43 (0,05)	0,44 (0,06)
<i>Adaptive Boosting</i>	0,50 (0,12)	0,48 (0,05)	0,47 (0,06)
<i>Gradient Boosting</i>	0,44 (0,10)	0,44 (0,05)	0,43 (0,06)
<i>Extreme Gradient Boosting</i>	0,50 (0,12)	0,43 (0,06)	0,46 (0,06)
<i>Light Gradient Boosting Machine</i>	0,51 (0,16)	0,46 (0,06)	0,45 (0,05)

Tabela 5.5: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão secundária. O número entre parênteses indica o desvio padrão.

5.3.5 Depressão grave

Na depressão grave, existe uma associação ao quadro de ideação ou planeja-mento suicida, ou mesmo impacto importante na capacidade de trabalhar, cuidar das coisas em casa, ou se relacionar com as pessoas. Os indivíduos que tinham respostas positivas nestes dois itens no questionário de avaliação do SCID foram incluídos neste subgrupo. A acurácia dos modelos gerados variou de 44% a 56%, independentemente do conjunto de variáveis utilizadas como entrada.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,44 (0,08)	0,49 (0,12)	0,52 (0,07)
<i>Support Vector Machine</i>	0,47 (0,10)	0,54 (0,11)	0,48 (0,07)
<i>Multilayer Perceptron</i>	0,44 (0,10)	0,50 (0,13)	0,50 (0,07)
<i>Logistic Regression</i>	0,50 (0,14)	0,48 (0,11)	0,42 (0,07)
<i>Naïve Bayes</i>	0,44 (0,07)	0,46 (0,07)	0,50 (0,04)
<i>K-Nearest-Neighbor</i>	0,46 (0,11)	0,55 (0,11)	0,52 (0,07)
<i>Random Forest</i>	0,45 (0,11)	0,53 (0,16)	0,50 (0,08)
<i>Adaptive Boosting</i>	0,50 (0,11)	0,51 (0,12)	0,46 (0,07)
<i>Gradient Boosting</i>	0,44 (0,10)	0,50 (0,14)	0,44 (0,07)
<i>Extreme Gradient Boosting</i>	0,50 (0,11)	0,52 (0,13)	0,47 (0,07)
<i>Light Gradient Boosting Machine</i>	0,50 (0,16)	0,56 (0,13)	0,51 (0,03)

Tabela 5.6: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão grave. O número entre parênteses indica o desvio padrão.

5.3.6 Depressão por álcool e drogas

A depressão por álcool e drogas, acontece concomitantemente ou após 1 mês ao uso abusivo de álcool ou drogas. Neste subgrupo foram incluídos os indivíduos com resposta positiva para sintomas de depressão relacionados ao uso de substâncias (álcool ou drogas) no questionário do SCID. Conforme a Tabela 5.7, a acurácia variou de 31% a 55%, independentemente do conjunto de variáveis utilizadas como entrada.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,45 (0,13)	0,47 (0,07)	0,47 (0,05)
<i>Support Vector Machine</i>	0,43 (0,13)	0,35 (0,13)	0,38 (0,12)
<i>Multilayer Perceptron</i>	0,49 (0,07)	0,45 (0,10)	0,55 (0,15)
<i>Logistic Regression</i>	0,35 (0,15)	0,36 (0,09)	0,36 (0,17)
<i>Naïve Bayes</i>	0,52 (0,06)	0,53 (0,09)	0,53 (0,09)
<i>K-Nearest-Neighbor</i>	0,36 (0,09)	0,40 (0,16)	0,36 (0,17)
<i>Random Forest</i>	0,43 (0,13)	0,48 (0,07)	0,45 (0,10)
<i>Adaptive Boosting</i>	0,38 (0,13)	0,42 (0,10)	0,47 (0,15)
<i>Gradient Boosting</i>	0,47 (0,10)	0,45 (0,13)	0,45 (0,10)
<i>Extreme Gradient Boosting</i>	0,35 (0,12)	0,33 (0,07)	0,31 (0,06)
<i>Light Gradient Boosting Machine</i>	0,50 (0,00)	0,40 (0,11)	0,40 (0,12)

Tabela 5.7: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão por álcool e droga. O número entre parênteses indica o desvio padrão.

5.3.7 Depressão NPI_1

Depressão NPI_1, assim nomeada neste trabalho, corresponde aos sintomas leves de depressão nos últimos 3 meses que antecederam a morte. Os indivíduos pertencentes a este grupo apresentavam NPI com pontuação total igual a 0. O resultado dos experimentos deste subgrupo estão na Tabela 5.8 abaixo. A acurácia variou de 42% e 52% entre os Experimento 1, Experimento 2 e Experimento 3.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,45 (0,06)	0,51 (0,07)	0,52 (0,07)
<i>Support Vector Machine</i>	0,44 (0,06)	0,49 (0,06)	0,48 (0,07)
<i>Multilayer Perceptron</i>	0,46 (0,07)	0,52 (0,06)	0,50 (0,07)
<i>Logistic Regression</i>	0,44 (0,06)	0,42 (0,06)	0,42 (0,07)
<i>Naïve Bayes</i>	0,49 (0,05)	0,50 (0,04)	0,50 (0,04)
<i>K-Nearest-Neighbor</i>	0,49 (0,05)	0,47 (0,04)	0,52 (0,07)
<i>Random Forest</i>	0,45 (0,06)	0,52 (0,07)	0,52 (0,05)
<i>Adaptive Boosting</i>	0,46 (0,05)	0,47 (0,08)	0,46 (0,07)
<i>Gradient Boosting</i>	0,45 (0,07)	0,46 (0,08)	0,44 (0,09)
<i>Extreme Gradient Boosting</i>	0,45 (0,05)	0,47 (0,07)	0,47 (0,07)
<i>Light Gradient Boosting Machine</i>	0,45 (0,06)	0,48 (0,05)	0,51 (0,03)

Tabela 5.8: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão NPI_1. O número entre parênteses indica o desvio padrão.

5.3.8 Depressão NPI 2

Depressão NPI_2, assim nomeada neste trabalho, corresponde aos sintomas intensos de depressão nos últimos 3 meses que antecederam a morte. Os indivíduos deste grupo apresentavam NPI com pontuação total entre 4 e 12. Neste grupo, a acurácia variou de 44% a 55%, independentemente do conjunto de variáveis utilizadas como entrada, conforme Tabela 5.9.

Algoritmos de AM	Experimento 1	Experimento 2	Experimento 3
<i>Decision Tree</i>	0,53 (0,07)	0,50 (0,08)	0,51 (0,07)
<i>Support Vector Machine</i>	0,55 (0,06)	0,48(0,10)	0,47 (0,10)
<i>Multilayer Perceptron</i>	0,53 (0,07)	0,52 (0,08)	0,51 (0,08)
<i>Logistic Regression</i>	0,45 (0,08)	0,50 (0,09)	0,47 (0,07)
<i>Naïve Bayes</i>	0,45 (0,05)	0,51 (0,05)	0,51 (0,05)
<i>K-Nearest-Neighbor</i>	0,55 (0,07)	0,54 (0,09)	0,55 (0,09)
<i>Random Forest</i>	0,54 (0,06)	0,51 (0,10)	0,50 (0,09)
<i>Adaptive Boosting</i>	0,45 (0,07)	0,46 (0,10)	0,44 (0,07)
<i>Gradient Boosting</i>	0,53 (0,06)	0,48 (0,10)	0,46 (0,10)
<i>Extreme Gradient Boosting</i>	0,54 (0,06)	0,48 (0,08)	0,48 (0,08)
<i>Light Gradient Boosting Machine</i>	0,46 (0,08)	0,48 (0,10)	0,50 (0,06)

Tabela 5.9: Acurácia dos algoritmos de classificação em relação aos resultados dos Experimento 1, Experimento 2 e Experimento 3 do subgrupo depressão NPI_2. O número entre parênteses indica o desvio padrão.

5.4 Discussão

Muito embora não existam estudos na literatura baseados em dados de autópsia para determinar causas de morte na depressão, estudos recentes de revisão associam a depressão ao aumento da mortalidade, porém, as evidências da causas de morte e depressão não são significativas [49, 52].

Na revisão sistemática realizada por Machado *et al.* [49], foi realizada a coleta e avaliação de várias revisões sistemáticas e meta-análises, a credibilidade de cada associação foi classificada conforme as seguintes categorias: convincente (classe I), altamente sugestivo (classe II), sugestivo (classe III), evidências fracas e associações não significativas. Nenhuma associação preencheu os critérios para evidências de classe I, enquanto apenas 4 associações, a saber, as de depressão e mortalidade por todas as causas em câncer, insuficiência cardíaca, ambientes diversos, bem como entre pacientes após infarto agudo do miocárdio, foram apoiadas apenas por evidências classe II. Análises de sensibilidade indicaram que as diferenças na apuração de casos de depressão, bem como a falta de ajuste adequado para variáveis de confusão e outros fatores de risco importantes, podem gerar várias associações apoiadas apenas por níveis mais baixos de evidência. Ajustar para idade e sexo, por exemplo, foi considerado essencial. Desta forma, quando apenas estudos que controlavam por idade e sexo foram considerados, a associação de depressão e mortalidade por todas as causas no câncer não foi mais apoiada por evidências classe II. Portanto, esta revisão sugere que as inferências causais entre depressão e

mortalidade por todas as causas em populações distintas não parecem ser tão conclusivas quanto se pensava. Além disso, nesta revisão nenhuma associação foi apoiada por evidências classe II quando apenas os estudos que empregaram entrevistas diagnósticas estruturadas/semiestruturadas foram considerados, o que foi o caso do presente estudo que utilizou o SCID para o diagnóstico de depressão maior e o NPI para o diagnóstico de sintomas depressivos nos meses que antecederam a morte do indivíduo.

Outras evidências que apoiam nossos resultados são o fato do tratamento da depressão não aumentar a sobrevivência na maioria dos estudos [49]. No entanto, existem exceções como a de um dos poucos estudos feitos em idosos com depressão grave, onde o tratamento pode reduzir a mortalidade [24]. Isso ressalta a imensa variabilidade nos resultados encontrados até o momento, tanto na amostra deste estudo como na literatura.

Estas divergências entre os estudos podem ser atribuídas a aspectos metodológicos, como também afirmam Miloyan, Beyon, and Eiko Fried [52]. Os autores realizaram um estudo de revisão sistemática sobre a relação entre depressão e mortalidade por todas as causas e apresentaram algumas razões metodológicas para a heterogeneidade entre os 293 estudos analisados. Apontaram que a divergência encontrada pode ser atribuída a alguns aspectos, tais como: tamanho e características da amostra, número de óbitos e períodos de acompanhamento, ajuste para transtornos mentais e comportamentos de saúde. Além disso, mais de 40 instrumentos diferentes foram usados para medir sintomas depressivos e, mesmo os estudos que usaram os mesmos questionários, frequentemente adotaram diferentes pontos de corte para um provável diagnóstico de transtorno depressivo. Isso atrapalha mais uma vez a comparação entre estudos dada a variabilidade das variáveis analisadas bem como a falta de ajuste por variáveis de confusão e análises de sensibilidade já citados por Machado *et al.* [49]. Desta forma, os resultados do presente estudo que divergem de parte da literatura podem ser devido a uma população distinta estudada – idosos provenientes da comunidade e não de centros especializados de tratamento de depressão, bem como falta de controle de variáveis de confusão. Neste estudo, foi feito um pareamento cuidadoso entre os indivíduos com depressão e seus controles comparativos evitando, em parte, a interferência de variáveis de confusão. No presente estudo, subgrupos de depressão (Figura 5.2) foram analisados e nenhum padrão distintivo de morte foi encontrado. A análise dos subgrupos de depressão permite avaliar não apenas a depressão maior, mas avaliar outros tipos de depressão que muitas vezes podem ficar não reconhecidas ou tratadas.

Dentre os inúmeros estudos da literatura com resultados similares aos nossos, destaca-se o estudo de Wu *et al.* [71] por incluir uma grande amostra de pacientes (1.087.125) com diabetes que tinham transtornos depressivos e pacientes com diabetes sem depressão para avaliar complicações cardiovasculares. Os pacientes que desenvolveram complicações macrovasculares e microvasculares bem como aqueles que morreram

foram analisados. Não houve associação de depressão com complicações microvasculares, mortalidade por doenças cardiovasculares e mortalidade por diabetes mellitus. O efeito da depressão nas complicações e mortalidade do diabetes foi mais proeminente entre os adultos jovens do que entre os adultos de meia-idade e idosos. Já no estudo de Sullivan *et al.* [65] examinaram a relação entre depressão, mortalidade e eventos cardiovasculares em uma amostra que recebia tratamento padronizado para diabetes. O estudo também incluiu critérios para definir complicações macrovasculares e microvasculares e a causa da morte. Examinaram o efeito da depressão na mortalidade entre aqueles com e sem doença cardiovascular prévia. Como resultado, indicaram que a depressão aumenta o risco de mortalidade por todas as causas e pode aumentar o risco de eventos macrovasculares entre adultos com diabetes tipo 2 com alto risco de eventos cardiovasculares.

A investigação para o tratamento da depressão e sua relação com as causas de morte, vale a pena por vários outros motivos, por exemplo, melhoria da qualidade de vida, mas não com a expectativa de que o risco de morte diminuirá. Além disso, as intervenções que visam promover um estilo de vida saudável, bem como o cuidado adequado de condições somáticas concomitantes em pessoas com depressão também podem levar a uma diminuição na mortalidade por todas as causas [49, 46]. No entanto, o impacto dessas intervenções em um nível individual, social e do sistema de saúde sobre a sobrevivência de todas as causas justifica uma investigação mais aprofundada [49].

Nesta investigação, foi explorado e analisado, por meio de algoritmos de AM, dados relacionados à morte de indivíduos com depressão. Algoritmos de AM têm como objetivo descobrir princípios gerais subjacentes a uma série de observações sem instruções explícitas [10]. Seus métodos são caracterizados por 1) fazer poucas suposições formais, 2) permitir que os dados “falem por si mesmos” e 3) a capacidade de extrair conhecimento estruturado de dados extensos [9]. Para estimar a previsão do modelo, foi considerado como dados de entrada para os algoritmos de aprendizado de máquina, as variáveis das Tabela 4.4 e 4.5. Além disso, a variável considerada como atributo alvo, determinou se um indivíduo pertencia ao grupo Depressão ou ao grupo controle (sem depressão). De acordo com os resultados dos diversos experimentos realizados, independente do modelo de aprendizado de máquina gerado, um baixo desempenho, refletido pela acurácia, foi obtido na classificação do grupo Depressão e seus subgrupos comparados ao grupo sem depressão. A acurácia de todos os experimentos realizados teve uma média de 48,4%. Portanto, a interpretação dos modelos gerados dada a acurácia obtida, não é significativa para se determinar um padrão de causa de morte dos indivíduos com depressão.

Conclusões

No presente estudo clínico-patológico transversal de base populacional, foram avaliadas causas de morte de 1.136 indivíduos segundo laudos de autópsia. Indivíduos com depressão foram comparados com indivíduos sem depressão para se determinar diferenças entre as causas de mortes. Foram consideradas nesta análise até 4 possíveis doenças diretamente relacionadas à morte de um determinado indivíduo conforme a informação contida nos relatórios de autópsia de corpo inteiro.

As causas de mortes foram agrupadas em diferentes níveis de categorias da CID-10. Por meio da Entrevista Clínica Estruturada para os Transtornos do DSM-IV (do inglês, SCID) os indivíduos foram agrupados primeiramente no grupo geral de depressão. Posteriormente, considerou-se subgrupos de depressão baseados tanto no SCID quanto no NPI.

Uma mineração de dados foi realizada seguindo o processo de Descoberta de Conhecimento em Bases de Dados (do inglês, KDD). Onze algoritmos de AM, bem estabelecidos na literatura e de diferentes vieses de aprendizado, foram aplicados, independentemente, em cada um dos grupos de depressão a fim de distinguir os indivíduos com e sem depressão baseados somente nas causas relacionadas à morte. No entanto, a baixa acurácia dos modelos gerados não permitiu que fossem extraídos padrões relevantes destes modelos.

Embora este estudo tenha feito uma ampla investigação na população em geral e em grupos específicos de pacientes, os resultados obtidos pelos modelos gerados não indicam diferenças de padrões nas causas de morte em indivíduos com e sem depressão. Este resultado corrobora com estudos anteriores da literatura onde as evidências das causas de morte por todas as causas e causas específicas e depressão não são significativas.

6.1 Contribuições

Dentre as principais contribuições deste trabalho pode-se citar a originalidade da análise de causas de morte baseado em laudo de autópsia de corpo inteiro, com o objetivo de se buscar associação entre as doenças que causaram a morte em pessoas

que tiveram depressão utilizando algoritmos de AM. Outro ponto de destaque deste estudo é em relação a amostra da população que fez parte deste experimento. A amostra incluiu pessoas acima de 50 anos de uma comunidade não hospitalizada, isso reflete a realidade dos casos menos graves de depressão, que acontecem com grande parte da população, ou seja, muitos casos de depressão não ficam conhecidos, já que uma pequena parcela dos casos são tratadas em centros de referência. Em geral, estudos anteriores apresentam seus resultados em ambientes como serviços especializados ou em pacientes hospitalizados. Embora as informações da entrevista clínica deste estudo tenha sido realizada no *post-mortem*, o formulário da entrevista clínica foi utilizada em ambientes clínicos para validar sua confiabilidade [22]. Além disso, a amostra avaliada é sobre uma população diversificada em termos de etnia e escolaridade de um país com renda média baixa, que é uma característica do Brasil, que possui baixo desenvolvimento econômico e social, marcado pela desigualdade, o que torna esse estudo importante para a realidade brasileira. Vale ressaltar ainda que este estudo apresentou como metodologia de análise o uso de algoritmos de AM. Tal análise não foi encontrada na literatura com as mesmas características da amostra avaliada.

6.2 Limitações

O estudo realizado apresentou limitações importantes quanto à amostra analisada. Os dados desta pesquisa são referentes apenas à população da cidade de São Paulo. Não foram incluídos os casos de morte não natural e os casos onde já se sabe a causa da morte. Um outro ponto sobre a amostra analisada, é que os casos analisados se referem aos indivíduos com 50 anos ou mais, não abrangendo uma população com faixa etária de idade mais distribuída. Além disso, os indivíduos não foram acompanhados durante a vida, ou mesmo durante um determinado tempo antes da morte. As informações em relação aos indivíduos foram realizadas no *post-mortem* por meio de entrevista com um informante. Para aumentar a confiabilidade desses dados, foram incluídos apenas participantes que tivessem contato pelo menos semanalmente com o informante e foram excluídos os indivíduos quando o informante forneceu informações conflitantes durante a entrevista clínica. Por fim, o uso de dados relatados por informantes coletados retrospectivamente é uma preocupação, uma vez que os informantes podem não estar cientes de alguns tratamentos e distúrbios do falecido.

6.3 Trabalhos Futuros

Como trabalhos futuros e continuidade a esta pesquisa, é sugerido estudar a associação da depressão e causas de morte, não considerando apenas as doenças causadoras da morte, porém explorando outras variáveis que podem ser coletas em estudos longitudinais, como por exemplo, considerar as doenças que o indivíduo teve durante toda a sua vida até a morte. Além disso, é importante ampliar os experimentos explorando outros métodos de AM, tais como a técnica de Aprendizado Relacional e técnicas de agrupamentos. Na representação de Aprendizado Relacional, objetos podem ser descritos em termos de seus componentes e relações entre esses componentes, além de possuir uma alta expressividade para representar conceitos e a importante habilidade de representar conhecimento do domínio, o que pode auxiliar na representação hierárquica das causas relacionadas à morte em indivíduos com depressão ou mesmo, em caso de estudos longitudinais, representar de maneira hierárquica as doenças que um indivíduo teve ao longo da vida até a morte. A técnica de agrupamento examina relações de interdependência entre todo o conjunto de variáveis, com o objetivo de agrupar objetos semelhantes. Esta análise pode permitir analisar subgrupos específicos de doenças, de modo que indivíduos que apresentem menor distância entre si são mais semelhantes, ou seja, possuam características similares formando grupos homogêneos. Por outro lado, grupos mais distantes participam de conglomerados distintos e conseqüentemente características diferentes.

Referências Bibliográficas

- [1] AGGARWAL, C. C. **Data mining: the textbook**. Springer, 2015.
- [2] ALPAYDIN, E. **Introduction to machine learning**. MIT press, 2010.
- [3] ARRIBAS, J. I.; CALHOUN, V. D.; ADALI, T. **Automatic bayesian classification of healthy controls, bipolar disorder, and schizophrenia using intrinsic connectivity maps from fmri data**. *IEEE Transactions on Biomedical Engineering*, 57(12):2850–2860, 2010.
- [4] BHARATI, M.; RAMAGERI, B. **Data mining techniques and applications**. *Indian Journal of Computer Science and Engineering*, 1, 12 2010.
- [5] BONACCORSO, G. **Machine learning algorithms**. Packt Publishing Ltd, 2017.
- [6] BRODERSEN, K. H.; ONG, C. S.; STEPHAN, K. E.; BUHMANN, J. M. **The balanced accuracy and its posterior distribution**. In: *2010 20th international conference on pattern recognition*, p. 3121–3124. IEEE, 2010.
- [7] BUJA, L. M.; BARTH, R. F.; KRUEGER, G. R.; BRODSKY, S. V.; HUNTER, R. L. **The importance of the autopsy in medicine: perspectives of pathology colleagues**. *Academic pathology*, 6:2374289519834041, 2019.
- [8] BUTNORIENE, J.; BUNEVICIUS, A.; SAUDARGIENE, A.; NEMEROFF, C. B.; NORKUS, A.; CICENIENE, V.; BUNEVICIUS, R. **Metabolic syndrome, major depression, generalized anxiety disorder, and ten-year all-cause and cardiovascular mortality in middle aged and elderly patients**. *International journal of cardiology*, 190:360–366, 2015.
- [9] BZDOK, D.; ALTMAN, N.; KRZYWINSKI, M. **Points of significance: statistics versus machine learning**. *Nature Publishing Group*, 2018.
- [10] BZDOK, D.; MEYER-LINDENBERG, A. **Machine learning for precision psychiatry: opportunities and challenges**. *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, 3(3):223–230, 2018.

- [11] CARVALHO, A.; FACELI, K.; LORENA, A.; GAMA, J. **Inteligência artificial—uma abordagem de aprendizado de máquina**. Rio de Janeiro: LTC, 2011.
- [12] CHANG, C.-K.; HAYES, R. D.; PERERA, G.; BROADBENT, M. T.; FERNANDES, A. C.; LEE, W. E.; HOTOPE, M.; STEWART, R. **Life expectancy at birth for people with serious mental illness and other major disorders from a secondary mental health care case register in london**. *PloS one*, 6(5), 2011.
- [13] CHICCO, D.; JURMAN, G. **The advantages of the matthews correlation coefficient (mcc) over f1 score and accuracy in binary classification evaluation**. *BMC genomics*, 21(1):1–13, 2020.
- [14] CUIJPERS, P.; SMIT, F. **Excess mortality in depression: a meta-analysis of community studies**. *Journal of affective disorders*, 72(3):227–236, 2002.
- [15] DAMIAN, J.; PASTOR-BARRIUSO, R.; VALDERRAMA-GAMA, E.; DE PEDRO-CUESTA, J. **Association of detected depression and undetected depressive symptoms with long-term mortality in a cohort of institutionalised older people**. *Epidemiology and psychiatric sciences*, 26(2):189–198, 2017.
- [16] DAS-MUNSHI, J.; CHANG, C.-K.; SCHOFIELD, P.; STEWART, R.; PRINCE, M. J. **Depression and cause-specific mortality in an ethnically diverse cohort from the uk: 8-year prospective study**. *Psychological medicine*, 49(10):1639–1651, 2019.
- [17] DINIZ, B. S.; REYNOLDS III, C. F.; BUTTERS, M. A.; DEW, M. A.; FIRMO, J. O.; LIMA-COSTA, M. F.; CASTRO-COSTA, E. **The effect of gender, age, and symptom severity in late-life depression on the risk of all-cause mortality: The bambuí cohort study of aging**. *Depression and anxiety*, 31(9):787–795, 2014.
- [18] DUARTE, J. C. **O Algoritmo Boosting at Start e suas Aplicações**. PhD thesis, Tese (Doutorado)—PUC-RJ, 2009. 45, 2009.
- [19] ELLINGER, F.; BEZERRA, K. **Manual de Procedimentos do Serviço de Verificação de Óbitos de Marília**. FAMEMA, 2011.
- [20] FAYYAD, U.; PIATETSKY-SHAPIRO, G.; SMYTH, P. **From data mining to knowledge discovery in databases**. *AI magazine*, 17(3):37–37, 1996.
- [21] FERNANDES, F.; CASTILLO, P.; BASSAT, Q.; QUINTÓ, L.; HURTADO, J. C.; MARTÍNEZ, M. J.; LOVANE, L.; JORDAO, D.; BENE, R.; NHAMPOSSA, T.; OTHERS. **Contribution of the clinical information to the accuracy of the minimally invasive and the complete diagnostic autopsy**. *Human pathology*, 85:184–193, 2019.

- [22] FERRETTI, R. E. D. L.; DAMIN, A. E.; BRUCKI, S. M. D.; MORILLO, L. S.; PERROCO, T. R.; CAMPORA, F.; MOREIRA, E. G.; BALBINO, É. S.; LIMA, M. D. C. D. A.; BATTELA, C.; OTHERS. **Post-mortem diagnosis of dementia by informant interview.** *Dementia & neuropsychologia*, 4:138–144, 2010.
- [23] FERRETTI-REBUSTINI, R. E. D. L.; JACOB-FILHO, W.; SUEMOTO, C. K.; FARFEL, J. M.; LEITE, R. E. P.; GRINBERG, L. T.; PASQUALUCCI, C. A.; NITRINI, R. **Factors associated with morphometric brain changes in cognitively normal aging.** *Dementia & neuropsychologia*, 9:103–109, 2015.
- [24] GALLO, J. J.; BOGNER, H. R.; MORALES, K. H.; POST, E. P.; TEN HAVE, T.; BRUCE, M. L. **Depression, cardiovascular disease, diabetes, and two-year mortality among older, primary-care patients.** *The American journal of geriatric psychiatry*, 13(9):748–755, 2005.
- [25] GAO, S.; CALHOUN, V. D.; SUI, J. **Machine learning in major depression: From classification to treatment outcome prediction.** *CNS neuroscience & therapeutics*, 24(11):1037–1052, 2018.
- [26] GAROFALO, M.; PICCOLI, L.; ROMEO, M.; BARZAGO, M. M.; RAVASIO, S.; FOGLIERINI, M.; MATKOVIC, M.; SGRIGNANI, J.; DE GASPARO, R.; PRUNOTTO, M.; OTHERS. **Machine learning analyses of antibody somatic mutations predict immunoglobulin light chain toxicity.** *Nature communications*, 12(1):1–10, 2021.
- [27] GILL, J. R. **From death to death certificate: what do the dead say?** In: *Journal of Medical Toxicology*, volume 13, p. 111–116. Springer, 2017.
- [28] GODWIN, T. A. **End of life: natural or unnatural death investigation and certification.** *Disease-a-Month*, 4(51):218–277, 2005.
- [29] GRINBERG, L. T.; NITRINI, R.; SUEMOTO, C. K.; DE LUCENA FERRETTI-REBUSTINI, R. E.; LEITE, R. E.; FARFEL, J. M.; SANTOS, E.; ANDRADE, M. P. G. D.; ALHO, A. T. D. L.; LIMA, M. D. C.; OTHERS. **Prevalence of dementia subtypes in a developing country: a clinicopathological study.** *Clinics*, 68:1140–1145, 2013.
- [30] GROSS, A. L.; GALLO, J. J.; EATON, W. W. **Depression and cancer risk: 24 years of follow-up of the baltimore epidemiologic catchment area sample.** *Cancer causes & control*, 21(2):191–199, 2010.
- [31] HAN, J.; PEI, J.; KAMBER, M. **Data mining: concepts and techniques.** Elsevier, 2011.

- [32] HO, C. S.; JIN, A.; NYUNT, M. S. Z.; FENG, L.; NG, T. P. **Mortality rates in major and subthreshold depression: 10-year follow-up of a singaporean population cohort of older adults.** *Postgraduate medicine*, 128(7):642–647, 2016.
- [33] HUYS, Q. J.; MAIA, T. V.; FRANK, M. J. **Computational psychiatry as a bridge from neuroscience to clinical applications.** *Nature neuroscience*, 19(3):404, 2016.
- [34] ISHITANI, L. H.; FRANÇA, E. **Uso das causas múltiplas de morte em saúde pública.** *Informe Epidemiológico do SUS*, 10(4):163–175, 2001.
- [35] JAMES, G.; PATTON, R. E.; HESLIN, A. S. **Accuracy of cause-of-death statements on death certificates.** *Public health reports*, 70(1):39, 1955.
- [36] KESSLER, R. C.; BROMET, E. J. **The epidemiology of depression across cultures.** *Annual review of public health*, 34:119–138, 2013.
- [37] KING, L. S.; MEEHAN, M. C. **A history of the autopsy. a review.** *The American journal of pathology*, 73(2):514, 1973.
- [38] KOHAVI, R. **A study of cross-validation and bootstrap for accuracy estimation and model selection.** In: *IJCAI'95: Proceedings of the 14th international joint conference on Artificial intelligence*, volume 2, p. 1137–1143. Morgan Kaufmann Publishers Inc., 1995.
- [39] KRITTANAWONG, C.; VIRK, H. U. H.; KUMAR, A.; AYDAR, M.; WANG, Z.; STEWART, M. P.; HALPERIN, J. L. **Machine learning and deep learning to predict mortality in patients with spontaneous coronary artery dissection.** *Scientific reports*, 11(1):1–10, 2021.
- [40] KÜBLER, U. **Structured Clinical Interview for DSM-IV (SCID).** In: *Encyclopedia of Behavioral Medicine*, chapter Structured, p. 1919–1920. Springer New York, New York, NY, 2013.
- [41] LAURENTI, R. **As manifestações de sofrimento mental mais frequentes na comunidade.** *SMAD, Revista Eletrônica em Salud Mental, Alcohol y Drogas*, 3(2):0, 2007.
- [42] LAURENTI, R.; JORGE, M. H. P. **O atestado de óbito: aspectos médicos, estatísticos, éticos e jurídicos.** In: *O atestado de óbito: aspectos médicos, estatísticos, éticos e jurídicos*, p. 154–154. 2015.
- [43] LAURSEN, T. M.; MUSLINER, K. L.; BENROS, M. E.; VESTERGAARD, M.; MUNK-OLSEN, T. **Mortality and life expectancy in persons with severe unipolar depression.** *Journal of affective disorders*, 193:203–207, 2016.

- [44] LIBRENZA-GARCIA, D.; KOTZIAN, B. J.; YANG, J.; MWANGI, B.; CAO, B.; LIMA, L. N. P.; BERMUDEZ, M. B.; BOEIRA, M. V.; KAPCZINSKI, F.; PASSOS, I. C. **The impact of machine learning techniques in the study of bipolar disorder: a systematic review.** *Neuroscience & Biobehavioral Reviews*, 80:538–554, 2017.
- [45] LIN, S.-K.; HSIU, H.; CHEN, H.-S.; YANG, C.-J. **Classification of patients with alzheimer’s disease using the arterial pulse spectrum and a multilayer-perceptron analysis.** *Scientific reports*, 11(1):1–14, 2021.
- [46] LIU, N. H.; DAUMIT, G. L.; DUA, T.; AQUILA, R.; CHARLSON, F.; CUIJPERS, P.; DRUSS, B.; DUDEK, K.; FREEMAN, M.; FUJII, C.; OTHERS. **Excess mortality in persons with severe mental disorders: a multilevel intervention framework and priorities for clinical practice, policy and research agendas.** *World psychiatry*, 16(1):30–40, 2017.
- [47] LUIZ, R. R.; STRUCHINER, C. J. **Inferência causal em epidemiologia: o modelo de respostas potenciais.** Editora Fiocruz, 2002.
- [48] MACÁRIO FILHO, V. **Um novo algoritmo de agrupamento semisupervisionado baseado no fuzzy c-means.** Master’s thesis, Universidade Federal de Pernambuco, 2009.
- [49] MACHADO, M. O.; VERONESE, N.; SANCHES, M.; STUBBS, B.; KOYANAGI, A.; THOMPSON, T.; TZOULAKI, I.; SOLMI, M.; VANCAMPFORT, D.; SCHUCH, F. B.; OTHERS. **The association of depression and all-cause and cause-specific mortality: an umbrella review of systematic reviews and meta-analyses.** *BMC medicine*, 16(1):112, 2018.
- [50] MALCHER, P. R. C.; FERREIRA, D. A. L.; OLIVEIRA, S. R. B.; VASCONCELOS, A. M. L. D.; OTHERS. **Um mapeamento sistemático sobre abordagens de apoio à rastreabilidade de requisitos no contexto de projetos de software.** 2015.
- [51] MENG, R.; YU, C.; LIU, N.; HE, M.; LV, J.; GUO, Y.; BIAN, Z.; YANG, L.; CHEN, Y.; ZHANG, X.; OTHERS. **Association of depression with all-cause and cardiovascular disease mortality among adults in china.** *JAMA network open*, 3(2):e1921043–e1921043, 2020.
- [52] MILOYAN, B.; FRIED, E. **A reassessment of the relationship between depression and all-cause mortality in 3,604,005 participants from 293 studies.** *World Psychiatry*, 16(2):219, 2017.

- [53] NAKAGAWA, E. Y.; SCANNAVINO, K. R. F.; FABBRI, S. C. P. F.; FERRARI, F. C. **Revisão sistemática da literatura em engenharia de software: teoria e prática.** 2017.
- [54] ORGANIZATION, W. H.; OTHERS. **Depression and other common mental disorders: global health estimates,** 2017.
- [55] OSELKA, G. **Atestado médico – prática e ética.** Technical report, Conselho Regional de Medicina do Estado de São Paulo, 2013.
- [56] PENNINX, B. W.; BEEKMAN, A. T.; HONIG, A.; DEEG, D. J.; SCHOEVERS, R. A.; VAN EIJK, J. T.; VAN TILBURG, W. **Depression and cardiac mortality: results from a community-based longitudinal study.** *Archives of general psychiatry*, 58(3):221–227, 2001.
- [57] REUTFORS, J.; ANDERSSON, T. M.; BRENNER, P.; BRANDT, L.; DIBERNARDO, A.; LI, G.; HÄGG, D.; WINGÅRD, L.; BODÉN, R. **Mortality in treatment-resistant unipolar depression: A register-based cohort study in sweden.** *Journal of affective disorders*, 238:674–679, 2018.
- [58] RUSSEL, S.; NORVIG, P. **Inteligência artificial. 2ª edição.** Rio de Janeiro: Campus, 2004.
- [59] RUTLEDGE, R. B.; CHEKROUD, A. M.; HUYS, Q. J. **Machine learning and big data in psychiatry: toward clinical applications.** *Current opinion in neurobiology*, 55:152–159, 2019.
- [60] SAINT ONGE, J. M.; KRUEGER, P. M.; ROGERS, R. G. **The relationship between major depression and nonsuicide mortality for us adults: the importance of health behaviors.** *Journals of Gerontology Series B: Psychological Sciences and Social Sciences*, 69(4):622–632, 2014.
- [61] SCHUH, G.; REINHART, G.; PROTE, J.-P.; SAUERMAN, F.; HORSTHOFER, J.; OPPOLZER, F.; KNOLL, D. **Data mining definitions and applications for the management of production complexity.** *Procedia CIRP*, 81:874–879, 2019.
- [62] STEWART, R. A.; NORTH, F. M.; WEST, T. M.; SHARPLES, K. J.; SIMES, R. J.; COLQUHOUN, D. M.; WHITE, H. D.; TONKIN, A. M. **Depression and cardiovascular morbidity and mortality: cause or consequence?** *European heart journal*, 24(22):2027–2037, 2003.

- [63] SUEMOTO, C. K.; DAMICO, M. V.; FERRETTI, R. E.; GRINBERG, L. T.; FARFEL, J. M.; LEITE, R. E.; NITRINI, R.; LAFER, B.; JACOB-FILHO, W.; PASQUALUCCI, C. A.; OTHERS. **Depression and cardiovascular risk factors: evidence from a large postmortem sample.** *International journal of geriatric psychiatry*, 28(5):487–493, 2013.
- [64] SUEMOTO, C. K.; LEITE, R. E.; FERRETTI-REBUSTINI, R. E.; RODRIGUEZ, R. D.; NITRINI, R.; PASQUALUCCI, C. A.; JACOB-FILHO, W.; GRINBERG, L. T. **Neuropathological lesions in the very old: results from a large brazilian autopsy study.** *Brain Pathology*, 29(6):771–781, 2019.
- [65] SULLIVAN, M. D.; O'CONNOR, P.; FEENEY, P.; HIRE, D.; SIMMONS, D. L.; RAISCH, D. W.; FINE, L. J.; NARAYAN, K. V.; ALI, M. K.; KATON, W. J. **Depression predicts all-cause mortality: epidemiological evaluation from the accord hrql substudy.** *Diabetes care*, 35(8):1708–1715, 2012.
- [66] TAN, P.-N.; STEINBACH, M.; KUMAR, V. **Introdução ao datamining: mineração de dados.** Ciência Moderna, 2009.
- [67] WALKER, E. R.; MCGEE, R. E.; DRUSS, B. G. **Mortality in mental disorders and global disease burden implications: a systematic review and meta-analysis.** *JAMA psychiatry*, 72(4):334–341, 2015.
- [68] WALL, M. M.; HUANG, J.; OSWALD, J.; MCCULLEN, D. **Factors associated with reporting multiple causes of death.** *BMC Medical Research Methodology*, 5(1):4, 2005.
- [69] WEI, J.; HOU, R.; ZHANG, X.; XU, H.; XIE, L.; CHANDRASEKAR, E. K.; YING, M.; GOODMAN, M. **The association of late-life depression with all-cause and cardiovascular mortality among community-dwelling older adults: systematic review and meta-analysis.** *The British Journal of Psychiatry*, 215(2):449–455, 2019.
- [70] (WHO), W. H. O.; OTHERS. **Icd-11 implementation or transition guide. 2019,** 2019.
- [71] WU, C.-S.; HSU, L.-Y.; WANG, S.-H. **Association of depression and diabetes complications and mortality: a population-based cohort study.** *Epidemiology and psychiatric sciences*, 29, 2020.
- [72] WULSIN, L. R.; VAILLANT, G. E.; WELLS, V. E. **A systematic review of the mortality of depression.** *Psychosomatic medicine*, 61(1):6–17, 1999.

Descrição dos dados de autópsia e dos dados da entrevista clínica

Descrição dos dados de autópsia

Na Tabela A.1 são apresentados os dados de autópsia utilizados neste trabalho.

Variável	Descrição	Valores
nsvo	identificação padronizada do indivíduo	valor do tipo string
cid_obito	CID da causa imediata do óbito	código de 3 dígitos formado por 1 letra e 3 números
causa_obito	descrição livre da causa imediata do óbito	valor do tipo string
cid_principal	CID da causa básica do óbito	código de 3 dígitos formado por 1 letra e 3 números
doenca_principal	descrição livre da causa principal do óbito	valor do tipo string
cid_basica1	CID da causa intermediária 1 do óbito	código de 3 dígitos formado por 1 letra e 3 números
causa_basica1	descrição livre da causa intermediária 1	valor do tipo string
cid_basica2	CID da causa intermediária 2 do óbito	código de 3 dígitos formado por 1 letra e 3 números
causa_basica2	descrição livre da causa intermediária 2	valor do tipo string
age	idade do óbito	valor do tipo real
gender	sexo biológico do indivíduo	feminino masculino
height	altura do indivíduo	valor do tipo real

Tabela A.1: Descrição dos dados do laudo de autópsia.

Descrição dos dados da entrevista clínica

Na Tabela A.2 são apresentados os dados sócio-demográficos, clínicos, NPI e do SCID.

Tabela A.2: Descrição dos dados da entrevista clínica.

Variável	Descrição	Valores
nsvo	identificação do indivíduo	valor do tipo string
idade	idade do óbito	valor do tipo real
sexo	sexo biológico do indivíduo	feminino masculino
etnia	características fenotípicas	caucasiano não caucasiano
civil	estado civil	solteiro casado viúvo amasiado divorciado
situacao	situação em relação ao trabalho	ainda trabalhava aposentado aposentado por invalidez outro
escolaridade	nível de escolaridade	analfabeto educação formal (1 - 4 anos) educação forma (5 anos ou mais)
has	avaliação da hipertensão do indivíduo	sim não
dm	avaliação da diabetes mellitus do indivíduo	sim não
dac	avaliação sobre doença arterial coronária	sim não
dlp	avaliação sobre dislipidemia	sim não
aids	se o indivíduo tinha a síndrome da imunodeficiência adquirida	sim não
sifilis	se o indivíduo tinha sífilis	sim não

Continua na próxima página

Variável	Descrição	Valores
neoplasias	se o indivíduo tinha algum tumor	sim não
ativiFisica	se o indivíduo praticava atividade física	caminhada doméstica trabalho outra
acamado	se o indivíduo estava acamado antes do óbito	sim não
peso	peso do indivíduo	valor do tipo real
altura	altura do indivíduo	valor do tipo real
imc_auto	índice de massa corporal	valor do tipo real
tabagismo	se o indivíduo fumava	nunca fumou fumava atualmente parou
etilismo	se o indivíduo consumia bebida alcoólica	nunca bebeu bebia socialmente alcoólatra parou
antidepressivo	se o indivíduo fazia uso de antidepressivos	sim não
iqcode_total	escala que mede a piora cognitiva nos últimos 10 anos	de 3 a 5 (contínuo)
cdr	avaliação do declínio cognitivo/demência	sim não
np1delerios	pontuação dos sintomas de delírios	0 a 12
np2alucinacoes	pontuação dos sintomas de alucinações	0 a 12
np3agitacao	pontuação dos sintomas de agitação ou agressão	0 a 12
np4depressao	pontuação dos sintomas de depressão ou disforia	0 a 12
np5ansiedade	pontuação dos sintomas de ansiedade	0 a 12

Continua na próxima página

Variável	Descrição	Valores
npi6elacao	pontuação dos sintomas de elação	0 a 12
npi7apatia	pontuação dos sintomas de apatia ou indiferença	0 a 12
npi8desinibicao	pontuação dos sintomas de desinibição	0 a 12
npi9irritabilidade	pontuação dos sintomas de irritabilidade	0 a 12
npi10distmot	pontuação dos sintomas do distúrbio motor	0 a 12
npi11compnot	pontuação dos sintomas de comportamento noturno	0 a 12
npi12apetite	pontuação dos sintomas de apetite e alimentação	0 a 12
npitotal	pontuação total dos sintomas dos NPIs	0 a 144
scid1_depre	avaliação do humor	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid2_depre	avaliação da perda de interesse ou prazer pelas coisas que gostava	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid3_depre	avaliação do peso ou apetite	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid4_depre	avaliação do sono	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid5_depre	avaliação da agitação	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.

Continua na próxima página

Variável	Descrição	Valores
scid6_depre	avaliação da energia/cansaço	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid7_depre	avaliação sobre a utilidade/valor por coisas que fazia ou deixava de fazer	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid8_depre	avaliação sobre a concentração/decisões	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid9_depre	avaliação sobre ideação/tentativa/pensamento suicida	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid10_depre	avaliação sobre prejuízos funcionais e de relacionamentos	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid11_depre	avaliação sobre excesso de bebida e/ou drogas	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid12_depre	avaliação sobre o estado físico	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid13A_depre	avaliação de sintomas de depressão explicados por luto	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.
scid13B_depre	avaliação de sintomas de depressão por duração de tempo por luto	1 - ausente ou falso 2 - subliminar (duvidoso) 3 - limiar ou verdadeiro “?” - informação inadequada.

Continua na próxima página

Variável	Descrição	Valores
outroepdepressao	outro episódio de depressão	sim não
inicioDepre	idade do início da depressão	valor do tipo real
scoreSocEcon	pontuação sócio econômica	baixo médio alto

Algoritmos de Aprendizado de Máquina

A seguir são descritos resumidamente os algoritmos de Aprendizado de Máquina que fizeram parte dos experimentos desta pesquisa. Maiores detalhes destes algoritmos podem ser consultados em [1, 2, 11, 66]. Os algoritmos serão apresentados de acordo com o viés de aprendizagem considerado, assim abrangendo: métodos baseados em busca, métodos baseados em otimização, métodos probabilísticos, métodos baseados em distâncias e métodos *ensemble*.

Métodos Baseado em Busca

Nesta categoria, o problema de aprendizado é formulado como um problema de busca num espaço de possíveis soluções, como por exemplo, os algoritmos de Árvores de Decisão.

Árvores de Decisão (do inglês, *Decision Tree* - AD), são algoritmos de classificação de aprendizado supervisionado, onde um conjunto de decisões hierárquicas são organizadas em uma estrutura de árvore [1]. A partir da raiz da árvore, o valor de uma determinada variável independente é avaliado e decide se o próximo nó será o da direita ou o da esquerda. Esse processo é repetido até que se chegue a um nó folha que indica a classe que será dada como resultado. Uma de suas principais vantagens é que sua representação é intuitiva e de fácil compreensão.

Métodos Baseado em Otimização

Métodos de otimização realizam busca por uma hipótese que descreve os dados recorrendo à otimização de alguma função. Dessa maneira, um problema de aprendizado é formulado como um problema de otimização, o qual consiste em minimizar ou maximizar uma determinada função objetivo [11]. Os algoritmos de otimização utilizados neste estudo foram Máquina de Vetores de Suporte (do inglês, *Support Vector Machine* - SVM) e Perceptron Multicamadas (do inglês, *Multilayer Perceptron* - MLP).

SVMs, são algoritmos supervisionados que naturalmente são utilizados para tarefas de classificação binária, embora possa ser estendido para problemas multiclases

[66, 1]. Utiliza um hiperplano como separador de classes, este hiperplano é descoberto usando os vetores de suporte (conjunto de treinamento) e funciona como um suporte para o limite da decisão ao classificar.

As Redes Neurais Artificiais (RNAs) são modelo baseados no sistema nervoso humano, com o objetivo de simular a capacidade de aprendizado do cérebro humano na aquisição de conhecimento [1, 11]. O principal bloco de construção de um cérebro humano é o neurônio, também chamados de células nervosas. Os neurônios estão conectados uns aos outros por meio de sinapses, onde ocorre a transmissão de impulsos nervosos de uma célula para outra. Normalmente, a força das conexões mudam em resposta a estímulos externos [1, 11]. Há diferentes abordagens para as RNAs. O modelo perceptron é a forma mais básica de uma RNA resolvendo apenas problemas simples, contém uma camada de entrada e uma camada de saída. A função aprendida pelo modelo perceptron é um modelo linear simples. A proposta do perceptron foi evoluída para as redes do tipo MLP, apresentando uma ou mais camadas intermediárias (ou camadas ocultas) de neurônios e uma camada de saída. MLPs têm uma camada oculta, além das camadas de entrada e saída. Os nós da camada oculta podem ser conectados a diferentes tipos de topologias, que envolve padrões de conexões e grau de conectividade [11]. Por exemplo, a própria camada oculta pode consistir em várias camadas e nós em uma camada pode alimentar os nós da próxima camada [1].

Métodos Probabilísticos

Os classificadores probabilísticos constroem um modelo que quantifica a relação entre as variáveis independentes e a variável classe (variável dependente) como uma probabilidade [1]. Existem diversas metodologias probabilísticas de classificação que têm por finalidade analisar dados binários, duas delas são usadas neste trabalho, Regressão Logística e *Naïve Bayes*.

Regressão Logística (do inglês, *Logistic Regression* - LR) são modelos de aprendizado supervisionado. Os algoritmos de LR, modela diretamente as probabilidades de pertencimento à uma classe (variável dependente) em termos das variáveis (independentes) por meio de uma função discriminante [1].

Redes Bayesianas também podem ser geradas por meio de algoritmos de aprendizado supervisionado. Um classificador Bayesiano é baseado no teorema de Bayes. Esse teorema quantifica a probabilidade condicional de uma variável de classe, supondo que os atributos do conjunto de dados sejam condicionalmente independentes [1]. O algoritmo *Naïve Bayes* - (NB), é um algoritmo baseado em Redes Bayesianas. Arribas et al. [3] apresentam em seu artigo que para fins de classificação para dados médicos, é especialmente desejável o uso de classificadores bayesianos probabilísticos, pois podem fornecer

uma estimativa probabilística da decisão que está sendo tomada.

Métodos Baseados em Distância

Métodos baseados em distância são caracterizados por proximidade entre os dados na realização de previsões. Tais métodos assumem que os dados podem ser representados como pontos em um espaço euclidiano. Portanto, dados similares tendem a estar concentrados em uma região neste espaço e dados não similares estarão distantes entre si [11]. Ao contrário de métodos de aprendizagem que constroem uma descrição explícita de uma função (padrão) a partir dos exemplos de treinamento, os métodos de aprendizagem baseados em distâncias armazenam os exemplos de treinamento. A generalização acontece quando um novo exemplo deve ser classificado e um conjunto de instâncias similares são buscadas na memória e referenciadas para classificar o novo exemplo. Nesta categoria, inclui o classificador de vizinho mais próximo, *K-Nearest-Neighbor (K-NN)*.

K-NN é um algoritmo de aprendizagem supervisionado e é uma extensão do algoritmo 1-NN. No algoritmo 1-NN cada objeto representa um ponto em um espaço definido pelos atributos, denominado espaço de entrada [11]. A partir de uma métrica definida, é possível calcular a distância entre dois objetos. No K-NN é possível definir o valor de k (número de vizinhos próximos), que influencia no processo de classificação de novos dados, que são os k objetos do conjunto de treinamento mais próximos do ponto de teste x_i . Para cada ponto de teste, são obtidos k vizinhos e cada vizinho vota em uma classe, considerando um problema de classificação, o objeto de teste é classificado na classe mais votada.

Métodos *Ensemble*

Métodos *ensemble* é uma abordagem que combina os resultados de vários classificadores em vez de usar um único modelo. O aprendizado se dá por aplicar o algoritmo ao conjunto de treinamento várias vezes usando modelos diferentes ou usando o mesmo modelo em diferentes subconjuntos de treinamento de dados. Os resultados de diferentes classificadores são então combinados em uma única previsão robusta [1]. Alguns algoritmos nesta categoria utilizados neste trabalhos são descritos abaixo:

O algoritmo de Floresta Aleatória (do inglês *Random Forest -RF*), é uma coleção de diferentes Árvores de Decisão, que são vetores aleatórios independentes e identicamente distribuídos. A aleatoriedade é explicitamente inserida no processo de construção do modelo de cada árvore de decisão [1]. O treinamento é realizado nas diversas árvores de decisão criadas a partir das variáveis preditoras, por meio de uma função de custo é

possível encontrar o melhor ponto de divisão. No caso de classificação, o resultado que mais vezes foi apresentado será o escolhido.

Boosting é um meta algoritmo de aprendizado de máquina que combina um conjunto de classificadores-base “fracos” ou com baixo desempenho, criando um classificador “forte” ou com alto desempenho [18]. Este método funciona em iterações. A cada iteração um algoritmo-base é chamado para gerar um classificador simples, utilizando uma versão diferente do conjunto de dados de treinamento, estas versões são obtidas através da variação do peso associado a cada um dos exemplos, obtendo diferentes versões ponderadas do conjunto de dados. Após um número determinado de iterações, o *Boosting* combina os diversos classificadores parciais, gerando um classificador único, que possui um desempenho melhor comparado ao classificador parcial [18]. Algoritmos derivados do *Boosting* apresentam o mesmo esquema geral, embora tenham formas diferentes de ponderação do conjunto de dados em cada iteração e a forma de combinação dos classificadores parciais. O *Boosting* assume várias formas, incluindo *Adaptive Boosting (AdaBoost)*, *Gradient Boosting (GB)* e o *Extreme Gradient Boosting (XGBoost)*.

O algoritmo *AdaBoost*, inicialmente, seleciona um subconjunto de treinamento aleatoriamente e atribui peso igual a cada observação. Se a previsão estiver incorreta, então dará maior peso à observação que foi mal prevista. O conjunto de dados usado para o treinamento é continuamente adaptado para forçar o modelo a se concentrar nas amostras que são classificadas incorretamente. Além disso, os classificadores são adicionados sequencialmente, então um novo classificador impulsiona o anterior, melhorando o desempenho nas áreas onde não foi tão preciso quanto o esperado [5]. O modelo é treinado selecionando o conjunto de treinamento com base na previsão precisa do último treinamento.

O (GB) difere do *AdaBoost* em relação à maneira com a qual os modelos são treinados com relação aos anteriores. Ao invés de estabelecer pesos, o GB treina novos modelos considerando os erros dos modelos anteriores. Como o GB treina os modelos com base nos erros residuais do modelo anterior, obtém-se a predição final somando a predição de todos os modelos, que são obtidas a partir de uma versão mais corrigida da primeira predição.

O algoritmo *Extreme Gradient Boosting (XGBoost)* é considerado um aprimoramento do algoritmo GB. O XGBoost faz uso de árvores de decisão com gradiente otimizado, proporcionando velocidade e desempenho aprimorados e depende muito da velocidade computacional e do desempenho do modelo de destino. XGBoost é usado para problemas de aprendizagem supervisionada, onde os dados de treinamento faz predição para uma variável alvo (variável de classe).

Light Gradient Boosting Machine (LGBM) é uma outra variação do GB que usa algoritmos de aprendizagem baseados em árvore. O LGBM usa duas técnicas, *Gradient-*

based One Side Sampling (GOSS) e *Exclusive Feature Bundling* (EFB), que fazem o modelo funcionar de forma eficiente e fornecendo vantagem sobre outras estruturas GB. As observações incorretas terão gradientes maiores e contribuirão mais para o ganho de informações. O GOSS tem como objetivo manter cada observação com gradientes maiores e apenas descarta aleatoriamente as observações com gradientes pequenos para manter a precisão da estimativa de ganho de informação. Este tratamento pode levar a uma estimativa de ganho mais precisa do que a amostragem uniformemente aleatória, com a mesma taxa de amostragem alvo, especialmente quando o valor do ganho de informação tem uma grande faixa [2]. EFB permite agrupar dados de alta dimensão, que geralmente são esparsos, quase sem perder para reduzir o número de variáveis.