

Universidade Federal de São Carlos – UFSCar
Centro de Ciências Exatas e de Tecnologia
Programa de Pós-Graduação em Ciência da Computação
Aníbal Tomás Silva Osses

Análise da Predição da Violência Infantil por Meio de Árvores de Decisão e Regras de Associação

São Carlos - SP

Junho/2020

Aníbal Tomás Silva Osses

Análise da Predição da Violência Infantil por Meio de Árvores de Decisão e Regras de Associação

Dissertação apresentada ao Programa de Pós-Graduação em Ciência da Computação da Universidade Federal de São Carlos, como parte dos requisitos para a obtenção do título de Mestre em Ciência da Computação, área de concentração: Inteligência Artificial.

Orientador: Prof. Dr. Ricardo Augusto Souza Fernandes

São Carlos - SP

Junho/2020



UNIVERSIDADE FEDERAL DE SÃO CARLOS

Centro de Ciências Exatas e de Tecnologia
Programa de Pós-Graduação em Ciência da Computação

Folha de Aprovação

Defesa de Dissertação de Mestrado do candidato Aníbal Tomás Silva Osses, realizada em 02/06/2020.

Comissão Julgadora:

Prof. Dr. Ricardo Augusto Souza Fernandes (UFSCar)

Prof. Dr. Fábio Anderson Silva Borges (UESPI)

Prof. Dr. Vinicius Ponte Machado (UFPI)

O presente trabalho foi realizado com apoio da Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil (CAPES) - Código de Financiamento 001.

O Relatório de Defesa assinado pelos membros da Comissão Julgadora encontra-se arquivado junto ao Programa de Pós-Graduação em Ciência da Computação.

Resumo

Segundo o Fundo das Nações Unidas para a Infância (*United Nations International Children's Emergency Fund*, UNICEF), atualmente cerca de 300 milhões de crianças sofrem de diversos tipos de violência, podendo ser: psicológica, física, abuso sexual, ou negligência. Considerando a gravidade dessa problemática, adiciona-se que analisá-la torna-se difícil por alguns aspectos: existem diferentes definições de violência utilizadas em diversos contextos, não havendo padrões que facilitem o estudo; determinar qual é o tipo de violência que acontece em cada caso, podendo ser um ou mais; os informes das próprias vítimas e as estatísticas oficiais de violência não sempre possuem uma qualidade apropriada; entre outros. Ainda, cabe ressaltar que em tais levantamentos não é considerado o fato de que em muitas ocasiões as únicas pessoas que sabem que a violência acontece são as crianças e os agressores, sendo a situação invisível e tornando muito mais difícil a tarefa de previsão da agressão. O objetivo deste trabalho é gerar e avaliar modelos baseados em técnicas de aprendizado de máquina que possam estimar em que casos está acontecendo ou irá acontecer algum ato de violência infantil através de modelos que sejam representados por regras simples de entender por qualquer humano. O método científico utilizado no projeto é *ex-post-facto* usando dois conjuntos de dados estruturados, um supervisionado e outro não supervisionado, ambos construídos por organizações do Chile e que possuem atributos categóricos e numéricos. Foram aplicadas técnicas de seleção de atributos para trabalhar com os elementos mais relevantes, para depois utilizar respectivamente os algoritmos C4.5 e Apriori em cada conjunto. Avaliou-se o primeiro com as áreas sob as curvas *receiver operating characteristic* e *precision-recall*, e o segundo com as métricas *lift*, *conviction*, e *leverage*. Sobre os resultados, para a técnica de classificação treinaram-se modelos com desempenhos próximos ao 0,9 para as métricas; e para as regras de associação, em todas as execuções encontraram-se sentenças com valores superiores ao limiar para estabelecer a implicância entre seus antecedentes e consequentes.

Palavras-chave: violência infantil, violência física, abuso sexual, violência psicológica, negligência, fatores de risco, fatores de proteção, aprendizado de máquina, árvores de decisão, regras de associação.

Abstract

According to the United Nations International Children's Emergency Fund (UNICEF), currently around 300 million children around the world suffer from various types of abuse, including: psychological, physical, sexual or negligence. Considering the severity of the problem, the analysis gets more difficult given elements such as: the existence of different definitions of violence used for different contexts, without patterns that ease the study; determine the type of abuse that's happening in each case, being more than one; victims reports and official abuse statistics not always have the expected quality; among others. It should even be noted that, while collecting data, the fact that most of the time the only people that are aware of the abuse situation are the children and their aggressors, rendering the situation invisible and making the abuse's prevention a much more difficult task. The aim of this work is to generate and evaluate models based on machine learning techniques that can estimate in which cases a situation of child abuse is currently happening or it could happen, via models represented by rules easily understandable by humans. The scientific method utilized in this project is *ex-post-facto* based on two structured datasets, one supervised and the other unsupervised, both built by Chilean organizations and that possess numeric and categorical attributes. Feature selection techniques were applied in order to work with the most relevant elements, and then use the C4.5 and Apriori algorithms on each dataset respectively. The first one was evaluated with the areas under the receiver operating characteristic and precision-recall curves, and the second one with the lift, conviction and leverage metrics. About the results, for the classification technique were built models with performances close to 0.9 for each metric; and for the association rules, in all executions the sentences found have higher values than the thresholds that define the implication between their antecedents and consequents.

Keywords: child abuse, physical abuse, sexual abuse, psychological abuse, negligence, risk factors, protective factors, machine learning, artificial intelligence.

Resumen

Según el Fondo de las Naciones Unidas para la Infancia (*United Nations International Children's Emergency Fund*, UNICEF), actualmente en el mundo cerca de 300 millones de niños sufren de diversos tipos de maltrato, pudiendo ser: psicológico, físico, abuso sexual, o negligencia. Considerando la gravedad del problema, se suma que el análisis se torna difícil por elementos tales como: existencia de diferentes definiciones de violencia usadas para distintos contextos, no habiendo patrones que faciliten el estudio; determinar cuál es el tipo de maltrato que está ocurriendo en cada caso, pudiendo ser más de uno; los informes de las propias víctimas y las estadísticas oficiales de maltrato no siempre poseen la calidad esperada; entre otros. Incluso se debe destacar que en dichos levantamientos no es considerado el hecho de que en muchas ocasiones las únicas personas que saben que el maltrato está ocurriendo son los niños y sus agresores, siendo la situación invisible y haciendo mucho más difícil realizar la tarea de previsión de la agresión. El objetivo de este trabajo es generar y evaluar modelos basados en técnicas de aprendizaje de máquina que puedan estimar en qué casos está aconteciendo o acontecerá algún acto de maltrato infantil a través de modelos que sean representados por reglas simples de entender por cualquier humano. El método científico utilizado en este proyecto es *ex-post-facto* utilizando dos conjuntos de datos estructurados, uno supervisado y otro no supervisado, ambos construidos por organizaciones de Chile y que poseen atributos categóricos y numéricos. Fueron aplicadas técnicas de selección de atributos para trabajar con los elementos más relevantes, para después utilizar los algoritmos C4.5 y Apriori en cada conjuntos respectivamente. El primero fue evaluado con las áreas bajo las curvas *receiver operating characteristic* y *precision-recall*, y el segundo con las métricas *lift*, *conviction*, y *leverage*. Sobre los resultados, para la técnica de clasificación se construyeron modelos con desempeños cercanos al 0,9 para cada métrica; y para las reglas de asociación, en todas las ejecuciones se encontraron sentencias con valores superiores a los umbrales para establecer la implicancia entre sus antecedentes y consecuentes.

Palabras clave: maltrato infantil, maltrato físico, abuso sexual, maltrato psicológico, negligencia, factores de riesgo, factores de protección, aprendizaje de máquina, árboles de decisión, reglas de asociación.

Lista de Ilustrações

Figura 1 – Representação dos tipos de resultados que um teste pode fazer sobre um problema de classificação binária.	35
Figura 2 – Exemplo de curva ROC (esquerda) e curva PR (direita) com seus respectivos valores de AUC.	37
Figura 3 – Modelo de desenvolvimento ecológico da violência infantil (BEGLE; DUMAS; HANSON, 2010).	46
Figura 4 – Árvore de decisão com 2 fatores de risco (LITTLE; RIXON, 1998).	49
Figura 5 – Mapa da estimação do riscos de violência infantil: imagem da esquerda corresponde à estimação para o ano de 2014 (baseada nos dados do ano 2013) e a imagem da direita representa à estimação do ano de 2014 com os casos reais (pontos são levemente distorcidos para manter a privacidade) (DALEY et al., 2016).	52
Figura 6 – Etapas e fluxo que compõem a metodologia utilizada.	61
Figura 7 – Fluxograma da metodologia proposta à tarefa de classificação, o que dividi-se em dois caminhos com conjuntos de dados diferentes.	73
Figura 8 – Árvore de decisão obtida ao aplicar o algoritmo J48 na predição de presença de violência sobre o conjunto de dados de crianças maiores de 12 anos.	76
Figura 9 – Considerações e trabalhos feitos previa aplicação do algoritmo Apriori sobre os conjuntos de dados não supervisionados.	79

Lista de Tabelas

Tabela 1 – Classificação dos tipos de violência infantil.	24
Tabela 2 – Desempenho dos algoritmos usados para prever o risco de violência infantil de uma notificação (CHOULDECHOVA et al., 2018).	56
Tabela 3 – Resumo de artigos sobre trabalhos relacionados.	57
Tabela 4 – Continuação de resumo de artigos sobre trabalhos relacionados.	58
Tabela 5 – Resultados obtidos ao aplicar três configurações distintas do algoritmo J48 aos diferentes conjuntos de dados, especificando o número de classes e registros de cada um e destacando os melhores resultados para cada atributo preditor.	75
Tabela 6 – Quantidade total de regras conseguidas nas diferentes execuções do algoritmo Apriori sobre os conjuntos de dados não supervisionados, utilizando um suporte de 0,05 e confiança de 0,75.	80
Tabela 7 – Quantidade total de regras conseguidas com o algoritmo Apriori onde os seus valores das métricas <i>lift</i> , <i>conviction</i> , e <i>leverage</i> são maiores às medias de cada uma das execuções.	80
Tabela 8 – Regras de associação geradas com o algoritmo Apriori para cada um dos diferentes conjuntos de dados diferenciados por vitimizações em vida e ano e qual é agregação dos atributos que apresenta detalhando o suporte e a confiança, ademais das métricas utilizadas para avaliá-las.	81

Lista de Abreviaturas e Siglas

ANN	<i>Artificial Neural Network</i>
ANOVA	<i>ANalysis Of VAriance</i>
AUC	<i>Area Under Curve</i>
CAPI	<i>Child Abuse Potential Inventory</i>
CART	<i>Classification And Regression Tree</i>
CEAD	<i>Centro de Estudos y Análisis del Delito</i>
CHAID	<i>CHi-squared Automatic Interaction Detector</i>
CPS	<i>Child Protective Services</i>
CV	<i>Cross-Validation</i>
DSRS	<i>Depression Self-Rating Scale</i>
FP	<i>False Positive</i>
FN	<i>False Negative</i>
IA	Inteligência Artificial
ID3	<i>Iterative Dichotomiser 3</i>
JVQ	<i>Juvenile Victimization Questionnaire</i>
LASSO	<i>Least Absolute Shrinkage and Selection Operator</i>
MH	Média Harmônica
ML	<i>Machine Learning</i>
MLP	<i>Multilayer Perceptron</i>
NB	<i>Naive Bayes</i>
NIS	<i>National Incidence Study</i>
OCA	<i>Organismo Colaborador Acreditado</i>
OMS	Organização Mundial da Saúde

OR	<i>Odds Ratios</i>
PR	<i>Precision-Recall</i>
PRM	<i>Predictive Risk Model</i>
PSI	<i>Parenting Stress Index</i>
RF	<i>Random Forest</i>
ROC	<i>Receiver Operating Characteristic</i>
RSES	<i>Rosenberg Self-Esteem Scale</i>
RTM	<i>Risk Terrain Modeling</i>
SciELO	<i>Scientific Electronic Library Online</i>
SENAME	<i>Servicio Nacional de Menores</i>
SMOTE	<i>Synthetic Minority Over-sampling TEchnique</i>
SOM	<i>Self-Organizing Maps</i>
STAXI	<i>State-Trait Angry Expression Inventory</i>
SVM	<i>Support Vector Machine</i>
TN	<i>True Negative</i>
TP	<i>True Positive</i>
UNICEF	<i>United Nations International Children's Emergency Fund</i>
VPI	<i>Violência entre Parceiros Íntimos</i>
WEKA	<i>Waikato Environment for Knowledge Analysis</i>
XGBoost	<i>eXtreme Gradient Boosting</i>

Sumário

1	Introdução	17
1.1	Contexto da Violência Infantil	17
1.2	Problema Identificado	18
1.3	Motivação e Justificativa	19
1.4	Proposta de Pesquisa	20
1.5	Hipótese e Objetivos	20
1.6	Organização do Trabalho	21
2	Fundamentação Teórica	23
2.1	Violência Infantil	23
2.1.1	Fatores Relacionados à Violência Infantil	25
2.1.2	Polivitimização	26
2.2	<i>Machine Learning</i>	26
2.2.1	Tipos de Aprendizado	27
2.2.1.1	Supervisionado	27
2.2.1.2	Não Supervisionado	27
2.2.1.3	Semi-Supervisionado	28
2.2.1.4	Por Reforço	28
2.2.2	Pré-Processamento de Dados	29
2.2.3	Seleção de Atributos	29
2.2.3.1	Métodos de Filtro	30
2.2.3.2	Métodos <i>Wrapper</i>	31
2.2.3.3	Métodos <i>Embedded</i>	31
2.2.4	Algoritmos Utilizados no Âmbito dessa Dissertação	32
2.2.4.1	Supervisionado	32
2.2.4.2	Não Supervisionado	37
3	Revisão Bibliográfica	41
3.1	Estudos Relacionados	41
3.1.1	Técnicas de Regressão	42
3.1.2	Técnicas de Classificação	47
3.1.3	Demais Técnicas	51
3.1.4	Análises Comparativas entre Distintos Modelos	52
3.2	Resumo dos Aspectos Relevantes dos Trabalhos Relacionados	56
3.3	Posicionamento da Presente Pesquisa	59

4	Metodologia Proposta	61
4.1	Obtenção e Exploração dos Dados	62
4.1.1	Centro de Proteção de Crianças	62
4.1.2	Questionário sobre Polivitimização	65
4.2	Pré-processamento de Dados	66
4.2.1	Conjunto de Dados Utilizado para Método Supervisionado	66
4.2.2	Conjuntos de Dados Utilizado para Método Não Supervisionado	68
4.3	Seleção de Atributos	69
4.4	Geração de Modelos	70
4.4.1	Algoritmo C4.5 – Classificação	70
4.4.2	Algoritmo <i>Apriori</i> – Regras de Associação	71
4.5	Validação de Resultados	71
4.5.1	Modelo de Classificação	71
4.5.2	Modelo de Regras de Associação	72
5	Resultados e Discussões	73
5.1	Centro de Proteção de Crianças	73
5.2	Questionário sobre Polivitimização	79
6	Conclusões	85
6.1	Trabalho Realizado	85
6.2	Trabalhos Futuros	87
APÊNDICE A	Resultados	89
A.1	J48	89
A.2	Apriori	92
Referências		99

1 Introdução

Nesse capítulo são apresentados o contexto e algumas estatísticas sobre a violência infantil e o problema de sua predição, além de expor a motivação, a justificativa, e a proposta dessa pesquisa. Ademais, é descrita a hipótese e os objetivos do trabalho, terminando com um breve resumo sobre a organização desse documento.

1.1 Contexto da Violência Infantil

Em 1989, o Fundo das Nações Unidas para a Infância (*United Nations International Children's Emergency Fund*, UNICEF) estabeleceu a Convenção sobre os Direitos das Crianças, que reconhece os direitos humanos das mesmas e, diferente dos adultos, precisam de atenção e cuidados especiais. Essa convenção foi consolidada na legislatura chilena¹ e brasileira² no ano 1990, a qual possui quatro princípios que a definem: não discriminação, o interesse da criança, desenvolvimento e proteção, e a participação em decisões que a afetem (UNICEF, 2014). Quando algum desses direitos são violados, se diz que acontece a vulneração ou violação dos direitos das crianças, podendo existir violência infantil.

A UNICEF tem feito uma sumarização estatística da violência infantil no mundo destacando alguns resultados (UNICEF, 2017), os quais estão relacionados com a primeira infância (dos 0 até os 5 anos):

- 300 milhões de crianças entre 2 e 4 anos (3 a cada 4) são habitualmente vítimas de algum tipo de disciplina violenta por parte das pessoas que têm os seus cuidados (cuidadores), não necessariamente sendo os pais delas;
- 250 milhões de crianças (6 a cada 10) são castigadas fisicamente; e
- 1.100 milhões de cuidadores (1 a cada 4 aproximadamente) dizem que o castigo físico é necessário para educar adequadamente as crianças.

Outras estatísticas estão relacionadas com a violência na escola, a saber:

- 130 milhões de estudantes (1 a cada 3) sofrem de assédio escolar;
- 17 milhões de adolescentes na Europa e na América do Norte (3 a cada 10) admitem que agredem outros estudantes na escola; e

¹ www.leychile.cl/Navegar?idNorma=15824

² www.planalto.gov.br/ccivil_03/decreto/1990-1994/D99710.htm

- 732 milhões de crianças em idade escolar (1 de cada 2) vivem em países nos quais o castigo físico não é completamente proibido.

Outrossim, podem ser verificados dados de violências que relacionam-se com o abuso sexual na infância e adolescência:

- 15 milhões de mulheres adolescentes foram vítimas de relações sexuais forçadas em algum momento de suas vidas;
- 17 milhões de mulheres adultas, em 38 países subdesenvolvidos, afirmam haver tido relações sexuais forçadas na infância; e
- 2,5 milhões de mulheres jovens afirmam ter sido vítimas de violência sexual antes dos 15 anos.

Além desses dados, outra estatística importante exposta pela UNICEF está relacionada às mortes causadas entre adolescentes. Somente no ano 2015, ocorreram 82.000 homicídios de adolescentes no mundo, sendo América Latina e o Caribe onde se registraram quase a metade desses falecimentos, mesmo essa região tendo menos de 10% dos adolescentes no mundo.

1.2 Problema Identificado

Um dos principais problemas da violência infantil é quando está acontecendo ou predizer quando irá suceder. Milner (1993) afirma que não existe um modelo que indique quais são as principais variáveis que identificam o problema da violência. Além disso, os fatores de risco (elementos que determinam a presença da violência) que ele detalha (sociais, psicofisiológicos, neuropsicológicos, afetivos e comportamentais) não podem predizer a probabilidade da violência hoje, bem como nem sempre serão os melhores preditores de ataques futuros.

Gómez (2011) discute vários modelos de violência, explicando que ainda não se tem uma caracterização universal que descreva essa problemática. Porém, é dito que alguns fatores de risco podem ser mais preditivos na determinação da violência infantil, sendo esses: a repetição do abuso por gerações passadas, a juventude dos pais, o estresse, a pobreza, o uso de drogas, a superlotação na moradia, e o isolamento social. Em um estudo desenvolvido por Farren (2007) com estudantes de escolas de diferentes idades e contextos socioeconômicos, são validadas hipóteses que avaliam diversas variáveis como fatores críticos da violência infantil, sendo essas: os comportamentos problemáticos da criança, consumo de drogas ou álcool por parte dos pais, a comunicação familiar, e o nível socioeconômico. Baseado nesses estudos, é possível estabelecer que não existem variáveis

nem modelos que sempre vão prever a violência infantil, mas sim tendências de quais delas poderiam ser relevantes. Destaca-se que tal assunto será aprofundado no [Capítulo 2](#).

1.3 Motivação e Justificativa

Em diferentes contextos, profissionais de muitas áreas trabalham diretamente com crianças ou adolescentes com dificuldades diversas e, muitas das vezes, não existe uma confiança real entre o adulto e o menor para que ele fale com sinceridade de seus problemas. Porém, em muitas dessas ocasiões, o adulto tem uma noção de que algo não está bem com essa criança.

Baseado na informação exposta nas seções supramencionadas, é um fato que a problemática da violência infantil é muito séria e esta presente no mundo todo. É por isso que esse trabalho visa contribuir na melhora da situação atual da violência em crianças, mas através de uma abordagem não convencional. Serão usadas técnicas computacionais de inteligência artificial sobre conjuntos de dados relacionados à violência infantil com o intuito de criar modelos que entreguem padrões relacionados ao abuso (discutidos em maiores detalhes no [Capítulo 2](#)) e que ajudem na determinação de quando uma criança está sendo ou poderá ser violentada. Se espera com esse trabalho obter uma proposta confiável para diminuir a dificuldade que se apresenta ao trabalhar com a violência infantil, podendo auxiliar qualquer pessoa que esteja habitualmente em contato com crianças e que creia que alguma delas pode estar em perigo.

Witten (2016) menciona que o aprendizado de máquina (*Machine Learning*, ML), pode gerar um conhecimento útil ao analisar dados provenientes de algum fato de interesse. Com isso, tendo um conjunto de dados que representa um estudo de caso, é possível descobrir informações que permitirão compreender melhor alguns aspectos relevantes ou prever um resultado. Considerando áreas relacionadas a problemáticas sociais, o ML tem sido vastamente aplicado. Como exemplo se têm as aplicações na área da saúde, onde existe uma grande variedade de assuntos estudados com a aplicação de ML, como o reconhecimento de padrões sobre as atividades que a demência produz (GAYATHRI et al., 2015), a predição do consumo de álcool em jovens (LINDSAY et al., 2017), entre outros. Uma área diferente na qual ML é utilizado é na segurança pública, onde atos criminosos são preditos por meio de indicadores urbanos (ALVES et al., 2018) ou com dados gerados por telefonia móvel (BOGOMOLOV et al., 2014). A educação também tem sido vinculada com essa subárea da inteligência artificial, principalmente no desempenho escolar e acadêmico. Exemplo disso são os trabalhos feitos por Goes e Steiner (2016), onde diferentes escolas do Brasil são rotuladas de acordo com seu desempenho ou, conforme a proposta de Mohd et al. (2013), é predito a performance acadêmica de estudantes universitários. Uma das pesquisas realizadas por Rusell (2015) coleta informações sobre como a análise

preditiva tem sido aplicada em diferentes áreas sociais enfatizando a proteção infantil. Portanto, pesquisas semelhantes à proposta dessa dissertação já têm sido desenvolvidas, conforme apresentado no [Capítulo 3](#).

1.4 Proposta de Pesquisa

Em concordância com o estabelecido na [seção 1.3](#), essa pesquisa faz uso de ML para encontrar e analisar padrões que permitam ajudar na predição da violência infantil. Portanto, há a necessidade de uma grande quantidade de dados para poder criar os modelos que respondam se a criança está em um contexto de abuso ou não. Assim, para esse trabalho foram obtidos, no Chile, dois conjuntos de dados diferentes relacionados à violência infantil. O primeiro foi conseguido junto a um centro de proteção que trabalha com crianças que apresentam variados tipos de problemas (principalmente de violência, consumo de substâncias e delitos), o qual possui 9.771 registros e 30 atributos relacionados à criança e seu entorno, tendo três possíveis variáveis preditoras: se tem sofrido violência, qual tipo de violência, e quem é o principal agressor. O segundo conjunto de dados corresponde à Primeira Enquete Nacional de Polivitimização em Crianças e Adolescentes ([SUBSECRETARÍA DE PREVENCIÓN DEL DELITO, 2018](#)), que avalia a presença de 32 tipos de vitimizações na vida das crianças, sendo 19.684 as que responderam a enquete (o conceito de polivitimização será exposto em detalhes no [Capítulo 2](#)).

Considerando o fato de que modelos gerados a partir de técnicas de ML na maioria das vezes não são simples de entender para pessoas que não têm conhecimentos em computação ou estatística, buscou-se aplicar algoritmos que entreguem resultados que o senso comum seja capaz de compreender sem precisar de estudos prévios. Essas técnicas são árvores de decisão ([ROKACH; MAIMON, 2008](#)) e regras de associação ([AGRAWAL; IMIELIŃSKI; SWAMI, 1993](#)), as quais foram usadas no conjunto de dados do centro de proteção e da enquete sobre polivitimização, respectivamente. Estritamente, os modelos que esses algoritmos criam são baseados em regras simples que definem sob quais circunstâncias a criança está sofrendo violência infantil. No [Capítulo 2](#) serão definidas as ferramentas computacionais utilizadas.

1.5 Hipótese e Objetivos

Baseado no contexto da violência infantil no Chile e o problema da sua predição, além dos resultados obtidos por ML nas diferentes áreas de pesquisa, é definida a seguinte hipótese para validar-se: modelos baseados em regras gerados por meio de técnicas de ML com dados de instituições chilenas que trabalham com crianças abusadas podem entregar padrões que permitam estimar em que casos algum ato de violência infantil está acontecendo ou irá acontecer?

Sobre os objetivos, o principal propósito do estudo é analisar a utilidade que podem ter árvores de decisão e regras de associação em gerar modelos baseados em regras para a predição da violência infantil. Para estimar o desempenho dos resultados, serão utilizadas métricas pertinentes para cada um dos conjuntos e algoritmos aplicados. Como objetivo secundário espera-se detectar quais são as principais características (variáveis) do entorno das crianças que entregam maior informação para a predição da violência infantil. Ademais, destacam-se outros dois objetivos importantes relacionadas a essa pesquisa, mas que pelo alcance desse projeto são considerados como trabalho futuro:

- Validar os modelos com profissionais que estejam imersos na problemática da violência infantil, como: sociólogos, psicólogos, trabalhadores sociais, entre outros; e
- Gerar uma ferramenta simples de predição, verificando a ocorrência da violência infantil em dados que não foram empregados na construção do modelo de preditivo.

1.6 Organização do Trabalho

Este documento tem a seguinte estrutura: o [Capítulo 2](#) apresentará a fundamentação teórica dos elementos que serão usados nesse trabalho; no [Capítulo 3](#) será discutido o estado da arte de como o ML e a estatística têm sido usados para dar respostas a diferentes problemas relacionados à predição da violência infantil; o [Capítulo 4](#) é destinado a definir a metodologia proposta da dissertação, em que serão discutidas as características presentes nos conjuntos de dados, bem como as configurações usadas para gerar os modelos e as formas em que foram avaliados; no [Capítulo 5](#) serão apresentados os resultados obtidos e as discussões que podem ser geradas a partir desses; e por último, o [Capítulo 6](#) expõe as conclusões de todo o estudo, além de possíveis trabalhos futuros.

2 Fundamentação Teórica

Nesse capítulo são expostos os elementos fundamentais da violência que acomete crianças, bem como a definição, os fatores que se relacionam a ela e o conceito de polivitimização. Ademais, ainda são apresentados os aspectos que fundamentam, de forma geral, os algoritmos de ML utilizados e as técnicas que visam a validação dos resultados.

2.1 Violência Infantil

A UNICEF (UNICEF, 2005) define a violência infantil como o ato ocasional ou habitual de violência física, sexual ou emocional em crianças e adolescentes entre 0 e 18 anos, provocado por um familiar ou alguma pessoa do entorno. Esses elementos são definidos como:

- Violência física – agressão física que pode ou não ter como resultado uma lesão;
- Violência emocional – divide-se em dois, assédio verbal (insultos, críticas negativas ou ridicularizar) e indiferença ou rechaço para a criança; e
- Abuso sexual – atividade sexual entre um adulto e uma criança.

Além dessas classificações, no documento menciona-se a negligência, fato que acontece quando as necessidades básicas das crianças não são satisfeitas por seus cuidadores.

No Chile, uma instituição que encarrega-se das crianças que sofrem algum tipo de violência é o *Servicio Nacional de Menores* (SENAME), o que define a violência infantil como toda ação que prejudique o desenvolvimento biológico, psicológico e social da criança (SENAME, 2016). Além disso, considera-se violência deixar de realizar alguma ação que irá prejudicar a criança. A classificação dos tipos de violência que o SENAME define são:

- Abuso sexual – atividade sexual que um adulto realiza em uma criança, onde o adulto tem uma posição de poder e a criança não possui a capacidade de parar a situação;
- Violência física – dano provocado à criança que põe em perigo sua integridade física, seja pela força (ativo) ou por um descuido casual ou intencionado (passivo); e
- Violência psicológica – atos de adultos que prejudicam a autoestima das crianças (como ameaças ou insultos), ou também não dar afeto ou rejeitá-las.

Outra instituição do Chile que pertence ao Ministério de Desenvolvimento Social é o *Crece Contigo*, que utiliza a mesma definição de violência infantil que a UNICEF (CRECECONTIGO, 2015). Portanto é possível estabelecer que diferentes organizações fazem uso de classificações semelhantes para o problema da violência infantil, tanto em contextos nacionais como internacionais.

Baseado nas definições apresentadas, a Tabela 1 sumariza e estrutura os tipos de violência infantil supracitados.

Tabela 1 – Classificação dos tipos de violência infantil.

Tipo de Violência	Ativo	Passivo (Negligência)
Física	Violência física	Abandono físico
Emocional (Psicológica)	Violência emocional	Abandono emocional
Sexual	Abuso sexual	Não definido

Fonte: Elaborado pelo autor.

Além do fato do sofrimento que a violência infantil produz nas crianças quando a experimentam, existem múltiplos tipos de consequências que conseguem ser geradas pelos abusos, podendo ser classificadas em 3 grupos diferentes: físicas, psicológicas e do comportamento (CHILD WELFARE INFORMATION GATEWAY, 2013). Alguns desses efeitos negativos são:

- Determinadas partes do cérebro deixem de funcionar corretamente, com consequências cognitivas, socioemocionais e de saúde mental (CHILD WELFARE INFORMATION GATEWAY, 2001);
- Jovens adultos que tenham sofrido de violência podem desenvolver depressão, ansiedade, transtornos alimentares e tentativas de suicídio (SILVERMAN; REINHERZ; GIACONIA, 1996);
- Incremento do risco de abusar de substâncias como drogas e álcool (FELITTI et al., 1998);
- Maior possibilidade de participar em problemas relacionados com delinquência, gravidez na adolescência e baixo desempenho acadêmico (B. THORNBERRY T., 1997);
- Diminuição da capacidade de ter e manter relacionamentos íntimos saudáveis na idade adulta (COLMAN; WIDOM, 2004); entre outras consequências.

De forma geral e considerando as possíveis consequências da violência infantil, pode-se dizer que cria um trauma complexo nas crianças (LECANNELIER, 2018), o qual refere-se à “experiência de ter sofrido múltiplas traumatizações¹ especificamente de

¹ Entende-se traumatização como o ato ou efeito de traumatizar ou gerar um trauma.

origem interpessoal, com resultados nefastos para o desenvolvimento da criança” (VAN DER KOLK, 2005).

Para entender como se produz ou se evita a violência infantil, existem alguns fatores que aumentam ou diminuem a probabilidade de que aconteça, os quais são discutidos na seção (CHILD WELFARE INFORMATION GATEWAY, 2004). Posteriormente, é apresentada a definição do conceito polivitimização (FINKELHOR; ORMROD; TURNER, 2007), o que é necessário para entender um dos conjuntos de dados usados no trabalho.

2.1.1 Fatores Relacionados à Violência Infantil

Existem dois tipos de elementos que podem estar relacionados com a violência infantil, os fatores de risco e os fatores de proteção.

Os fatores de risco referem-se às circunstâncias estressantes na vida da criança ou do seu entorno, que acrescentam a possibilidade de que aconteça a violência infantil (CHILD WELFARE INFORMATION GATEWAY, 2014). É difícil definir quais são os fatores de risco de forma exata, mas já foram feitos vários estudos sobre os principais elementos que geram a violência infantil, sempre convergindo em componentes semelhantes. Os primeiros desses trabalhos (MCDONALD; MARKS, 1991) e outros mais recentes (STITH et al., 2009) destacam diversos fatores, porém agrupados nas mesmas categorias: características da criança, interação entre os cuidadores e a criança, características dos cuidadores, e características da família e do entorno. Instituições como a Organização Mundial da Saúde (OMS) descrevem modelos de fatores de risco com atributos parecidos, a saber: características da criança e de seu relacionamento (principalmente com a família), fatores comunitários e fatores sociais (BUTCHART; HARVEY, 2009). Organizações como o Departamento de Saúde e Serviços Humanos (*Department of Health and Human Services*) dos Estados Unidos divide em três grupos os fatores de risco: da criança, dos pais e família, e do entorno e sociais (CHILD WELFARE INFORMATION GATEWAY, 2004).

Sobre os fatores de proteção, pode-se estabelecer que correspondem a elementos que mitigam o risco da presença da violência infantil e promovem o desenvolvimento saudável da criança, podendo estar presentes no contexto individual, familiar, da comunidade ou da sociedade; os que não têm sido estudados com a mesma intensidade que os fatores de risco (CHILD WELFARE INFORMATION GATEWAY, 2014). Nos anos 90, destacaram-se alguns fatores relacionados principalmente aos pais (MASTEN; GARMEZY, 1985), para que pesquisas posteriores expandissem a elementos como: comunidade e suporte social, escola, colegas, família, e aspetos individuais da criança (DURLAK, 1998). Estudos mais recentes continuam validando fatores semelhantes na prevenção da violência infantil (AFIFI; MACMILLAN, 2011), que são reafirmados pela OMS (BUTCHART; HARVEY, 2009).

2.1.2 Polivitimização

O termo polivitimização foi estabelecido por Finkelhor et al. (2007) e refere-se ao fato da criança sofrer múltiplos tipos de vitimizações (é importante destacar que o conceito “vitimizar” do português é diferente de *victimize* do inglês, o que corresponde a quando uma pessoa faz uma má ação em outra²). Nessa pesquisa foram coletados os dados de 2.030 crianças por meio da enquete *Juvenile Victimization Questionnaire* (JVQ) (FINKELHOR et al., 2005), que é composta por 34 perguntas que avaliam a presença das vitimizações com a sua respectiva cronicidade, as quais são divididas em 4 grupos: crimes violentos e de propriedade; violações do bem-estar da criança; violência da guerra e distúrbios civis; e *bullying*. O questionário também tem perguntas que avaliam alguns aspectos da saúde mental das crianças, especificamente: ansiedade, sintomas depressivos e raiva/agressão. Ademais, estabeleceu-se que quando o número de tipos vitimizações que sofre uma criança for maior que 4, existe o fenômeno da polivitimização, a qual ainda é classificada em baixa (entre 4 e 6) ou alta (maior que 6) (FINKELHOR; ORMROD; TURNER, 2007).

Sobre a importância da polivitimização, Finkelhor et al. (2011) aprofundam nesse assunto. Pesquisas anteriores também referem-se aos efeitos negativos que produz a acumulação de diversos tipos de violência ao longo de tempo, especialmente em crianças (DONG et al., 2004). Ademais, tem-se estudado o fato de que as vitimizações não são eventos isolados, pelo que a maioria dos casos de violência serão de polivitimização (FINKELHOR; ORMROD; TURNER, 2007). Além dessas informações, pode-se complementar que as crianças que tenham experimentado por muito tempo diferentes tipos de vitimizações podem desenvolver um trauma complexo (COOK et al., 2005).

2.2 *Machine Learning*

Muitas vezes tenta-se analisar ou resolver problemas que estão associados ao fato de encontrar relações entre diferentes características sobre um conjunto de dados, o que pode ser difícil de conseguir de acordo com a natureza do assunto. Assim, ML é aplicado com sucesso para resolver esse tipo de problema (MAGLOGIANNIS ILIAS, 2007). O conceito de “aprendizado” está relacionado com o fato de empregar dados históricos para encontrar padrões que entreguem informação a comportamentos presentes e futuros (WITTEN et al., 2016). Na prática, ML é uma subárea da inteligência artificial que constrói modelos matemáticos mediante um conjunto de dados que descrevem um fenômeno (o que é chamado de conjunto de treinamento), tentando fazer previsões de algum assunto de interesse (BISHOP, 2006).

² <https://dictionary.cambridge.org/dictionary/english/victimize>

2.2.1 Tipos de Aprendizado

Em ML existem quatro principais tipos de aprendizado, a saber: supervisionado, não supervisionado, semi-supervisionado e por reforço. A escolha de qual usar é feita baseada nas propriedades do conjunto de dados e na natureza do problema (SHAVLIK et al., 1990).

2.2.1.1 Supervisionado

Algoritmos supervisionados criam modelos a partir de dados que contêm as entradas e saídas esperadas (rotulados). Assim, quando o modelo gerado a partir desses dados estiver pronto, permitirá prever a saída de um novo registro que não foi usado antes e que não tem uma saída conhecida (RUSSELL; NORVIG, 2016). As técnicas supervisionadas podem se dividir em 2 grandes grupos, de classificação e regressão (ALPAYDIN, 2009).

- Algoritmos de Classificação – São usados principalmente quando se tem uma coleção limitada de saídas (categorias) e precisa-se atribuir uma delas ao novo registro, valor que corresponde à predição (ALPAYDIN, 2009). Alguns dos principais algoritmos de classificação são: árvores de decisão (*Decision Tree*) (ROKACH; MAIMON, 2008), redes neurais artificiais (*Artificial Neural Network*, ANN) (ROJAS, 2013), *Support Vector Machine* (SVM) (WANG, 2005), entre outras técnicas.
- Algoritmos de Regressão – Tais algoritmos são usados quando a saída corresponde a um valor numérico contínuo, e têm como objetivo calcular o valor de saída de uma nova instância (ALPAYDIN, 2009). Ainda assim, existem variedades dessa técnica que permitem trabalhar com dados categorizados, como a regressão logística (PAMPEL, 2000).

2.2.1.2 Não Supervisionado

Nesse tipo de aprendizado não se tem um valor final com o qual comparar o resultado, tendo como objetivo detectar as propriedades ou padrões de um conjunto de dados sem um “supervisor” (não rotulados). Duas das principais técnicas de aprendizado não supervisionado são o *clustering* (agrupamento) e as regras de associação (HASTIE; TIBSHIRANI; FRIEDMAN, 2009).

- Algoritmos de Agrupamento – Métodos nos quais são construídos *clusters* ou grupos de objetos (registros de dados), em que os elementos de um mesmo conjunto possuem características muito semelhantes entre eles, mas diferem das de outros grupos (GAN; MA; WU, 2007). Algumas das técnicas mais populares de *clustering* são baseadas

nos centróides dos agrupamentos, como o *k-means* (WU, 2012) ou os mapas auto-organizados (*Self-Organizing Maps*, SOM) (KOHONEN, 1990).

- Regras de Associação – Permitem descobrir fortes relações (regras) entre elementos de um conjunto de dados, detectando principalmente associações entre seus atributos (AGRAWAL; IMIELIŃSKI; SWAMI, 1993). O principal algoritmo desse tipo de aprendizado de máquina é o *Apriori* (AGRAWAL; SRIKANT et al., 1994).

2.2.1.3 Semi-Supervisionado

As técnicas de aprendizado semi-supervisionado fazem uso de dados rotulados e não rotulados, tendo como objetivo a classificação desses últimos para obter modelos com melhores resultados do que aqueles que apenas fazem uso de dados rotulados (SØGAARD, 2013). Existem diferentes tipos de algoritmos de aprendizado semi-supervisionado, os quais podem ser agrupados em: *wrapper* (LIU et al., 2016), de agrupamento (*clustering*) (BAIR, 2013), e do vizinho mais próximo (*nearest neighbor*) (NISBET; ELDER; MINER, 2009).

- *Wrapper* – A finalidade desses métodos é gerar modelos supervisionados a partir de conjuntos de dados parcialmente rotulados (LIU et al., 2016). Alguns dos mais populares são: *self-training*, *co-training*, e *expectation maximization* (CHAPELLE; SCHOLKOPF; ZIEN, 2009).
- Agrupamento – Conforme definido no aprendizado não supervisionado na [subseção 2.2.1.2](#), as técnicas de *clustering* criam grupos de instâncias que possuem características semelhantes. No caso do aprendizado semi-supervisionado, isso é feito a partir de um conjunto de dados que têm os agrupamentos parcialmente atribuídos (BAIR, 2013).
- *Nearest Neighbor* – Outro conjunto de técnicas de aprendizado semi-supervisionado são as que fazem uso dos algoritmos de busca do vizinho mais próximo (NISBET; ELDER; MINER, 2009), das quais pode-se destacar *label propagation* e algoritmos *editing* e *condensed* (SØGAARD, 2013).

2.2.1.4 Por Reforço

O aprendizado por reforço tem como objetivo definir que ação deve fazer um agente (software) em um ambiente (problema) determinado, tentando maximizar algum ganho estabelecido. A diferença dos algoritmos supervisionados, é que neste tipo de aprendizado não existem os elementos de entrada e saída. Portanto, as ações que o agente execute podem não ser completamente corretas ou errôneas (JAKSCH; ORTNER; AUER, 2010). O aprendizado por reforço não possui técnicas bem definidas como os outros tipos, porém

é possível destacar o algoritmo *q-learning* (SUTTON; BARTO, 1998) como a principal técnica e que dá origem a outras ferramentas nesse tipo de aprendizado.

2.2.2 Pré-Processamento de Dados

O pré-processamento ou preparação dos dados corresponde ao processo pelo qual os registros são trabalhados para que as ferramentas de ML possam ser usadas mais facilmente sobre eles (PYLE, 1999). Não é um procedimento fixo e pode ser ou não composto por diferentes etapas, as que são agrupadas nas seguintes categorias (GARCÍA; LUENGO; HERRERA, 2015).

- Limpeza de Dados – Corresponde a operações de correção de dados errados, eliminação de registros incorretos ou *outliers* (valores atípicos), redução de elementos desnecessários, ou qualquer tipo de ruído que possa ter o conjunto (KIM et al., 2003).
- Transformação de Dados – Etapa em que os dados são modificados de forma apropriada para a extração de informação. Pode realizar-se construção de novos atributos baseados em outros já existentes, sumarização ou agregação de dados, normalização, discretização de variáveis, entre outras tarefas (LIN, 2002).
- Integração de Dados – Processo em que são unidas diferentes fontes de dados procurando ter informação homogênea. É efetuado principalmente através da unificação de variáveis (DETOURS et al., 2003).
- Complementação de Valores – Completar registros vazios ou errados de forma intuitiva, ou estimando valores que serão melhores que deixar o campo em branco (ALLISON, 2001).

Às vezes, nos problemas de classificação, pode ser necessário equilibrar o número de instâncias de cada classe da coleção de dados para obter resultados mais confiáveis, tendo que aplicar os processos de *undersampling* e *oversampling* (HE; GARCIA, 2008). O primeiro consiste em selecionar apenas alguns registros do conjunto para a análise. Nesse caso, buscando que exista um número semelhante de instâncias de cada uma das classes do problema. O segundo corresponde ao processo inverso, ou seja, a partir dos elementos iniciais criam-se outros permitindo de igual forma equilibrar o conjunto de dados (YAP et al., 2014).

2.2.3 Seleção de Atributos

Esse procedimento tenta escolher um subconjunto de atributos da coleção inicial que se tem de variáveis. Alguns dos fatores pelos quais é necessário fazer esse trabalho

são (LEE; VERLEYSEN, 2007):

- O tempo de execução dos algoritmos diminui;
- Simplificação dos modelos gerados para sua interpretação;
- Redução do sobre-ajuste nos modelos (*overfitting*) (HAWKINS, 2004);
- Evitar a presença da “maldição da dimensionalidade” (*curse of dimensionality*) (BELLMAN, 2013); entre outros.

O conceito de *overfitting* refere-se ao treino excessivo de um algoritmo de ML para um conjunto de dados determinado, implicando que quando é avaliado com outros dados não conseguem-se os resultados corretos. Pode-se dizer que um modelo sobre-ajustado é o oposto de um modelo geral, o qual tem um desempenho adequado tanto para os dados de treino como para dados de validação ou de teste (HAWKINS, 2004). Sobre o conceito da maldição da dimensionalidade, isso corresponde ao fato de ter demasiadas dimensões (atributos) em um conjunto de dados, o que produz dois problemas principais: (1) o número de combinações possíveis das instâncias é excessivamente grande, pelo qual precisam-se uma enorme quantidade de registros para poder estimar todos os resultados possíveis; e (2) que em alguns casos o tempo de processamento torna-se inviável (LEE; VERLEYSEN, 2007).

Existem diferentes técnicas para realizar a seleção de atributos, que podem ser agrupadas em três principais categorias: de filtro, métodos *wrapper*, e métodos *embedded* (CHANDRASHEKAR; SAHIN, 2014).

2.2.3.1 Métodos de Filtro

Essas técnicas avaliam a correlação dos diferentes atributos de um conjunto com a variável que representa a classe através de algum teste estatístico (GUYON et al., 2008), como por exemplo: correlação de Pearson (BENESTY et al., 2009), qui-quadrado (MCHUGH, 2013) e ANOVA (*ANalysis Of VAriance*) (RUTHERFORD, 2001). Baseado nos resultados do teste aplicado, é que são escolhidos os atributos que têm melhor desempenho ou que possuem maior relação com a variável preditora.

As principais vantagens desses métodos são que não precisam de um algoritmo de ML para ser aplicados, e que diminuem a presença de *overfitting* nos modelos gerados depois do treinamento. Sobre as desvantagens, pode-se mencionar que esses métodos fazem apenas uma análise univariada, pela qual não consideram o fato de como um grupo de características pode estar relacionado com o valor a ser predito, além de não identificar a colinearidade entre os atributos para eliminar variáveis redundantes (GUYON et al., 2008).

2.2.3.2 Métodos *Wrapper*

Criam diferentes sub-conjuntos de características a partir dos dados iniciais para posteriormente ser avaliados com alguma técnica de ML, tentando encontrar quais são os atributos que permitem gerar o modelo com o melhor desempenho. Existem dois mecanismos principais para gerar os sub-conjuntos, *feature selection* e *feature extraction* (OLVERA-LÓPEZ et al., 2010).

- *Feature Selection* (Seleção de Características) – Começa-se com um sub-conjunto vazio e agrega-se um atributo de cada vez, para depois criar um modelo com esses elementos. Em cada iteração o modelo treinado é avaliado, e o processo é feito até que ao adicionar uma nova variável o desempenho do algoritmo não melhora (LIU; MOTODA, 2012).
- *Feature Extraction* (Extração de Características) – Começa-se com o conjunto de todos os atributos e é removido um atributo de cada vez, para depois criar um modelo com o novo sub-conjunto. Em cada iteração o modelo treinado é avaliado, e o processo é feito até que ao remover uma nova variável o desempenho do algoritmo não melhora (COELHO; RICHERT, 2015).

A principal vantagem desse tipo de técnica (ao contrário dos métodos de filtro) é que permite avaliar o desempenho de vários atributos ao mesmo tempo. De maneira oposta, considerando que a cada vez que é criado um novo sub-conjunto de atributos precisa-se treinar um modelo, a principal desvantagem dos métodos *wrapper* é o tempo de processamento que necessitam para obter um resultado, fato que pode tornar inviável a aplicação dessas técnicas (OLVERA-LÓPEZ et al., 2010).

2.2.3.3 Métodos *Embedded*

Elas selecionam quais são os atributos que entregaram os melhores resultados enquanto o modelo é criado (CHANDRASHEKAR; SAHIN, 2014). Alguns dos principais métodos *embedded* são: LASSO (*Least Absolute Shrinkage and Selection Operator*) (TIBSHIRANI, 1996), *ridge regression* (MARQUARDT; SNEE, 1975), *elastic net* (ZOU; HASTIE, 2005), CART (*Classification And Regression Trees*) (GUYON; ELISSEEFF, 2003), entre outros.

Considerando a natureza dessas técnicas, as principais vantagens que têm é que permitem a interação entre os atributos e o modelo, e computacionalmente são menos custosos que os métodos *wrapper*. Sobre desvantagens, pode-se mencionar que conceitualmente são mais complexos pelo fato de ter juntos os processos de seleção e treinamento (MALDONADO; WEBER, 2012).

2.2.4 Algoritmos Utilizados no Âmbito dessa Dissertação

Para esse trabalho foram usados dois conjuntos de dados, um supervisionado e outro não supervisionado, conseqüentemente precisam-se técnicas de ML adequadas para cada um (ver [subseção 2.2.1](#)). Em seguida, serão apresentados os algoritmos escolhidos para a análise dessas coleções de dados, além das técnicas utilizadas para a validação dos modelos treinados.

2.2.4.1 Supervisionado

O conjunto de dados rotulados que se tem, permite classificar três fatos diferentes: (1) se a criança sofre de algum tipo de violência; (2) qual é o tipo de violência que a criança sofre; e (3) quem é o principal agressor desse abuso. Ademais, buscando um procedimento metodológico que seja compreensível a pessoas leigas em inteligência artificial, porém trabalham com crianças, optou-se pelas árvores de decisão ([ROKACH; MAIMON, 2008](#)). Essa escolha foi motivada pela facilidade que uma árvore de decisão fornece, em termos de compreensão, como modelo de classificação.

Árvores de decisão (ou de classificação) são usadas para catalogar uma instância ou objeto dentro de um grupo pré-definido de classes com base nos atributos que ele tem, sendo uma técnica útil na exploração dos conjuntos de dados ([ROKACH; MAIMON, 2008](#)). Algumas das principais vantagens das árvores de decisão sobre outras técnicas de ML são ([JAMES et al., 2013](#)):

- A simplicidade para que qualquer pessoa possa compreender eles como modelo preditivo através das regras que o compõem;
- Podem utilizar conjuntos de dados combinados por variáveis numéricas e categóricas;
- Não precisa de uma preparação dos dados elaborada;
- Podem lidar facilmente com preditores qualitativos sem a necessidade de criar variáveis de difícil compreensão.

Uma árvore de decisão é formada por um conjunto de nós unidos direcionadamente, distinguindo 3 elementos diferentes: a raiz, os nós internos e as folhas. A raiz corresponde a um nó único que não possui vértices de entrada e pode-se entender como o nó inicial da árvore; os nós internos possuem vértices de entrada e saída e dividem o espaço em dois ou mais sub-espacos de acordo com alguma função discreta aplicada às instâncias do conjunto de dados; os nós folhas (também chamados terminais) somente têm vértices de entrada e são eles os que atribuem uma classe à instância ([ROKACH; MAIMON, 2008](#)). Para exemplificar essa estrutura, na [Figura 4](#) o nó “habilidades do cuidador” é a raiz da árvore, os dois nós “atitudes do cuidador” são internos e os demais nós correspondem às

folhas, que classificam o nível de risco que as crianças têm de sofrer violência. Além disso, todos os vértices têm valores que determinam a continuidade da árvore. Por exemplo, se as habilidades do cuidador fossem positivas, o risco será pequeno, mas se fossem negativas e as atitudes do cuidador muito negativas, o risco será elevado. Existem múltiplas variações de árvores de decisão, como: ID3 (*Iterative Dichotomiser 3*) (QUINLAN, 1979), CART (BREIMAN, 2017), CHAID (*CHi-squared Automatic Interaction Detector*) (KASS, 1980), entre outras; mas considerando a natureza mista do conjunto de dados usados nessa parte da pesquisa (categóricos e numéricos), escolheu-se o algoritmo C4.5 (QUINLAN, 1993) para criar o modelo preditivo da violência infantil.

C4.5 é uma técnica de ML que permite gerar árvores de decisão para resolver problemas de classificação. Algumas características importantes desse algoritmo são: permite trabalhar com atributos que tenham campos vazios, permite trabalhar com variáveis contínuas e categóricas, e utiliza um processo de poda para substituir galhos que não aportam na precisão final de classificação por nós folha (QUINLAN, 1993). Sobre seu funcionamento, C4.5 faz uso da métrica *information gain* (ou ganho de informação) para determinar qual atributo é o que melhor divide o conjunto de dados em cada um dos nós da árvore. O conceito de informação refere-se à pureza de um grupo de instâncias; no caso de árvores de decisão, a pureza pondera o fato de que um nó esteja representado por apenas uma classe, ou seja, quando todos os elementos de um nó são da mesma categoria se diz que é puro (QUINLAN, 1986). Para calcular o valor de *information gain* precisa-se antes obter a entropia (H) (SHANNON, 1948) dos dados, que corresponde à imprevisibilidade do estado dos registros (pode-se entender como a medida oposta à pureza). Considerando que se tem um conjunto D com alguma variável aleatória com possíveis estados $\{d_1, d_2, \dots, d_n\}$, a entropia é calculada por:

$$H(D) = - \sum_{i=1}^n P(d_i) \log_b P(d_i),$$

onde $P(d_i)$ é a probabilidade de que a variável tenha o valor d_i e b é a base do logaritmo. Assim, o valor da métrica *information gain* corresponde a uma variação da entropia calculada antes e depois de fazer a divisão de um nó da árvore, o que pode-se representar por:

$$\text{information_gain}(D, A) = H(D) - \sum_{i=1}^m P(a_i) H(a_i),$$

em que A é um atributo qualquer com valores $\{a_1, a_2, \dots, a_m\}$ (QUINLAN, 1986).

O Algoritmo 1 apresenta o pseudo-código do procedimento para obter uma árvore de decisão através da técnica C4.5. Na linha 2 mencionam-se os critérios de parada do algoritmo, que corresponde a quando a árvore não tem que se dividir mais, o que acontece quando: todas as instâncias são da mesma classe ou nenhuma das características do atributo proporciona ganho de informação (QUINLAN, 1993). O ciclo entre as linhas 3 e 4 refere-se ao cálculo da métrica *information gain* para todos os atributos A do conjunto D , para

depois na variável A_{melhor} calcular qual é o que entrega a maior informação. Já nas linhas 7 e 8, cria-se um nó $Node$ que divide os dados através do atributo A_{melhor} e adiciona-se à árvore $Tree$. Em seguida divide D em subconjuntos com base em $Node$ e os guarda em S . Por último, cada um dos elementos s de S , de forma recursiva, entra no procedimento C4.5 para encontrar a sub-árvore $Tree_{sub}$ e adicioná-la à árvore original $Tree$.

Algoritmo 1 - Pseudo-código do procedimento C4.5(D) para gerar árvores de decisão (QUINLAN, 1993).

Input: D ▷ D : conjunto de dados.

- 1: $Tree \leftarrow \{\}$; ▷ $Tree$: conjunto de nós que representam a árvore.
- 2: **if** D não cumpre algum dos critérios de parada **then**
- 3: **for all** $A \in D$ **do** ▷ A : atributo do conjunto de dados.
- 4: | $ig_A \leftarrow information_gain(D, A)$; ▷ ig_A : valor *information gain* de A .
- 5: **end for**
- 6: $A_{melhor} \leftarrow max(ig_A)$; ▷ A_{melhor} : atributo com melhor ganho de informação.
- 7: $Node \leftarrow$ Nó com A_{melhor} como divisor; ▷ $Node$: novo nó de decisão.
- 8: $Tree \leftarrow Tree \cup Node$;
- 9: $S \leftarrow$ Sub-conjuntos de D baseados em $Node$;
- 10: **for all** $s \in S$ **do**
- 11: | $Tree_{sub} \leftarrow C4.5(s)$;
- 12: | $Tree \leftarrow Tree \cup Tree_{sub}$;
- 13: **end for**
- 14: **return** $Tree_{sub}$;
- 15: **end if**

Sobre a validação dos algoritmos, frequentemente as técnicas de classificação fazem uso do método *holdout* (ou alguma das suas variações) na criação dos modelos. *Holdout* divide aleatoriamente a amostra de dados em dois subconjuntos, um com o qual os modelos são criados (treinamento) e outro para estimar o erro (validação), sendo a maioria das vezes segmentado em porcentagens de 70 e 30 ou 80 e 20 registros (BLATTBERG; KIM; NESLIN, 2008). Outra técnica que surge desse procedimento é *cross-validation* ou validação cruzada, destacando duas variações. A primeira é o *k-fold*, a qual tem o mesmo funcionamento que *holdout*, mas faz o processo k vezes (com k um número definido pelo usuário), considerando todo o conjunto de dados no final para depois estimar o desempenho geral de todas as iterações; e a segunda é *leave-one-out*, que itera tantas vezes o número de registros existentes no conjunto de dados, tomando apenas um deles em cada iteração para a validação, enquanto os demais são usados para o treinamento do modelo (KOHAVI et al., 1995). Para esse estudo foi selecionado o método *cross-validation* porque considera-se apropriado ter a flexibilidade de selecionar o número de iterações que serão executadas para gerar os modelos.

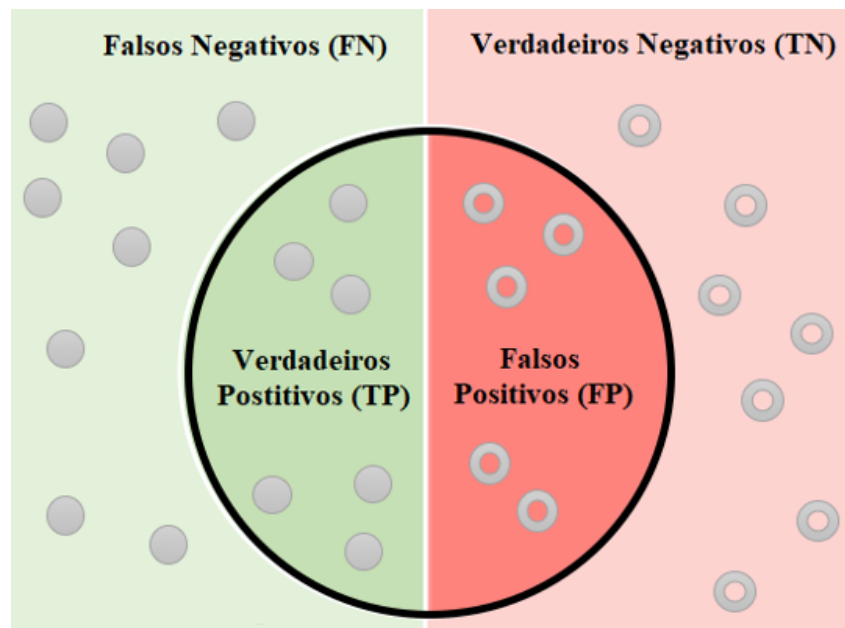
Quando um modelo já está pronto, precisa ser avaliado o desempenho que tem. No caso das técnicas de classificação, portanto saber se é competente na tarefa de detectar a classe correta de uma nova instância do problema. Isso é possível através de três métricas

frequentemente usadas em classificação: precisão, sensibilidade, e especificidade. Essas métricas são definidas por meio de outros quatro elementos (ALTMAN; BLAND, 1994):

- Verdadeiro Positivo (*True Positive*, TP) – Instância que foi classificada corretamente como positiva ou no grupo correto;
- Verdadeiro Negativo (*True Negative*, TN) – Instância que foi classificada corretamente como negativa ou que não pertence a um grupo específico;
- Falso Positivo (*False Positive*, FP) – Instância que foi classificada incorretamente como positiva em um grupo que não corresponde a sua classe; e
- Falso Negativo (*False Negative*, FN) – Instância que foi classificada incorretamente como negativa em um grupo que corresponde a sua classe.

A Figura 1 apresenta um esquema dos tipos de classificações possíveis: na faixa verde estão as instâncias que são realmente positivas e na rosa as negativas, enquanto as que ficam dentro do círculo preto correspondem às instâncias que o modelo classificou como positivas, sejam verdadeiras ou falsas.

Figura 1 – Representação dos tipos de resultados que um teste pode fazer sobre um problema de classificação binária.



Fonte: Adaptado de <https://transparint.com/blog/2016/04/01/false-negatives-a-serious-danger-in-your-aml-program>.

Portanto, as métricas precisão, sensibilidade, e especificidade são definidas como (ALTMAN; BLAND, 1994):

- Precisão – Taxa dos valores que o modelo classificou corretamente como positivos sobre o total dos que foram preditos como positivos, definida como:

$$\frac{TP}{TP + FP};$$

- Sensibilidade – Taxa dos valores que o modelo classificou corretamente como positivos sobre o total de valores positivos, definida como:

$$\frac{TP}{TP + FN};$$

- Especificidade – Taxa dos valores que o modelo classificou corretamente como negativos sobre o total de valores negativos, definida como:

$$\frac{TN}{TN + FP}.$$

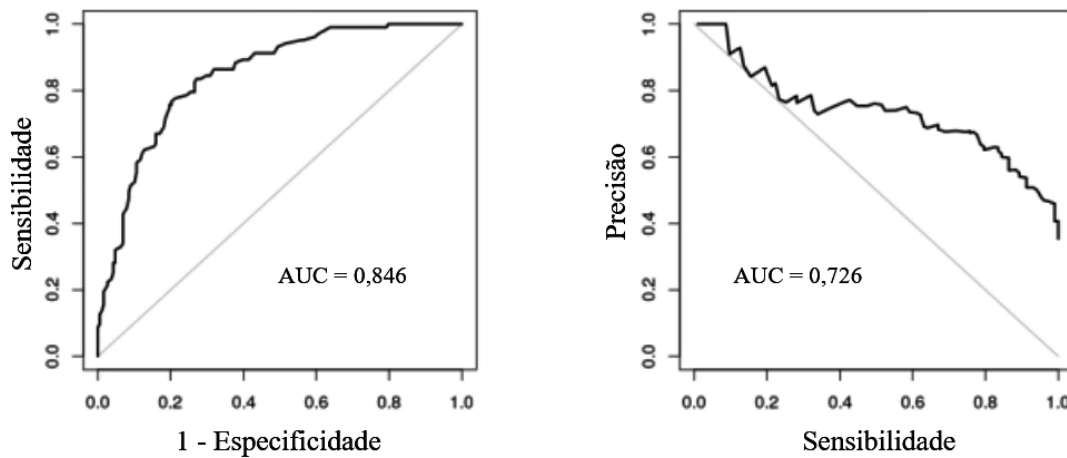
Baseado nas métricas previamente definidas, existem outras técnicas que permitem avaliar o desempenho dos modelos, como as curvas *Receiver Operating Characteristic* (ROC) e *Precision-Recall* (PR) (DAVIS; GOADRICH, 2006). O primeiro (comumente conhecido como curva ROC) faz a análise gráfica para diferentes limiares que pode ter o modelo da sensibilidade contra o complemento da especificidade ($1 - \text{especificidade}$) ou a taxa de FP para conseguir calcular a área sob a curva (*Area Under Curve*, AUC), número que se encontra no intervalo $[0, 1]$. Quanto maior for esse número, o modelo terá uma maior capacidade de diferenciar elementos positivos de negativos (FAWCETT, 2006). O método PR é semelhante à curva ROC, basicamente porque para diferentes limiares também contrasta a precisão contra a sensibilidade, obtendo um gráfico similar à curva ROC. Sua AUC também está contida no intervalo $[0, 1]$, e quanto maior for, significa que tanto a precisão como a sensibilidade do modelo têm valores elevados. Entretanto, quando um desses elementos é elevado, normalmente o outro é baixo. Portanto, busca-se que ambas métricas possuam valores elevados (BOYD; ENG; PAGE, 2013). A Figura 2 apresenta um exemplo de cada curva com seus valores de AUC.

Sobre essas métricas de validação de modelos de classificação (curva ROC e PR), uma maneira adequada de poder avaliar as duas uniformemente é calcular algum tipo de média que retorne apenas um valor final para medir a performance. Uma opção é utilizar a Média Harmônica (MH), a que define-se para o grupo de números $\{x_1, x_2, \dots, x_i\}$ como

$$MH = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_i}}.$$

A MH outorga diferentes pesos aos valores dependendo da sua distribuição, tendo mais presença no resultado final os números mais pequenos (MOND; PECARIĆ, 1996). No geral, a curva ROC tem melhores resultados que a curva PR porque não considera o

Figura 2 – Exemplo de curva ROC (esquerda) e curva PR (direita) com seus respectivos valores de AUC.



Fonte: Adaptado de <https://www.r-bloggers.com/area-under-the-precision-recall-curve>.

overfitting do modelo, podendo gerar uma avaliação positiva de um modelo que não seja real. É por isso que é pertinente utilizar a MH quando tenta-se estimar um valor único entre essas duas métricas, entregando uma maior ponderação ao valor mais pequeno.

2.2.4.2 Não Supervisionado

O conjunto de dados sobre polivitimização é não supervisionado, pois não existe um atributo que determine algum valor esperado ou classe para as instâncias. Técnicas de agrupamento e de regras de associação (HASTIE; TIBSHIRANI; FRIEDMAN, 2009) são bem aplicadas sobre esse tipo de registros, mas considerando a natureza binária de todos os atributos (se a criança tem sofrido ou não cada uma das vitimizações), regras de associação são adequadas.

O aprendizado baseado em regras de associação tenta descobrir relações de interesse entre os elementos de um conjunto de dados através de métricas específicas para cada algoritmo (AGRAWAL; IMIELIŃSKI; SWAMI, 1993). O método foi desenvolvido para encontrar relações nas transações de compras de produtos. Por exemplo, a relação $\{farofa, feijão\} \Rightarrow \{arroz\}$ diz que quando leva-se farofa e feijão juntos, é provável ter arroz também (a presença do primeiro conjunto implica a presença do segundo).

Para gerar uma regra de associação precisa-se representar os dados através de conjuntos, sendo necessário dois principais: o primeiro corresponde ao grupo de itens $I = \{i_1, i_2, \dots, i_n\}$, com n como o número total de elementos (variáveis); e o segundo equivale ao de transações $T = \{t_1, t_2, \dots, t_m\}$, com m como o número total de instâncias da amostra. Com isso, é possível definir de forma geral uma regra de associação como $X \Rightarrow Y$, sendo $X, Y \subseteq I$. Além desses conjuntos, essa técnica possui dois elementos fundamentais para operar, o suporte e a confiança. Para que seja válida, é necessário que

essas duas métricas sejam maiores que o limiar que o usuário define (ZHANG; ZHANG, 2002).

- Suporte – Indica a frequência com que um grupo de elementos aparece no conjunto de dados completo. O suporte de um conjunto X pode ser calculado como:

$$\text{suporte}(X) = \frac{|\{X \subseteq t; t \in T\}|}{|T|},$$

em que t representa alguma transação de T .

- Confiança – Considerando a regra $X \Rightarrow Y$, pode-se interpretar ou entender a confiança como a probabilidade com que o conjunto Y aparece quando se tem o conjunto X , sendo representada por:

$$\text{confiança}(X \Rightarrow Y) = \frac{\text{suporte}(X \cup Y)}{\text{suporte}(X)}.$$

Dependendo das características dos dados e requisitos específicos do problema, existe uma determinada técnica para gerar as regras de associação (KOTSIANTIS; KANELLOPOULOS, 2006). No caso dessa pesquisa, utilizou-se o algoritmo *Apriori*, que opera sobre dados categóricos binários. Para obter as regras de associação de um conjunto de dados através do *Apriori*, o primeiro passo é encontrar as frequências com que aparecem os diferentes conjuntos de itens (ZHANG; ZHANG, 2002). O Algoritmo 2 apresenta o pseudo-código dessa técnica, o qual requer como parâmetros de entrada os conjuntos I e T previamente definidos, além do suporte mínimo de interesse *minsup*. A variável A corresponde ao grupo de conjuntos que serão avaliados em uma determinada iteração. Por exemplo, se existe o conjunto inicial I formado por $\{i_1, i_2, i_3\}$, então na primeira iteração $A = \{\{i_1\}, \{i_2\}, \{i_3\}\}$ (linha 2) e na segunda e terceira iterações $A = \{\{i_1, i_2\}, \{i_1, i_3\}, \{i_2, i_3\}\}$ e $A = \{\{i_1, i_2, i_3\}\}$, respectivamente (linha 33). O conjunto C encarrega-se de fazer a contagem das vezes que os elementos de A aparecem nas transações (t) de T , isso entre as linhas 5 e 16. A validação dos elementos que cumprem com o suporte mínimo é feita na linha 20 e são guardados na variável V , para depois uni-los com o conjunto final de saída F . Entre as linhas 27 e 35 apenas se faz a validação de que o conjunto A não tenha elementos que se sabem que no futuro não cumprirão com o suporte mínimo. Por exemplo, tem-se que $\{i_1\}$ não atinge-os, então $\{i_1, i_2\}$ também não irão satisfazê-los.

O segundo passo para obter as regras de associação é calcular a confiança de todas as sentenças do tipo $X \Rightarrow Y$ (onde X e Y podem ser qualquer elemento do conjunto F) que satisfazem o limiar estabelecido pelo usuário (AGRAWAL; IMIELIŃSKI; SWAMI, 1993), obtendo todas as regras de interesse baseadas no conjunto de dados inicial.

Diferente dos algoritmos supervisionados, os não supervisionados não têm métricas gerais para a avaliação dos resultados, principalmente porque não se tem um valor final

Algoritmo 2 Pseudo-código do algoritmo apriori de regras para associação (ZHANG; ZHANG, 2002).

Input: I, T, minsup ▷ I : itens; T : instancias; minsup : suporte mínimo.

1: $F \leftarrow \{\}$; ▷ F : conjuntos que cumprem com o suporte mínimo.

2: $A \leftarrow \{\{i_1\}, \{i_2\}, \dots, \{i_n\}\}$; ▷ A : conjuntos que serão avaliados.

3: **while** $A \neq \{\}$ **do**

4: $C \leftarrow \{\}$; ▷ C : conjunto de elementos candidatos com suas frequências.

5: **for all** $t \in T$ **do**

6: **for all** $a \in A$ **do**

7: **if** $a \subseteq t$ **then**

8: **if** $a \in C$ **then**

9: $c_a.\text{contagem} \leftarrow c_a.\text{contagem} + 1$;

10: **else**

11: $C \leftarrow C \cup \{a\}$;

12: $c_a.\text{contagem} \leftarrow 1$;

13: **end if**

14: **end if**

15: **end for**

16: **end for**

17: $A \leftarrow \{\}$;

18: $V \leftarrow \{\}$; ▷ V : Conjuntos dos quais será validado seu suporte.

19: $\text{total} \leftarrow 0$; ▷ total : quantidade de conjuntos com suporte mínimo.

20: **for all** $a \in C$ **do**

21: **if** $c_a.\text{contagem}/T.\text{tamanho} \geq \text{minsup}$ **then**

22: $V \leftarrow V \cup c_a$;

23: $\text{total} \leftarrow \text{total} + 1$;

24: **end if**

25: **end for**

26: **if** $\text{total} \neq 0$ **then**

27: $F \leftarrow F \cup V$;

28: $i \leftarrow 1$;

29: $j \leftarrow 2$;

30: **while** $i < \text{total}$ **do**

31: $i \leftarrow i + 1$;

32: **while** $j \leq \text{total}$ **do**

33: $j \leftarrow j + 1$;

34: $A \leftarrow A \cup \{V_i \cup V_j\}$;

35: **end while**

36: **end while**

37: **end if**

38: **end while**

que obter ou prever (RUSSELL; NORVIG, 2016). Contudo, existem diversos trabalhos que propõem técnicas de validação para algoritmos não supervisionados, como o trabalho feito por Halkidi et al. (2001) para modelos de agrupamento, ou por Ordoñez (2006) para regras de associação.

Especificamente nas regras de associação, o mesmo procedimento feito para segui-las é uma forma de validar sua qualidade, isso através do suporte e da confiança (ZHANG; ZHANG, 2002). Na medida que essas duas métricas tenham valores próximos a 1, as regras serão mais frequentes e precisas. Cabe mencionar que um $suporte(X)$ elevado implica que uma grande parte das instâncias totais estejam incluídas no conjunto de itens frequentes X , enquanto uma $confiança(X \Rightarrow Y)$ elevada resulta em uma alta probabilidade de ao se ter X também se terá Y . Existem também outras métricas que permitem avaliar a qualidade das regras de associação, destacando neste trabalho os valores *lift*, *conviction* (BRIN et al., 1997), e *leverage* (PIATETSKY-SHAPIRO, 1991).

- *Lift* – Taxa de vezes que estão presentes o antecedente e consequentes juntos ($suporte(X \cup Y)$) assumindo que são elementos independentes ($suporte(X) \times suporte(Y)$). Sua expressão matemática corresponde a:

$$lift(X \Rightarrow Y) = \frac{suporte(X \cup Y)}{suporte(X) \times suporte(Y)}.$$

Lift é um valor que vai no intervalo $[0, \infty[$; quando é maior a 1 existe uma dependência entre X e Y , quando é igual a 1 são elementos independentes, e quando é menor a 1 o antecedente e o consequente são mutuamente exclusivos (BRIN et al., 1997).

- *Conviction* – Assumindo que X e a ausência de Y são elementos independentes ($suporte(X) \times suporte(\neg Y)$), corresponde ao quociente obtido entre esse valor e a taxa de vezes que o antecedente não aparece junto ao consequente ($suporte(X \cup \neg Y)$). Sua expressão matemática corresponde a:

$$conviction(X \Rightarrow Y) = \frac{suporte(X) \times suporte(\neg Y)}{suporte(X \cup \neg Y)} = \frac{1 - suporte(Y)}{1 - confian\c{a}(X \Rightarrow Y)}.$$

Conviction é um valor que vai no intervalo $[0, \infty[$; quando é maior a 1 existe uma dependência entre X e Y , quando é igual a 1 são elementos independentes, e quando é menor a 1 o antecedente e o consequente são mutuamente exclusivos (BRIN et al., 1997).

- *Leverage* – Diferença entre a taxa de vezes que estão presentes o antecedente e consequentes juntos ($suporte(X \cup Y)$) e o valor resultante ao considerar que são elementos independentes ($suporte(X) \times suporte(Y)$). Sua expressão matemática corresponde a:

$$leverage(X \Rightarrow Y) = suporte(X \cup Y) - suporte(X) \times suporte(Y).$$

Leverage é um valor que vai no intervalo $[-1, 1]$; quando é maior a 0 existe uma dependência entre X e Y , quando é igual a 0 são elementos independentes, e quando é menor a 0 o antecedente e o consequente são mutuamente exclusivos (PIATETSKY-SHAPIRO, 1991).

3 Revisão Bibliográfica

Desde o final do século passado, ML tem sido usado para ajudar a entender melhor o problema da violência infantil. Neste capítulo, diferentes trabalhos sobre esse assunto serão expostos com o intuito de entender como lidam com esse fato, destacando alguns elementos como: descrição da amostra, quais ferramentas foram utilizadas, a avaliação e resultados dos modelos gerados, os fatores de risco selecionados, entre outras qualidades importantes ao contexto deste estudo. Além disso, no final do capítulo situa-se a presente pesquisa entre os estudos analisados.

3.1 Estudos Relacionados

Para encontrar trabalhos semelhantes ao proposto, utilizaram-se os motores de busca bibliográfica Google Scholar¹ e Scopus², além da base bibliográfica SciELO³ (*Scientific Electronic Library Online*). A pesquisa foi feita nas línguas: espanhol, inglês, e português por igual; isso através de cadeias de *strings* formadas pelas palavras-chave: *child abuse*, *child maltreatment*, *predicting child abuse*, e *predicting child maltreatment* para referenciar a violência infantil; e para aludir o ML utilizaram-se: *prediction*, *machine learning*, *artificial intelligence*, *neural network*, *decision tree*, *regression*, e *apriori algorithm*. Assim, as cadeias criadas usando os operadores lógicos OR e AND foram:

1. (*child abuse* OR *child maltreatment*) AND (*prediction* OR *machine learning* OR *artificial intelligence* OR *neural network* OR *decision tree* OR *regression* OR *apriori algorithm*);
2. (*prediction*) AND (*child abuse* OR *child maltreatment*) AND (*machine learning* OR *artificial intelligence* OR *neural network* OR *decision tree* OR *regression* OR *apriori algorithm*);
3. (*predicting child abuse* OR *predicting child maltreatment*) AND (*machine learning* OR *artificial intelligence* OR *neural network* OR *decision tree* OR *regression* OR *apriori algorithm*); e
4. (*child abuse* OR *child maltreatment*) AND (*neural network* OR *decision tree* OR *apriori algorithm*).

¹ <https://scholar.google.com>

² <https://www.scopus.com/home.uri>

³ <https://www.scielo.org>

Como resultado dessa pesquisa, obteve-se 19 trabalhos relacionados, todos abordando o problema da violência infantil a partir de um enfoque estatístico ou computacional.

A análise desses estudos foi dividida em quatro grupos, os quais estão relacionados com as técnicas de ML utilizadas para analisar o problema da violência infantil. O primeiro corresponde aos que usaram algum tipo de regressão; o segundo aos estudos com algoritmos de classificação; o terceiro aos que misturaram diferentes tipos de técnicas (principalmente estudos comparativos); e o quarto aos que utilizaram outras técnicas não convencionais. Cada uma das seguintes subseções correspondem a uma dessas classificações, as que destacam as principais características dos artigos estudados.

3.1.1 Técnicas de Regressão

O primeiro trabalho que tenta prever a violência infantil por meio de uma análise matemática foi desenvolvido por Altermeier et al. (1984), que abordou essa problemática investigando mulheres grávidas (1.400 no total). Foi aplicada uma enquete com os seguintes aspectos (atributos de entrada): saúde da criança, apoio de pessoas próximas, atitude frente à gravidez e saúde própria (onde destaca consumo de drogas ou álcool). Para selecionar os grupos de variáveis preditivas, foi usado um modelo de regressão hierárquico (GELMAN; HILL, 2006), enquanto, para estabelecer se existiu violência ou não, foram usados os reportes de abuso das mesmas mulheres para seus outros filhos, estabelecendo que 273 casos foram de alto risco. O período de avaliação deles foi de 2 anos após a entrevista, onde somente 22% das famílias classificadas com alto risco tiveram reportes de violência infantil, tendo muitos FP. Esse problema implicou que os pesquisadores fizeram uma segunda entrevista com as mulheres composta por 20 perguntas (preditores mais importantes) escolhidas por meio dos resultados da relação com o abuso, onde quase todas tiveram um valor do coeficiente de correlação (R) maior que 0,34. Desta forma, gerou-se uma lista de questões que permitiram prever a violência infantil, porém a maioria das perguntas tinham uma relação intrínseca com o problema de pesquisa. Um exemplo disso é a questão “tendência agressiva durante a entrevista”.

Algumas considerações que podem ser estabelecidas desse trabalho são: a análise matemática foi somente estatística, não permitindo gerar modelos que realmente possam prever a violência; o fato de determinar se houve ou não abuso através dos reportes, em muitos casos pode não ser real pois nem sempre eles são denunciados (explicação dos pesquisadores para os FP).

Outro caminho para avaliar a presença de violência infantil é realizar entrevistas padronizadas, como fez Burrell et al. (1994). O estudo teve dois objetivos principais: o primeiro deles visa determinar se existe relação entre o potencial abuso infantil com o estresse parental, apoio social e recursos familiares; e o segundo busca entender se modelos preditivos de violência de famílias com filhos com e sem necessidades especiais têm

diferenças. A pesquisa foi feita com 113 mães, 53 com filhos que apresentaram alguma necessidade especial e 60 que não. A primeira parte do estudo baseou-se em aplicar 4 diferentes entrevistas às mães: inventário do potencial abuso infantil (*Child Abuse Potential Inventory*, CAPI), estresse parental, percepção dos recursos, e apoio social. Essas entrevistas buscaram algum padrão de correlação entre a CAPI e as demais enquetes através de uma análise de regressão. Dessa forma, a maior correlação foi com o estresse parental com um R de aproximadamente 0,53. Na segunda etapa realizou-se o mesmo estudo, porém dividindo a amostra em dois grupos, um com as mães de crianças com necessidades especiais e outro sem. O resultado não foi significativo, destacando somente a desigualdade de correlação dos grupos no escopo de estresse parental com uma diferença de 0,121, sendo o valor mais alto o da amostra de mães com crianças com necessidades especiais. Para avaliar os resultados, simplesmente foi feita uma análise qualitativa deles, validando as duas hipóteses: estresse parental tem relação direta com a violência e crianças com necessidades especiais são mais propensas a receber o abuso.

No geral, a pesquisa tem uma boa abordagem do problema ao propor utilizar ferramentas que já existiam (entrevistas sociais) para relacionar a violência com outras características familiares, mas ainda assim não utiliza fatores de risco para prever a violência infantil. É por isso que estabelece-se que o estudo é orientado ao aspecto teórico do problema, e não tenta uma resolução via métodos matemáticos.

Rodriguez e Green (1997) fizeram outra pesquisa com a mesma metodologia que o estudo anterior. Assim, por meio de enquetes padronizadas relacionaram o estresse parental e a expressão de raiva com a violência infantil potencial, utilizando regressão hierárquica. Foram empregadas três ferramentas para esse estudo, a CAPI como variável dependente (valor a ser predito), o *Parenting Stress Index* (PSI) e o *State-Trait Angry Expression Inventory* (STAXI) como variáveis independentes. Dessa forma, tentou-se prever o valor de CAPI através de PSI e STAXI, dividindo a pesquisa em dois estudos:

1. Voluntariamente 39 pais fizeram parte da investigação, na qual tiveram que responder às três enquetes. Logo, foram obtidas as correlações de PSI e STAXI com CAPI, sendo os valores de R iguais a 0,67 e 0,69, respectivamente. Depois disso, ao modelo gerado entre CAPI e PSI foi adicionada a variável STAXI criando outro com um R de 0,83, estimando que o modelo de regressão hierárquica consegue ser um bom preditor da violência infantil potencial. Além disso, tentou-se encontrar alguma relação entre elementos demográficos (sexo, idade, etnia, ingressos, etc.) com os resultados das enquetes. Entretanto, o único resultado relevante foi a correlação negativa entre a idade dos pais com a expressão de raiva deles ($R=-0,41$);
2. Voluntariamente 84 pais responderam as enquetes, que receberam através das escolas de seus filhos. Foi feita a mesma análise, obtendo valores de R entre 0,53 e 0,44 para

CAPI e PSI e STAXI. Ao gerar o modelo hierárquico, o R final foi de 0,62. Também foi feita uma abordagem demográfica, obtendo a relação negativa entre o número de filhos e a expressão de raiva ($R=-0,33$).

Baseado nos resultados, não foi feita nenhuma validação além de considerar as correlações entre as enquetes. Ademais, todos os resultados foram fundamentados nas respostas do CAPI, pelo qual a pesquisa somente permanece na teoria e não consegue demonstrar que realmente houve violência a partir dos dados usados. Ainda vale destacar que, apesar de os estudos serem distintos, fizeram uso de um mesmo procedimento, alterando apenas a amostra. No primeiro, a maioria dos sujeitos tinham perfis econômicos bons, o que estimou que não era um grupo representativo da população toda; e no segundo, estabeleceu-se que os pais responderam as enquetes de forma defensiva, considerando que os questionários chegaram através de seus filhos (fato que os pesquisadores usaram para explicar a diminuição dos coeficientes de correlação entre os estudos). Por consequência, pode-se dizer que os dois estudos tiveram um viés ao selecionar as amostras, não sendo representativas da população.

DePanfilis e Zuravin (1999) também pesquisaram sobre o problema da violência infantil com uma abordagem estatística, mas tentando identificar a recorrência do abuso, estudo que utilizou uma população de 446 mães. Sobre os fatores de risco utilizados, foram considerados os tópicos: colocação da criança (se foi afastada da sua família); informação da violência (tipo, severidade e prioridade que tem a criança no centro de proteção); vulnerabilidade da criança (problemas mentais, de desenvolvimento infantil, e se houve a presença de alguma criança menor de 6 anos na moradia); problemas pessoais do cuidador (consumo de drogas ou álcool e *deficit* na resolução de problemas); Violência entre Parceiros Íntimos (VPI) sobre a mãe; estresse familiar (quantidade de filhos, mãe muito nova ao nascer o primeiro filho, e duração dos anos de criança); condições do lar (falta de recursos econômicos, falta de moradia, moradia em estado ruim e falta de recursos em atenção médica); *deficit* de ajuda social (famílias ou amigos, e sistemas de ajuda). Com todas essas categorias, tinham-se 52 variáveis para analisar. Sobre a análise, a mesma foi dividida em duas partes: (1) examinar as relações entre as variáveis usando uma função de sobrevivência com o modelo de Cox (COX, 2018) e estimar a probabilidade de que um sujeito consiga sobreviver em um tempo determinado (que não aconteça recorrência de violência nele); e (2) selecionar variáveis que tenham relações relevantes entre elas durante o tempo que dura o estudo de caso. O modelo final gerado foi o de riscos proporcionais de Cox (COX, 1972) que, da mesma maneira que os de regressão múltipla (AIKEN; WEST; RENO, 1991), permite analisar as relações entre variáveis dependentes com a independente, mas também facilita relacioná-las com o tempo que tem decorrido os fatos. Por exemplo, as variáveis que determinam a recorrência da violência em um período de 1 ano talvez não são as mesmas que em um período de 1 mês. Logo, dessa análise apenas foram

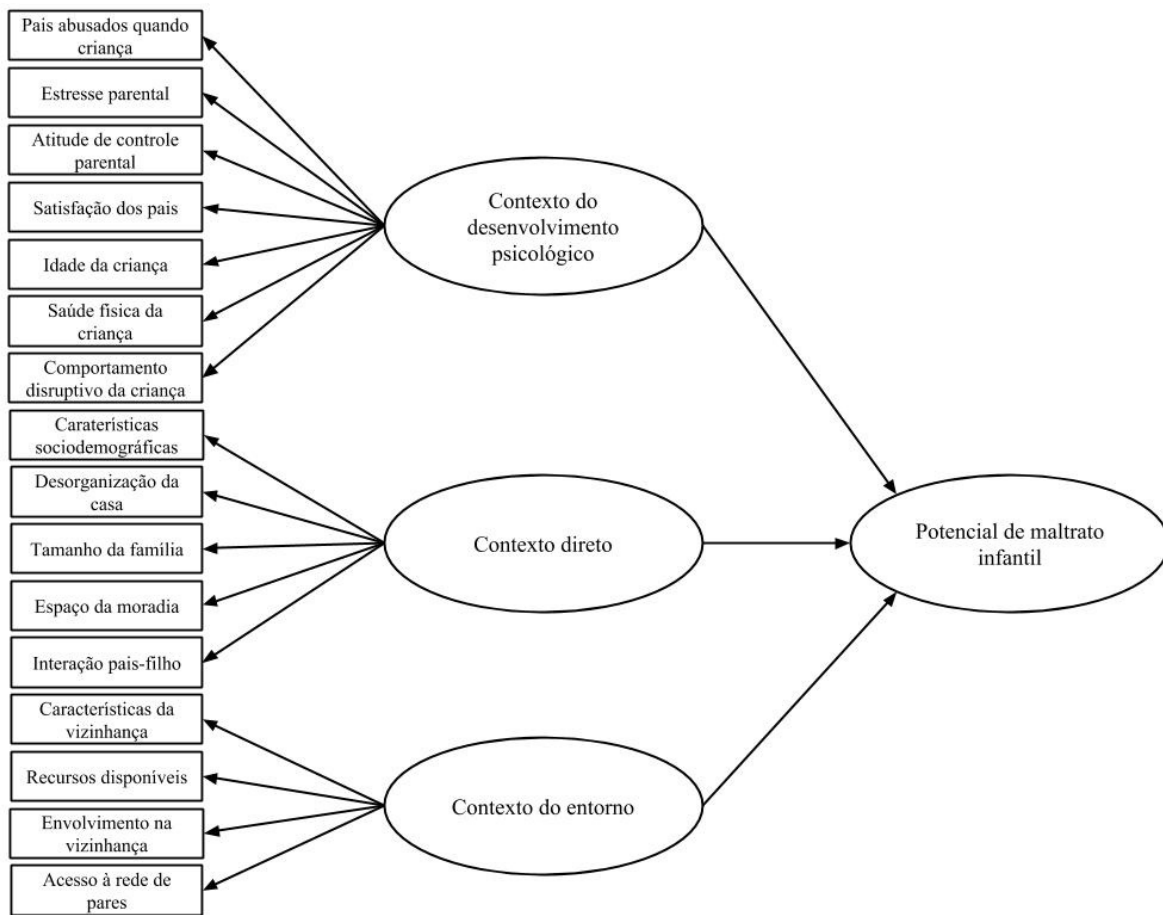
considerados 10 fatores críticos: colocação, prioridade, tipo e severidade do abuso, fatores da criança, fatores da mãe, VPI, estresse familiar, condições do lar e suporte social.

Posteriormente, gerou-se um modelo sem condicionamento entre as variáveis, que teve como principais fatores a colocação, a VPI e o suporte social, sendo esses os fatores que apresentam um maior risco para sobreviver sem a recorrência de violência. Com esse resultado, foi criado outro modelo baseado na relação hipotética que têm as variáveis VPI e suporte social, com o qual obteve-se 92% de probabilidade de que não aconteça recorrência se esses fatores não estão presentes na família. Alguns pontos interessantes do artigo: modelos de Cox são aplicáveis ao problema da predição da violência infantil, pensando em quando uma criança sobrevive ao abuso (neste caso, há a recorrência dele); a VPI não tinha sido considerado até esse estudo e há indícios de ser um fator crítico importante da violência, além do suporte social que a família tem. Sobre o que pode ser melhorado, os modelos foram criados com apenas uma medida de recorrência da violência (a confirmação futura de um novo reporte de abuso), que na prática não permitem realmente prever quando novamente irá acontecer o abuso (pode ser que um novo reporte nunca chegue).

Um trabalho diferente aos anteriores foi apresentado por Begle et al. (2010), no qual tentou-se validar de forma empírica os modelos teóricos de desenvolvimento ecológico e de risco acumulado, que descrevem a violência infantil. O primeiro, define que existem 3 elementos fundamentais que caracterizam o abuso (ver Figura 3): o desenvolvimento psicológico (*developmental-psychological*), contexto direto (*immediate*) e os recursos sociais (*broad*). No segundo, se estabelece que o risco da violência aumenta (acumula-se) quando o número de fatores de risco for maior. Sobre o conjunto de dados, foram entrevistadas 610 figuras parentais (pais, mães ou cuidadores) de crianças de 3 a 6 anos de idade. Sobre os fatores de risco, foram usadas 22 variáveis que descrevem os elementos do modelo ecológico, além de uma variável preditora da violência infantil potencial obtida através do CAPI. Os fatores de risco estiveram relacionados com: características da criança, dos cuidadores, de seu relacionamento, além das características do entorno e da família.

Para validar o modelo ecológico e o de risco acumulado, foram utilizados os métodos do teste de qui-quadrado (MCHUGH, 2013) e regressão linear múltipla (MONTGOMERY; PECK; VINING, 2012), respectivamente. Para o primeiro, não houve ajuste com os dados coletados, implicando que o modelo ecológico não pode se relacionar com os valores empíricos. Sobre o segundo, após realizar a regressão linear, notou-se que há uma relação já que apresentou 76% de precisão na classificação dos casos. Um ponto importante desse trabalho é a mistura que foi feita entre modelos teóricos e empíricos, ainda mais quando eles vêm de áreas de estudo muito diferentes (neste caso a psicologia e a estatística), conseguindo validações deles de outras maneiras. Porém, nesse estudo, os dados coletados não foram representativos do problema, visto que a média do resultado do teste CAPI foi muito baixa, o que implica que as famílias no geral não tinham problemas

Figura 3 – Modelo de desenvolvimento ecológico da violência infantil (BEGLE; DUMAS; HANSON, 2010).



Fonte: Adaptado de Begle, Dumas e Hanson (2010).

de abuso infantil, não permitindo gerar uma predição certa da violência infantil potencial.

Dubowitz et al. (2011) tentaram analisar o problema da violência infantil através de outra perspectiva. Realizaram um estudo longitudinal que avaliou a violência em um período de 10 anos com 224 mães de famílias de baixa renda para identificar quais fatores críticos podem produzir violências, para posteriormente predizê-los. Assim, foi feita uma análise de sobrevivência através da regressão de Cox, onde o fato de “sobreviver” para as famílias foi não ter reportes de violências nos centros de proteção de crianças. Os fatores críticos foram 12 no início, os quais estão relacionados com as crianças, suas mães, suas famílias e a comunidade onde moram.

A etapa de gerar os modelos teve duas partes, na primeira foi feita uma análise bivariada entre todas as variáveis independentes com a dependente (reporte de violência nos centros de proteção), obtendo-se um modelo inicial em que foram selecionados os fatores com maior relação com esta última: desenvolvimento cognitivo da criança, nível educacional da mãe, se a mãe apresenta depressão, número de filhos e consumo de drogas da mãe; fazendo o mesmo processo de análise com esses fatores como segunda fase. Esse

estabeleceu que o uso de drogas é a variável com a maior relação ao reporte de abuso; enquanto o nível educacional o de menor relação. Esse estudo permitiu detectar relações entre fatores críticos e a violência infantil que antes não havia sido esclarecido completamente (como o nível de escolaridade dos pais), mas para obter algumas informações e fazer toda essa análise foram necessárias entrevistas padronizadas, bem como para conhecer o estado de depressão da mãe ou o suporte social da família, tornando difícil replicar a pesquisa.

Considerando que o problema da violência infantil é global, no Japão também tem-se desenvolvido pesquisas sobre esse tema. Horikawa et al. (2016) realizaram um estudo que teve como objetivo criar um modelo multivariado que identifique as crianças com um alto risco de sofrer a primeira recorrência de violência utilizando dados governamentais. O conjunto de dados foi coletado entre os anos 1996 e 2011 com 716 registros, sendo necessário usar uma técnica de *bootstrapping* (KOHAVI et al., 1995) para acrescentar amostras. Além disso, no começo tinha-se 26 variáveis, que podem ser agrupadas em: tipo de violência e frequência, características da criança e do perpetrador, características do entorno e da família, entre outros elementos. Ao fazer um processo de *feature extraction*, restaram apenas 6 variáveis: idade da criança, idade do agressor, histórico de violência na infância do agressor, estabilidade financeira na moradia, presença de alguém que vigiasse a vítima e se a informação do caso vem de alguma organização oficial. A variável dependente foi a presença do reporte de violência após um ano do primeiro incidente registrado (binária). Utilizou-se regressão logística multivariada para a modelagem, resultado que foi avaliado com a curva ROC, obtendo uma AUC de 0,66. Posteriormente, categorizaram-se os casos através do modelo nos níveis de risco baixo, médio e alto de sofrer recorrência de violência, que permitiu determinar os passos a seguir nessa situação de abuso. Sobre as limitações da pesquisa, pode-se dizer que a quantidade de fatores críticos usados foi limitada, não variando-os entre análises nem aumentando o número deles. Ademais, considerando o intervalo de anos dos dados, não foi levado em conta que os fatores de risco podem mudar no tempo, o que tem uma implicação direta na validade do modelo final.

3.1.2 Técnicas de Classificação

O primeiro trabalho estudado sobre a predição da violência infantil através da classificação de instâncias foi feito por Frank et al. (1992), onde tentaram determinar diferenças em crianças hospitalizadas (entre as que sofrem ou não de violências). Além disso, propuseram-se estabelecer a frequência da violência mediante dois caminhos: com um reporte clínico que fazem os médicos e com um diagnóstico técnico feito com árvores de decisão. A amostra no total foi de 749 crianças, das quais 114 apresentaram sintomas de violência que foram classificados em três grupos: abuso físico, abuso sexual e negligência física; sendo no total 12 diferentes sintomas (variáveis) usados para a análise. Sobre os

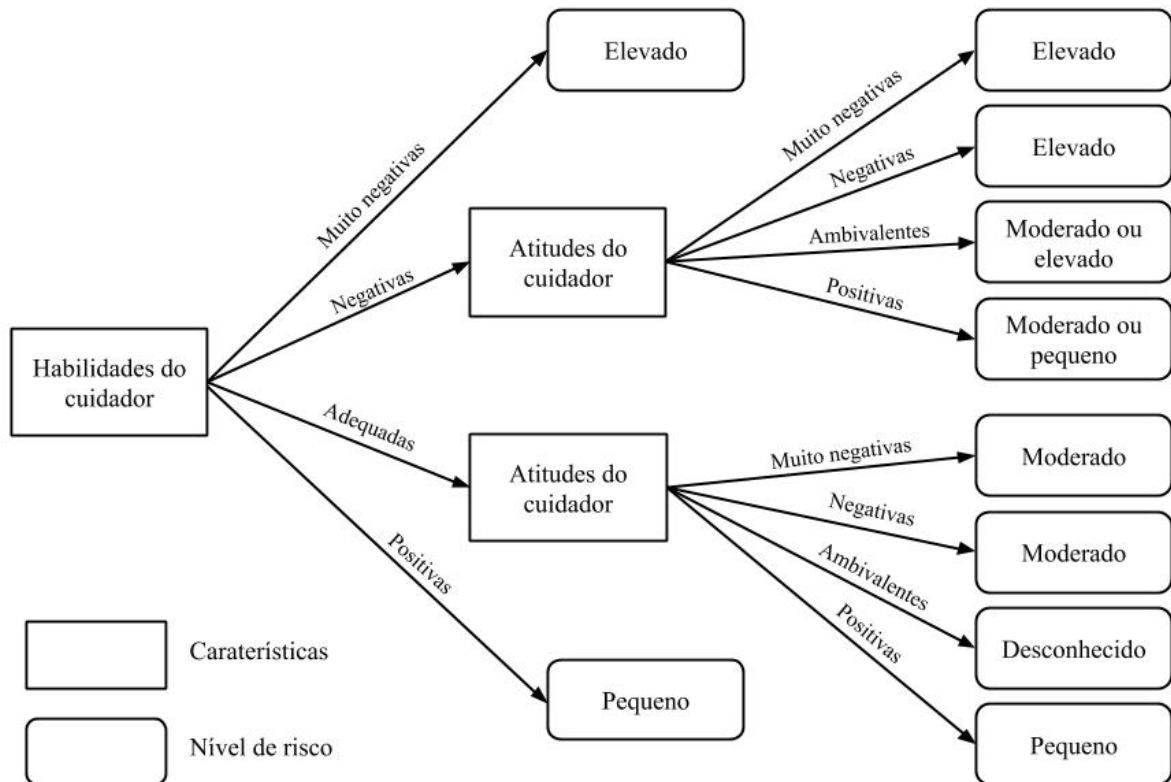
resultados, o reporte clínico determinou que 10 crianças foram violentadas e o diagnóstico técnico 7, sendo somente 4 as que estiveram presentes nos dois grupos. No final, os médicos definiram quais crianças estavam sendo abusadas e o procedimento de como trabalhar com elas. Portanto, a análise feita com árvores de decisão não foi muito considerada pelos especialistas, já que os resultados não eram parecidos aos deles. Essa pesquisa não teve um resultado positivo no enfoque da modelagem matemática do problema. Acredita-se que, semelhante ao estudo anterior (ALTEMEIER et al., 1984), não é um bom caminho estabelecer que uma criança está sendo abusada com uma entrevista ou um diagnóstico médico, já que os casos podem ser subjetivos ou podem não apresentar sintomas visíveis.

No trabalho realizado por Little e Nixon (1998) explorou-se a possibilidade da IA contribuir aos serviços de risco em proteção infantil, sendo o problema definido “estimar o nível de risco (baixo, médio ou alto) que tem as crianças que participam nos centros de proteção”, e a hipótese “o uso de árvores de decisão melhora a estimativa do risco de abuso que é feita por assistentes sociais”. Esta pesquisa foi realizada em conjunto com um grupo de assistentes sociais que, antes de qualquer análise, tiveram que realizar três tarefas. A primeira foi identificar quais são os fatores mais relevantes para determinar o nível de risco, a segunda diz respeito a desenvolver uma escala que permita classificar cada um dos fatores previamente encontrados e a terceira foi classificar cada um dos casos utilizados para o estudo com a escala criada. Foram utilizados 20 estudos de caso com 8 diferentes fatores a avaliar (tipo de violência infantil, atitude dos cuidadores, capacidade dos cuidadores, impacto no comportamento da criança, registros prévios de violência, severidade da violência, vulnerabilidade da criança e grau de proximidade do agressor na criança), além do classificador preditivo “nível de risco”. A Figura 4 apresenta o modelo resultante que tem apenas 2 dos 8 fatores. Essa árvore fornece regras que permitem prever ou refletir sobre o nível de risco que uma criança pode ter. Por exemplo, se “habilidade do cuidador” e “atitudes do cuidador” forem negativas, o risco da criança ser agredida é alto.

O primeiro elemento a notar sobre esse trabalho é que com um pequeno conjunto de dados foi possível entregar um modelo que atenda ao objetivo de classificar o nível de risco, além do fato de uma árvore de decisão garantir uma simplicidade em termos de visualização e compreensão para qualquer pessoa que não seja especialista em conceitos de IA. Porém, também pode ser considerado como um modelo ruim, principalmente porque apresenta uma quantidade de dados pequena. Outro ponto a destacar é a dificuldade que os assistentes sociais tiveram para definir a agressão com base nos níveis de risco dos fatores, precisando transformar um conceito qualitativo em quantitativo.

O primeiro uso de redes neurais no problema da violência infantil foi feito por Schwartz et al. (2004), tentando prever quais crianças apresentam um “dano padrão” (*harm standard*) com dados do Estudo Nacional de Incidência de Abuso e Negligência

Figura 4 – Árvore de decisão com 2 fatores de risco (LITTLE; RIXON, 1998).



Fonte: Adaptado de Little e Rixon (1998).

Infantil (*National Incidence Study of Child Abuse and Neglect*, NIS) dos Estados Unidos, usando 1.767 registros. No início foram 30 variáveis as escolhidas. A partir do pré-processamento de dados, principalmente da transformação de atributos categóricos em binários, foram obtidos 141 variáveis no total. Os autores não especificam quais são os fatores críticos usados, mas falam da variável preditora, um elemento binário que é representado por 0 quando não atinge os critérios de dano e 1 quando atinge. Somente foi implementada uma rede neural do tipo *MultiLayer Perceptron* (MLP), que teve a seguinte estrutura: 141 entradas, uma camada escondida com 71 neurônios e a camada de saída com 1 neurônio. Essa última camada entrega um valor final entre 0 e 1, ou seja, a classificação propriamente dita: no intervalo $[0,0;0,1]$ não há violência, entre $]0,1;0,9]$ não é possível determiná-lo e entre $]0,9;1,0]$ há violência. O conjunto de treinamento fez uso de 65% dos dados totais. A precisão da rede sobre os dados de teste foi de 89,6% de acertos com 64 casos classificados erroneamente, em que 48 foram considerados indeterminados. No geral, a pesquisa foi muito simples (apenas um algoritmo com uma execução) e no artigo não pode-se obter muitas informações sobre os dados usados, em especial dos fatores críticos. Os autores detalham que o ponto forte do trabalho foi ter empregado o NIS, porém o mesmo não é replicável em todos os contextos ou para diferentes países.

No trabalho de Amrit et al. (2017), aborda-se o problema da predição da violência

infantil através da mineração de texto, o que implicou o uso de dados não estruturados para realizar parte da análise, sendo este o seu problema de pesquisa. Assim, foi considerada a hipótese de que “com um volume maior de dados, que não estão em um formato estruturado, é possível obter melhores resultados na previsão do abuso infantil na aplicação de técnicas de ML”. Eles usaram um conjunto de dados do departamento de saúde da criança da Holanda, além de trabalhar com um grupo de profissionais dessa instituição. Esta base de dados tem 195.188 consultas, com uma média 14,82 por criança, sendo 13.170 crianças classificadas como agredidas. Três algoritmos diferentes foram aplicados e comparados: *Naive Bayes* (NB), *Random Forest* (RF) e SVM. Todos os algoritmos são destinados à classificação de dados. Para os dois primeiros algoritmos foram usados também as técnicas ANOVA e qui-quadrado para selecionar um conjunto de variáveis a serem analisadas. Na validação dos resultados foi utilizada a técnica da AUC. Sobre os dados usados, os estruturados correspondem a: características da criança, relacionamento familiar e presunções de violência. Em termos de dados não estruturados, foram consideradas anotações médicas feitas por enfermeiros ou pediatras que registravam os dados das crianças e da dinâmica familiar.

Depois de ajustar os três modelos, o algoritmo que obteve o melhor resultado foi o SVM, com um valor de AUC de 0,906. Esse resultado foi ainda melhor utilizando dados estruturados e não estruturados em conjunto, sendo a AUC de 0,914. Um destaque dessa pesquisa foi a consideração de fatores como a quantidade de consultas que uma criança pode ter e a extensão em palavras nos seus relatos, variáveis que poderiam prever melhor se houve ou não uma violência.

Schwartz et al. (2017) tiveram um objetivo diferente, pensaram em melhorar a validação dos centros de proteção de crianças tanto na predição do risco como em determinar qual serviço é o indicado para cada caso (ações judiciais, remover a criança da moradia, incluir em algum programa, ou suporte familiar). O conjunto de dados foi formado por 78.394 casos com aproximadamente 150 variáveis relacionadas à violência. Os algoritmos usados para gerar os modelos não são explicitados, somente afirma-se que foram utilizadas técnicas de classificação (é possível ver no documento um par de árvores de decisão que foram gerados e são usadas como exemplos); e a forma de validação foi através da curva ROC. A análise foi dividida em 3 partes e em cada uma considerou-se uma pergunta de pesquisa: (1) “qual é a probabilidade de que um reporte de violência infantil seja comprovado durante o processo de investigação do caso?”, obtendo do modelo uma AUC de 0,87; (2) “qual é a probabilidade de que um caso não seja encaminhado ao serviço apropriado?”, sendo a AUC de 0,81; e (3) “qual serviço tem a maior probabilidade de impedir que um caso novamente apresente violência infantil?”, para essa pergunta não utilizou-se nenhuma técnica e estabeleceu-se que se os dois modelos anteriores forem implementados em conjunto, as respostas em casos de violência infantil poderiam melhorar em até 30%.

No geral, o estudo apresenta bons resultados obtidos com os modelos de ML, além de por o foco em determinar qual serviço é o indicado para cada caso (*prescriptive analysis*), mas não detalha quais algoritmos foram utilizados ou como eles foram configurados, ou quais variáveis consideraram-se para cada um.

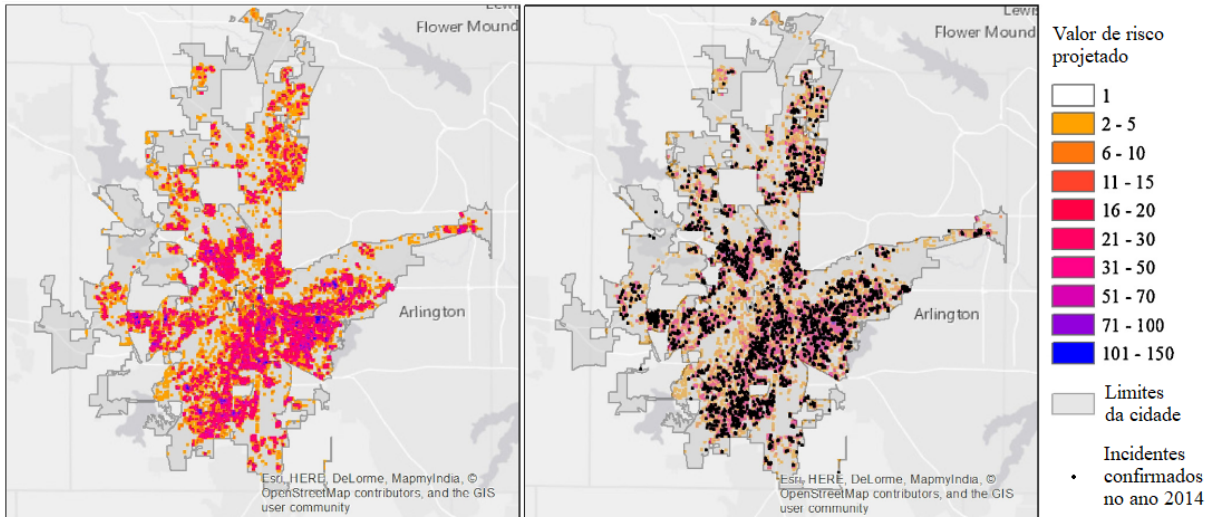
3.1.3 Demais Técnicas

Vaithianathan et al. (2013) fizeram uma pesquisa em que não detalham informação sobre qual algoritmo usaram, somente mostram que foi implementado um modelo preditivo de risco (*Predictive Risk Model*, PRM). Este trabalho teve como objetivo validar o desempenho desse mesmo modelo. Os dados usados foram obtidos dos centros de proteção de crianças e, no total, tinha 103.397 registros com 57.986 crianças. Além disso, o conjunto de dados era composto por 224 variáveis e usaram-se 132, as quais foram selecionadas através de um procedimento de *backward selection*. Tais variáveis foram agrupadas em 3 categorias: da criança, do cuidador principal e do seu casal. Sobre os algoritmos para gerar os modelos, não é dito nada além de ter usado um PRM e que a técnica de validação foi a curva ROC, usando 70% dos dados para treinamento e 30% para teste. O resultado obtido foi de 0,76 de AUC, não estabelecendo se é um resultado bom ou ruim. Aliás, são determinadas as probabilidades de que uma criança sofra violência para os grupos com menor e maior risco, sendo 25 vezes mais para esse último grupo. Outro resultado considerável é que o fator raça (ou grupo étnico) não contribui para o modelo final. Nesse estudo é mostrada uma limitação muito importante, visto que os dados relacionados ao contexto social não foram empregados nas análises.

Outro trabalho não convencional tentado aplicar uma técnica diferente com dados que possam prever a violência infantil foi o feito por Daley et al. (2016). Eles aplicaram a Modelagem de Terreno de Risco (*Risk Terrain Modeling*, RTM), método que permite determinar quais são os lugares com maior risco de acontecer algum fato que deseja-se estudar em uma área estabelecida (CAPLAN; KENNEDY, 2011). Nesse caso, é o primeiro estudo que tenta relacionar a probabilidade de que aconteça violência infantil em um lugar específico, aplicando-o em uma cidade dos Estados Unidos. Os dados foram coletados de diferentes centros públicos oficiais, com um total de 10 atributos. Por meio de um processo de regressão foram escolhidos os 6 com maiores coeficientes (mais significativos): pobreza, violência doméstica, agressões agravadas, fugitivos, assassinatos e delitos de drogas. Logo após gerar o modelo, esse foi testado usando 5.391 novos casos de violência infantil nessa mesma cidade, estabelecendo se estavam em uma zona de risco ou não, o que resultou em uma precisão do 52% nas áreas de alto risco e 90% nas de menor risco, conforme mostrado na Figura 5.

Pela natureza desse estudo, muitos fatores de risco comumente relacionados à violência infantil ficaram fora da pesquisa, principalmente pelo fato de relacioná-los com

Figura 5 – Mapa da estimação do riscos de violência infantil: imagem da esquerda corresponde à estimação para o ano de 2014 (baseada nos dados do ano 2013) e a imagem da direita representa à estimação do ano de 2014 com os casos reais (pontos são levemente distorcidos para manter a privacidade) (DALEY et al., 2016).



Fonte: Adaptado de Daley et al. (2016).

um ponto físico determinado da cidade. Além disso, para conseguir aplicar esse estudo em outros lugares é necessário obter informação da localização de todas as variáveis que possam prever o abuso, fazendo com que não seja replicável em todos os contextos.

3.1.4 Análises Comparativas entre Distintos Modelos

Marshall e English (2000) propuseram dois objetivos: (1) explorar a utilidade das redes neurais na análise de dados de ciências sociais com foco na avaliação de risco na proteção infantil; e (2) comparar esses resultados com os obtidos por modelos de regressão. Para realizar essa análise, foram utilizados os dados de um serviço de proteção da criança dos Estados Unidos, que já possuía uma estrutura de dados que considerava os fatores de risco. No total, esse conjunto possui 37 itens avaliados em uma escala de 0 até 5 pontos, resumindo para a pesquisa esses valores apenas em dois, entre 0 e 2 pontos e entre 3 e 5 pontos (valor binário). As variáveis utilizadas podem se agrupar em: características da criança e do cuidador, relação entre eles, severidade e cronicidade da violência infantil, fatores econômicos e acesso do agressor à criança. Nesta investigação, foram utilizados 12.978 casos cadastrados nos centros de proteção de crianças, sendo que 21% (2.695) possuem 3 ou mais pontos (alto risco). As técnicas usadas foram: regressão linear (SEBER; LEE, 2012), regressão logística (HOSMER; LEMESHOW; STURDIVANT, 2013), e Rede neural do tipo MLP (RUMELHART; HINTON; WILLIAMS, 1985).

O artigo conclui que a MLP supera ambas as regressões. Também é importante ressaltar que a sensibilidade dos dados da MLP é maior. Portanto, com diferentes valores

de entrada haverá maior variação na saída, o que pode ser prejudicial para a previsão. Mesmo obtendo estes resultados, precisa-se descartar que houveram problemas na determinação da classificação de cada caso. Por exemplo, como saber quando o risco é baixo, médio ou alto? Isso envolveu refletir sobre como abordar o problema. Na maioria dos casos, trabalhou-se apenas com riscos baixos (entre 0 e 2 pontos) e altos (entre 3 e 5 pontos), limitando os possíveis resultados. Uma análise exaustiva foi feita com base na MLP, modificando as características dos dados de entrada e estrutura do modelo para obter conclusões diferentes, destacando a versatilidade deste tipo de técnica de predição. Contudo, afirma-se que dependendo do problema, as estatísticas convencionais podem obter os mesmos ou melhores resultados do que outras que derivam de IA.

A pesquisa realizada por Flaherty e Patterson (2003) visa desenvolver e treinar uma rede neural bayesiana a partir de dados de violência infantil física por meio de CV e comparar os resultados com um modelo alternativo obtido através do algoritmo de regressão logística. As hipóteses foram: “um modelo de predição de rede neural será significativamente melhor que um modelo de regressão logística para o mesmo conjunto de variáveis” e “todas as variáveis do conjunto de dados permitirão prever se houve violência”. Este estudo foi realizado com os casos de agressão física em crianças de famílias da força aérea dos Estados Unidos, sendo 5.612 instâncias no total, das quais 6,3% são indicadas como violência física recorrente e o resto em uma ocasião. Além disso, 13 fatores de risco presentes nos diferentes contextos que as crianças participam foram utilizados para a análise, os quais podem ser agrupados em características: da criança, dos pais e agressor, do entorno e do tratamento. Para criar os modelos de regressão e rede neural com quantidades equivalentes de resultados de saída para atos de violência recorrentes e não recorrentes, uma amostragem aleatória de 492 registros foi feita, que corresponde a 70% do total de casos em que houve agressões mais de uma vez. Para validar a precisão da classificação de ambos algoritmos foi utilizada a curva ROC.

Sobre os resultados, as áreas sob a curva ROC da rede neural e da regressão logística foram de 0,644 e 0,650, respectivamente. Deve-se notar que muitas modificações foram feitas na estrutura da rede neural, mas nenhuma obteve melhores resultados que a regressão. As duas hipóteses propostas no início não poderiam ser validadas, pois apenas 4 fatores de risco são preditores de violência recorrente (raça da vítima, renda familiar, fonte do encaminhamento inicial do caso e tempo que o caso permaneceu aberto após o encaminhamento inicial) e o modelo da rede neural não obteve os melhores resultados.

Sledjeski et al. (2008) também fizeram uma análise comparativa entre modelos. Neste caso, utilizaram regressão logística e *Classification And Regression Tree* (CART), tendo como objetivo demonstrar que esse último é um melhor preditor na recorrência de violência. Para obter os dados, foram entrevistadas 244 famílias de um centro de proteção de crianças nos Estados Unidos, destacando 24 atributos que estão relacionados com a

repetição da violência infantil, em que profissionais avaliaram cada um deles com um grau de risco: nenhum, pouco, moderado ou alto. Os fatores de risco foram agrupados em: severidade da violência, características da criança e do cuidador e influência da intervenção e do entorno da criança. No estudo, primeiro foi feita uma análise bivariada (BABBIE, 2013) de regressão logística tentando descobrir quais elementos eram os melhores preditores da recorrência. Foram detectados 9 relevantes, sendo os 3 com maiores *Odds Ratios* (OR) (RUDAS, 1998): histórico de violência, visibilidade da criança e relação pai-filho. Assim, gerou-se um modelo multivariado de regressão logística, responsável por atingir 37% de sensibilidade e 87% de especificidade. Posteriormente, foi feita a análise com a técnica CART com todos os atributos, obtendo 88% e 36% de sensibilidade e especificidade, respectivamente. Na sequência, realizou-se a mesma análise, mas apenas com os registros que não eram casos provados como violência infantil (novas instâncias), tentando determinar quais variáveis são mais significativas para esse conjunto de dados. Esses atributos foram: plano ativo de tratamento do caso, histórico de violência e visibilidade da criança. Assim, foi possível atingir uma sensibilidade de 87% e especificidade de 65%.

Comparado com a regressão logística, o algoritmo CART obteve uma melhor sensibilidade sempre. Portanto, é possível estabelecer que a CART é um melhor preditor nesse cenário. Alguns pontos fracos do trabalho que puderam repercutir nos resultados são: (1) o conjunto de dados tinha poucos registros para fazer uma análise desse tipo; e (2) a avaliação feita pelos profissionais no começo da pesquisa foi subjetiva (por exemplo, um elemento talvez fosse classificado com risco moderado por alguém, mas outra pessoa poderia classificar como alto).

Thurston e Miyamoto (2018) realizaram uma pesquisa de caso-controle (BORGAN et al., 2018) tentando demonstrar que árvores de decisão podem melhorar os resultados das ferramentas que avaliam o risco de violência infantil, isso através de um trabalho feito com crianças menores de 6 anos. A base de dados foi composta por 233 casos e 467 controle (700 registros no total). Somente foram considerados casos de violência graves, onde o resultado possível é da criança ser hospitalizada ou ela morrer. O conjunto de variáveis total não foi detalhado, mas sim os fatores de risco selecionados para gerar os modelos, escolhidos através de análises univariada e bivariada, relacionando-os com a violência grave. Essas variáveis foram: gênero da criança, número de crianças menores de 5 anos presentes na família, idade da mãe, problemas de saúde dos cuidadores, se a criança recebe ajuda dos serviços de proteção infantil, superlotação na moradia, violência doméstica prévia, prisão de algum dos cuidadores por violência, se a criança tem seguro de saúde e porcentagem de dias que a criança ficou em cuidados externos no período de interesse. Usou-se árvores de decisão para gerar os modelos, podendo a variável de saída tomar os valores de caso ou controle. Para a pesquisa considerou-se importante o nível de pobreza das famílias e se elas apresentavam comprovações de violência anteriores, pelo que foi criado um modelo para cada um desses elementos, além de detectar fatores de risco vinculados. O primeiro

resultado considerou a porcentagem de dias que a criança passou em cuidados externos e se ela tinha seguro de saúde. Quando esse tempo foi maior que 2%, concluiu-se que não existia relação entre pobreza e violência (OR de 0,72). Quando o tempo foi menor que 2% e a criança não tinha o seguro de saúde, foi estabelecida a relação (OR de 1,97). O segundo modelo (que relaciona comprovação de violência infantil com violência grave) detectou os mesmos dois fatores de risco do primeiro, além de verificar que existe um histórico prévio de violência infantil.

Assim, dois resultados relevantes puderam ser obtidos: (1) quando o tempo de cuidado externo é maior que 2% não se estabelece a relação entre comprovação prévia de violência infantil com violência grave (OR de 0,59); e (2) quando o tempo foi maior que 2% e existe histórico de violência infantil e a criança não apresenta seguro de saúde, ou seja, há uma vinculação muito significativa entre a comprovação de violência infantil e violência grave (OR de 2,97). Os dois modelos foram validados através da sensibilidade e a especificidade, obtendo valores de 17%/90% e 71%/49%, respectivamente. Algumas das relações dos fatores críticos encontradas são interessantes e não tinham sido estabelecidas em pesquisas anteriores, como é o tempo que a criança passa em cuidados externos aos da família e se tem seguro de saúde. Porém, o estudo não obteve resultados significativos. A sensibilidade média dos modelos foi de 44% e a especificidade média de 70%.

Chouldechova et al. (2018) tentaram melhorar a ferramenta preditiva de um centro de proteção de crianças dos Estados Unidos que contava com um modelo de regressão logística (HOSMER; LEMESHOW; STURDIVANT, 2013) que usava 71 variáveis diferentes para prever se uma notificação (*recall*) de violência infantil era de alto risco ou não. Estimou-se que esse modelo estava superestimado por duas razões principais: (1) construção errada dos conjuntos de treinamento e de teste; e (2) fatores de risco escolhidos sem considerar a relação deles com a violência infantil. Essa pesquisa utilizou o mesmo conjunto de dados com o que foi feita a regressão logística anterior, composto por 46.503 registros e 800 variáveis. Os fatores de risco utilizados não são declarados no artigo mas sim as categorias deles, sendo essas: demografia, bem-estar público, dados da prisão da localidade e estado de saúde de todos os envolvidos nas notificações. Foram gerados três modelos diferentes para a predição: SVM, RF, e XGBoost (*eXtreme Gradient Boosting*) (CHEN; GUESTRIN, 2016), validando-os através da AUC. A Tabela 2 apresenta os resultados obtidos por essas técnicas, além da regressão logística feita anteriormente, onde TP *recall* corresponde às notificações de alto risco verdadeiras (*true positive recall*) e FP *recall* às notificações de alto risco falsas (*false positive recall*) classificadas pelos modelos. Para as três colunas, os valores das técnicas aplicadas nesse estudo foram melhores que os da regressão logística, sendo a melhor obtida por XGBoost. Porém, esse aumento não é significativo, visto que uma sensibilidade de 0,61 não é considerada adequada para um modelo de ML.

Tabela 2 – Desempenho dos algoritmos usados para prever o risco de violência infantil de uma notificação (CHOULDECHOVA et al., 2018).

Técnica	AUC	TP <i>Recall</i>	FP <i>Recall</i>
Regressão Logística	0,70	0,49	0,21
SVM	0,77	0,57	0,20
Árvores de Decisão	0,77	0,58	0,20
XGBoost	0,80	0,61	0,19

Fonte: Chouldechova et al. (2018).

Os autores desse trabalho discutem dois assuntos muito importantes. O primeiro é a possibilidade de somente usar técnicas de ML para a predição da violência infantil quando um novo caso é recebido por algum centro de proteção de crianças (sem ajuda de especialistas no tema), com o objetivo de melhorar a velocidade com que essas organizações entregam respostas e fazem encaminhamentos. Sobre esse fato, os resultados das pesquisas que envolvem violência infantil com ML não tem sido bastante precisas a ponto de confiar unicamente nelas. O segundo assunto é o viés que podem ter os assistentes que encarregam-se dos casos de violência infantil, nos quais a gravidade às vezes será decidida somente pelo critério profissional dessa pessoa. Nessa pesquisa comprovou-se que quando se avaliam crianças de diferentes etnias (especificamente cor de pele), existe uma tendência a classificar com um maior risco os casos que envolvem famílias negras, sendo que os demais fatores de risco para ambos grupos são os mesmos. Esse fato também é estudado na pesquisa anterior (THURSTON; MIYAMOTO, 2018), em que determinou-se que quando o nome do investigador que está levando o caso de violência infantil é conhecido, a precisão do modelo aumenta, o que reflete uma falta de seriedade por parte dos assistentes.

3.2 Resumo dos Aspectos Relevantes dos Trabalhos Relacionados

Após analisar todas as pesquisas relacionadas com a predição de violência infantil através de ML para projetar uma investigação completa, é necessário sumarizar as principais características que definem um estudo dessa natureza. Esses elementos são: objetivo, tamanho da amostra usada, principais fatores (variáveis) usados, técnicas para gerar e validar os modelos, além de destacar os principais resultados obtidos nos trabalhos estudados. Todas essas informações dos artigos são resumidas na Tabela 3 e na Tabela 4.

Além da informação anterior, podem-se destacar algumas estatísticas gerais considerando todas as investigações revisadas, essas são:

- Considerando todos os trabalhos, a média de registros usados é de 24.955;

Tabela 3 – Resumo de artigos sobre trabalhos relacionados.

Autores	Objetivo do Estudo	Tamanho da Amostra	Categorias dos Fatores da Violência Usados	Téc. de Modelagem e Validação	Resultados Gerais
Altermeier et al. (1984)	Predizer a violência em futuros filhos de mulheres grávidas.	1.400 mulheres grávidas.	Saúde da criança e da mãe, apoio externo, atitude frente à gravidez, e consumo de substâncias (20 no total).	Regressão hierárquica; análise de correlação.	Se encontraram alguns fatores de risco relacionados à violência, mas nenhum com um valor significativo.
Frank et al. (1992).	Predizer a violência em crianças hospitalizadas.	749 crianças, 114 com sintomas de violência.	Sintomas de abuso físico (6), de abuso sexual (4), e negligência (2).	Árvore de decisão e estudo clínico; análise de sensibilidade.	A precisão do estudo clínico foi melhor, implicando que o modelo de árvore de decisão não fosse considerada no final.
Burrell et al. (1994).	(1) Determinar relações entre a violência e estresse parental, e (2) entender se a violência é diferente em famílias com crianças com necessidades especiais.	113 mães, 53 de crianças com necessidades especiais.	Enquetes relacionadas a: violência, estresse parental, percepção dos recursos, e apoio social.	Análise de regressão; análise de correlação.	Determinou-se a relação significativa entre o estresse parental e a violência, mas não houveram diferenças entre os dois grupos de famílias.
Rodriguez e Green (1997).	Relacionar o estresse parental e a expressão de raiva com a violência.	123 pais (divididos em duas amostras).	Enquetes relacionadas a: violência, estresse parental, e expressão de raiva.	Regressão hierárquica; análise de correlação.	Modelo relaciona as duas variáveis com a violência com uma alta correlação ($R > 0,8$), mas não permite ter a certeza de que houve agressão (resultados da CAPI podem não ser verdadeiros).
Little e Nixon (1998).	Predizer o nível de risco da violência das crianças que residem em centros de proteção.	20 crianças.	Características da violência recebida, dos cuidadores, e do impacto na criança.	Árvore de decisão; não houve validação do modelo.	Árvore que predize o nível do risco, mas foi gerado apenas com 20 registros e não teve validação.
DePanfilis e Zuravin (1999).	Predizer a recorrência da violência.	446 mães.	Colocação da criança, informações da violência, vulnerabilidade da criança, do cuidador, estresse familiar, e características da moradia (10 no total).	Regressão de Cox; <i>log-likelihood ratio</i> .	Modelo de supervivência com variáveis bem definidas (VPI e suporte social).
Marshall e English (2000).	Predizer o risco da violência com redes neurais e comparar esse resultado com regressões.	12.978 registros de centros de proteção de crianças.	Características da criança e do cuidador, relação entre eles, severidade e cronicidade da violência, e fatores econômicos (37 no total).	Regressão linear, regressão logística, e MLP; análise de sensibilidade.	Precisão de um 79% para MLP, e de 87% para regressão logística, usando 5 e 27 fatores de risco respectivamente.
Schwartz et al. (2004).	Predizer quais crianças apresentam “dano padrão” da violência.	1.767 registros de estudo sobre violência.	Não é detalhado o conjunto de fatores críticos; somente fala-se sobre a variável binária dependente, 1 se atingem-se os critérios e 0 se não.	MLP; análise de sensibilidade.	Rede neural com precisão de quase 90%, mas usando um conjunto de dados que não pode ser obtido em outros contextos (NIS).
Flaherty e Patterson (2003).	Determinar se redes neurais tem melhores resultados que regressão logística na predição da violência.	5.612 crianças.	Características da criança, dos pais, do agressor, do entorno, e do tratamento (13 no total).	Rede neural bayesiana e regressão logística; análise de sensibilidade e AUC.	R. logística teve melhores resultados que R. neurais; apenas 4 fatores foram significantes (raça da criança, renda familiar, e início e tempo transcorrido do caso).
Sledjeski et al. (2008).	Determinar se CART tem melhores resultados na predição da recorrência da violência que regressão logística.	244 famílias de um centro de proteção de crianças.	Severidade da violência, características da criança e do cuidador, e influencia da intervenção da violência e do entorno da criança (12 no total).	CART e regressão logística; análise de sensibilidade.	CART obteve melhores resultados em todas as análises feitas com um alta sensibilidade. O fator “visibilidade da criança” foi bem relacionado à violência.

Fonte: Elaborado pelo autor.

Tabela 4 – Continuação de resumo de artigos sobre trabalhos relacionados.

Autores	Objetivo do Estudo	Tamanho da Amostra	Categorias dos Fatores da Violência Usados	Téc. de Modelagem e Validação	Resultados Gerais
Begle et al. (2010).	Validar empiricamente modelos teóricos da violência.	610 cuidadores.	Caraterísticas da criança, dos cuidadores, do seu relacionamento, da família, e do entorno.	Regressão linear; valor qui-quadrado e análise de correlação.	Somente o modelo de risco acumulado teve uma precisão alta ao comparar com dados.
Dubowitz et al. (2011).	Identificar fatores de risco em famílias de baixa renda e predir a violência.	224 mães.	Caraterísticas da crianças, da mães, da família, e da comunidade.	Análise bivariada, regressão de Cox; valor qui-quadrado.	Criação de modelo de supervivência da violência (perceber quando há), e identificação de fatores de risco importantes.
Vaithianathan et al. (2013).	Validar a performance de um modelo preditivo de risco (PRM).	103.397 registros, 57.986 crianças.	Caraterísticas da criança, do cuidador principal, e do seu casal.	PRM (não especificado); curva ROC e AUC.	O modelo obteve uma precisão de 76%, mas não é dito qual foi usado na análise. Determina-se que o fator “raça” não está relacionado com a violência.
Horikawa et al. (2016).	Predizer a primeira ocorrência de violência com dados dos governamentais.	716 casos <i>boots-trapping</i> para acrescentar a amostra).	Idade da criança e do agressor, histórico de violência do agressor, estabilidade financeira na moradia, presença de vigilante à vítima, e a fonte de informação do caso.	Regressão logística; curva ROC e AUC.	Desempenho do modelo não foi alto (AUC de 0,66), além de que poucos fatores de risco foram usados e não considerou-se o fato da variação deles ao passar o tempo.
Daley et al. (2016).	Determinar quais são os lugares com maior risco de violência através de RTM.	Não é detalhado.	Caraterísticas de uma localidade específica: pobreza, violência doméstica, agressões agravadas, fugitivos, assassínios e delitos de drogas.	RTM; analisou-se a precisão através de novos casos.	Gerou-se um modelo que permite visualmente ver quais partes de uma localidade apresentam maior risco de violência, mas dificilmente pode ser replicado em outros lugares (fonte de dados única).
Amrit et al. (2017).	Predição da violência através de mineração de texto.	195.188 consultas, 13.170 crianças.	Estruturados: caraterísticas da criança, relacionamento familiar, e presunções de violência. Não estruturados: dados da criança e da dinâmica familiar.	NB, RF, e SVM; ANOVA e análise qui-quadrado, curva ROC e AUC.	Modelos permitem o uso de dados estruturados e não estruturados, logrando resultados da predição altos (AUC>0,9 para SVM).
Schwartz et al. (2017).	Melhorar a validação dos centros de proteção, tanto na predição da violência como no encaminhamentos dos casos.	78.394 casos.	Mais de 150 variáveis, não detalhando informação delas além do nome.	Árvores de decisão; curva ROC e AUC.	Os modelos conseguem bons resultados, mas não é detalhado que técnicas de ML foram usadas e as configurações delas, nem que fatores de risco utilizou cada um.
Thurston e Miyamoto (2018).	Demonstrar que árvores de decisão melhoram as ferramentas que avaliam o risco de violência entregando mais informação.	700 registros de caso-controle.	Caraterísticas da criança, dos cuidadores, se houve violência prévia no grupo familiar, e antecedentes delituais na família.	Análise multivariado e árvores de decisão; análise de sensibilidade.	Precisão dos modelos não é alta, não sendo adequados para predir a violência grave. Fatores de risco como contar com seguro de saúde e o tempo em cuidados externos não tinham sido estudados antes.
Chouldechova et al. (2018).	Melhorar ferramenta preditiva da violência de um centros de proteção de crianças.	46.503 registros.	Caraterísticas demográficas, bem-estar público, dados da prisão da localidade, e estado de saúde de todos os envolvidos nas notificações de violência.	SVM, RF, e <i>XGBoost</i> ; curva ROC e AUC.	Melhora na precisão em comparação a modelos anteriores, mas não significativamente. Interessante descoberta sobre o viés dos trabalhadores ao coletar informação dos casos de violência.

Fonte: Elaborado pelo autor.

- A quantidade de variáveis para gerar os modelos é muito desigual entre estudos, portanto não é possível determinar um número padrão de fatores de risco e de proteção;
- Sobre o assunto dos dados usados, 63% dos artigos utilizaram informações da criança, 58% do entorno, 53% da família, 37% dos cuidadores, 16% do perpetrador, e 11% do tratamento para a violência infantil;
- Sobre as técnicas implementadas, 53% dos trabalhos utilizaram algum tipo de regressão, 37% alguma variação de árvore de decisão, 16% algum tipo de rede neural, 11% SVM, e 21% outras técnicas (PRM, RTM, NB, e XGBoost);
- Sobre as técnicas de validação, 58% das pesquisas usaram análise de sensibilidade para validar os modelos, e 32% ainda adicionou o cálculo da AUC. Além disso, 21% utilizou uma análise de correlação entre os fatores de risco e a violência infantil;
- Os desempenhos dos modelos variam entre 50% e 90%, portanto que não existe um padrão de resultados esperados na predição da violência infantil através de ML.

3.3 Posicionamento da Presente Pesquisa

Da mesma forma que as outras pesquisas estudadas, esse trabalho tenta ajudar na melhora da predição da violência infantil através de análises de dados com técnicas computacionais. Porém, nessa investigação propõem-se três principais diferenças: (1) até agora não existem registros de trabalhos feitos dessa problemática na América Latina e o Caribe e, portanto, qualquer avanço em países dessa zona são estudos nunca antes realizados; (2) o fato de aplicar técnicas descritivas (sobre dados não supervisionados) também não tinha sido feito anteriormente; e (3) a aplicação de técnicas de ML sobre um conjunto de dados relacionado ao fenômeno da polivitimização.

Como foi descrito na [seção 1.4](#), esse trabalho tem base em dois conjuntos de dados diferentes, um supervisionado e outro não supervisionado. O primeiro obteve-se através de um centro de proteção de crianças e tem 9.771 instâncias com 30 atributos diferentes relacionados com: a criança e a sua situação escolar, sua família, se tem vinculação com alguma instituição de ajuda e a agressão sofrida (se a criança experimentou violência, o tipo de violência, e quem é o principal agressor). O segundo conjunto corresponde aos resultados da aplicação da primeira enquete sobre polivitimização em crianças e adolescentes no Chile, que foi obtido do portal de dados livres, *Centro de Estudios y Análisis del Delito*⁴ (CEAD) do governo desse país. Esse segundo conjunto possui 19.684 registros com avaliações de presença de 32 tipos diferentes de violência (vitimizações). Para cada uma

⁴ <http://cead.spd.gov.cl>

perguntou-se se tinha acontecido no último ano e na vida toda (64 atributos no total). Ademais, foram coletados alguns dados da criança e do seu entorno.

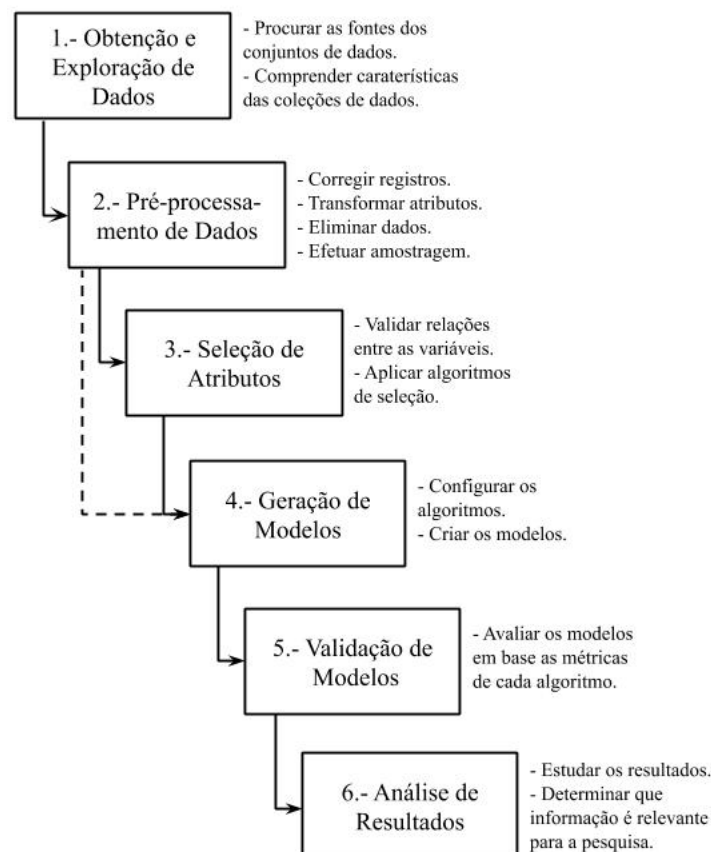
Considerando a natureza dos dados, precisaram ser usadas diferentes técnicas para poder trabalhar com ambos conjuntos. Para os registros supervisionados foi escolhido o algoritmo C4.5 (QUINLAN, 1993) e para os não supervisionados o algoritmo *Apriori* (AGRAWAL; IMIELIŃSKI; SWAMI, 1993). Os dois algoritmos foram selecionados com base no fato de que entregam regras simples como modelo de predição e descrição, respectivamente. Sobre os seguintes passos dessa pesquisa, no [Capítulo 4](#) é definida a exploração e pré-processamento dos dados, além das configurações usadas para os algoritmos e quais foram os experimentos feitos, continuando no [Capítulo 5](#) com a apresentação dos resultados e a sua respectiva análise.

4 Metodologia Proposta

Como foi visto nos capítulos anteriores, para desenvolver uma investigação de ML precisa-se de um conjunto de dados com o qual gerar os modelos. Ademais, é necessário um entendimento completo do problema de pesquisa e das variáveis disponíveis, o que permite determinar da melhor forma quais tipos de análises são possíveis de realizar e quais técnicas podem ser aplicadas.

Nesse capítulo será exposto o processo que permite obter os resultados dos modelos preditivos da violência infantil, detalhando cada uma das etapas realizadas. A [Figura 6](#) apresenta todas essas fases com os principais elementos de cada uma, as quais foram aplicadas para os dois conjuntos de dados utilizados.

Figura 6 – Etapas e fluxo que compõem a metodologia utilizada.



Fonte: Elaborado pelo autor.

A primeira etapa consiste na obtenção dos dados necessários para realizar um estudo desse tipo, sendo um modelo supervisionado e outro não supervisionado. Ademais, precisou-se compreender completamente a estrutura e quais deveriam ser as análises feitas para cada modelo. Entendendo qual seria a abordagem de trabalho sobre cada conjunto de

dados, foi necessário realizar a tarefa de pré-processamento para adequar os dados, sendo baseada em: correção e transformação de dados, eliminação de registros, amostragem, entre outras. Posteriormente, foi feita a seleção de atributos com a finalidade de escolher apenas os atributos que poderiam ser relevantes para a análise. Essa etapa foi aplicada somente para o conjunto supervisionado, o que explica o fluxo alternativo da etapa 2 da Figura 6. Em seguida, foram treinados os modelos para obter a maior quantidade de resultados possíveis e poder compará-los, o que realizou-se na fase de validação dos modelos. Por último, a metodologia conclui com a análise dos resultados.

A seguir, cada uma das seções desse capítulo corresponde a uma das fases desenvolvidas na investigação, detalhando os principais aspectos de cada uma. As etapas 3, 4 e 5 somente apresentam as configurações das técnicas de seleção, modelagem e de validação respectivamente, deixando os resultados para o seguinte capítulo. É importante mencionar também que a ferramenta utilizada para a geração e validação dos modelos foi o *Waikato Environment for Knowledge Analysis* 3.8 (WEKA)¹, principalmente por ter disponível todos os algoritmos necessários para essa pesquisa.

4.1 Obtenção e Exploração dos Dados

Essa fase corresponde à obtenção dos dados necessários para o estudo, bem como descobrir quais são as principais características de cada conjunto de dados. Foram obtidas duas coleções de dados diferentes: uma composta por registros armazenados por um centro de proteção de crianças do Chile e outra correspondente às respostas de uma enquete sobre polivitimização aplicada em crianças e adolescentes desse mesmo país. As duas subseções seguintes apresentam o desenvolvimento dessa primeira etapa.

4.1.1 Centro de Proteção de Crianças

Um dos centros de proteção de crianças que existe no Chile, equivalente a um *Organismo Colaborador Acreditado* (OCA) do governo desse país (instituições encarregadas do cuidado e tratamento de crianças e adolescentes que sofrem de violações de direitos ou apresentam problemas delituais), considerou interessante esse projeto e se dispôs a entregar uma parte dos registros sobre crianças que participaram dos programas que essa organização realiza.

O conjunto de dados fornecido para a pesquisa é composto por 9.771 registros e possui 30 variáveis diferentes, sendo essas:

1. Ano (numérica) – ano em que o caso foi atendido;

¹ www.cs.waikato.ac.nz/ml/index.html

2. Zona (categórica) – setor do Chile ao qual pertence a criança (norte, centro ou sul);
3. Linha de Intervenção (categórico) – motivo pelo qual a criança participa na organização (educativo, justiça, proteção de direitos, ou tratamento de drogas);
4. Nome do projeto (categórico) – nome do projeto que participa a criança (mais de 90 projetos diferentes);
5. Idade (numérica) – quantidade de anos da criança;
6. Sexo (categórica) – sexo da criança (masculino ou feminino);
7. Nacionalidade (categórica) – país de origem da criança (7 diferentes países);
8. Região (categórica) – estado no qual mora a criança (11 diferentes estados);
9. Etnia (categórica) – grupo étnico ao qual pertence a criança (9 diferentes etnias);
10. Apresenta necessidades especiais (binária) – se a criança tem alguma necessidade especial;
11. Zona na qual mora (categórica) – tipo de espaço no qual a criança mora (rural ou urbano);
12. Superlotação (binária) – se existe superlotação na moradia da criança;
13. Situação socioeconômica (categórica) – estado de pobreza da criança (não pobreza, pobreza ou pobreza extrema);
14. Tipo de família (categórica) – composição da família da criança (8 tipos diferentes);
15. Número de filhos (numérica) – quantidade de filhos que tem a criança;
16. Adulto responsável (categórica) – pessoa maior que é responsável pela criança (7 possíveis pessoas);
17. Número de programas de proteção nos quais tem participado (numérica) – quantidade de programas de proteção infantil que a criança já participou;
18. Vinculação com redes comunitárias (binária) – se a criança teve ou tem alguma vinculação com organizações comunitárias;
19. Vinculação com redes institucionais (binária) – se a criança teve ou tem alguma vinculação com organizações institucionais;
20. Último ano aprovado (categórica) – último ano de escola aprovado pela criança (16 diferentes níveis);

21. Situação escolar atual (categórica) – estado dos estudos da criança no momento de registrar ela na organização (6 possíveis estados);
22. Apresentou saída escolar (binária) – se a criança tem deixou os estudos em algum momento;
23. Repetência (binária) – se a criança repetiu algum ano da escola;
24. Projeções educacionais (categórica) – pretensões sobre educação da criança para ela mesma (7 diferentes projeções);
25. Socialização em rua (binária) – se a criança adquiriu hábitos sociais e culturais em ambientes de rua;
26. Apresenta consumo de substâncias (binária) – se a criança consome algum tipo de substância;
27. Tem cometido delitos (binária) – se a criança cometeu algum delito;
28. Sofreu violência (binária) – se a criança tem sofreu algum tipo de violência;
29. Tipo de violência (categórica) – tipos de violência que a criança pode ter sofrido (15 tipos diferentes); e
30. Principal agressor (categórica) – agressor da violência na criança (8 possibilidades).

Para ter uma melhor noção conceitual dos fatores de violência e também para poder compará-los com os dos estudos relacionados, os atributos do conjunto de dados foram agrupados nas categorias mencionadas na revisão feita no [Capítulo 3](#):

- Da criança – idade, sexo, nacionalidade, etnia, e apresenta necessidades especiais;
- Da família – zona, região, zona na que mora, superlotação, situação socioeconômica, tipo de família, número de filhos e adulto responsável;
- Do entorno – ano, linha de intervenção, nome do projeto, número de programas de proteção nos que tem participado, vinculação com redes comunitárias e vinculação com redes institucionais;
- Escolares da criança – último ano aprovado, situação escolar atual, apresentou saída escolar, repetência e projeções educacionais; e
- Violência que tem sofrido a criança – socialização (ou permanência) em rua, apresenta consumo de substâncias, tem cometido delitos, sofreu violência, tipo de violência e principal agressor.

Ao fazer uma exploração inicial dos dados, percebeu-se que todos os registros de crianças menores de 13 anos têm o valor *sim* no atributo preditor *sofreu violência*. Desse modo, decidiu-se dividir o conjunto de dados em dois: um com os casos entre 0 e 12 anos e outro com as crianças de 13 ou mais anos. Portanto, para o primeiro subconjunto não foi possível prever se uma criança sofreu violência, porém sim o tipo de violência e o principal agressor.

4.1.2 Questionário sobre Polivitimização

O segundo conjunto de dados corresponde à Primeira Enquete Nacional de Polivitimização em Crianças e Adolescentes ([SUBSECRETARÍA DE PREVENCIÓN DEL DELITO, 2018](#)) aplicada no Chile, que foi baseada no JVQ ([FINKELHOR et al., 2005](#)). Este questionário Avalia 32 tipos de vitimizações de forma binária (se aconteceu ou não esse tipo de violência), as quais foram categorizadas em 6 grupos.

- Crimes Comuns (CC) – Roubo sem uso da força (CC1), roubo com uso da força (CC2), estragar algum elemento pessoal intencionalmente (CC3), ameaça ou percepção de dano (CC4), ataque físico com objetos (CC5), ataque físico sem objetos (CC6) e ameaça por alguma característica própria (CC7).
- Efetuadas pelos Cuidadores (Cuidadores) – Ofender-se por insulto de um adulto próximo (Cuidadores1), ataque físico por um adulto próximo (Cuidadores2), sentir-se mal por descuido de adultos com quem vive (Cuidadores3) e mantido afastado ou escondido de seu pai ou mãe (Cuidadores4).
- Efetuadas pelos Pares (Pares) – Ataque físico de uma criança ou adolescente (Pares1), ataque físico de um grupo de crianças ou adolescentes (Pares2), sentir-se mal por insulto de um grupo de crianças ou adolescentes (Pares3), imposição para fazer coisas que não se querem (Pares4) e ataque físico do cônjuge (Pares5).
- Sexuais – Práticas sexuais com pessoas maiores de 18 anos de idade com consentimento (Sexuais1), ofender-se por *bullying* sexual não feito através da internet (Sexuais2), forçado a olhar para partes íntimas pela força ou surpresa (Sexuais3), forçado a fazer coisas de natureza sexual por uma criança ou adolescente (Sexuais4), tocado ou tentativa de toque de partes íntimas por um adulto estranho (Sexuais5), tocado ou tentativa de toque de partes íntimas por um adulto conhecido (Sexuais6) e tento ou participação de relações sexuais sem o seu próprio consentimento (Sexuais7).
- Indiretas – Presenciar roubo em casa (Indiretas1), presenciar violência (Indiretas2), presenciar discriminação (Indiretas3), presenciar ataque físico sem armas (Indiretas4), presenciar ataque físico com armas (Indiretas5), presenciar ataques físicos

entre seus próprios cuidadores (Indiretas6) e testemunhar agressões físicas dos pais a irmãos (Indiretas7).

- Digitais – *Cyberbulling* (Digitais1) e assédio sexual na internet (Digitais2).

Para todas as vitimizações foi perguntado se aconteceu alguma vez na vida e durante os últimos 12 meses (sim ou não para as duas temporalidades). Para essa última questão também foi perguntado a quantidade de vezes, tendo como respostas: nunca, 1 vez, 2 ou 3 vezes, ao menos uma vez por mês, ao menos 1 vez por semana, ou todos os dias. Além dos tipos de violência, o questionário avaliou o estado de autoestima e depressão através da Escala de Autoestima de Rosenberg (*Rosenberg Self-Esteem Scale*, RSES) (ROSENBERG, 1965) e a Escala de Autoavaliação de Depressão de Birlson (*Depression Self-Rating Scale*, DSRS) (BIRLESON, 1981) respectivamente, ambas adaptadas ao contexto do Chile. Também foram coletados alguns dados demográficos das crianças e das suas famílias.

Sobre a amostra, foram no total 19.867 crianças de 699 escolas dos 2 últimos anos do ensino fundamental e dos 3 primeiros do ensino médio, abarcando todo o território do Chile. Browne et al. (2018) fizeram um estudo multivariado sobre esse conjunto de dados, encontrando basicamente correlações entre os tipos de vitimizações e alguns aspectos psicológicos das crianças, mas não foi feita nenhuma análise preditiva nem descritiva.

4.2 Pré-processamento de Dados

A segunda etapa da metodologia foi a realização de diferentes tarefas de pré-processamento necessárias para aplicar os algoritmos selecionados sobre os conjuntos de dados de forma adequada. Essa etapa foi realizada para ambas as coleções de dados, porém houve maior trabalho no processamento dos registros do centro de proteção.

4.2.1 Conjunto de Dados Utilizado para Método Supervisionado

A princípio, buscou-se preencher os campos vazios (*missing data*) e corrigir algumas inconsistências, isso utilizando a informação que entregou o centro de proteção, além de simples suposições geradas a partir do mesmo conjunto de dados. Essas ações foram:

1. Nos registros com *superlotação = sim*, preencheu-se o valor de *situação socioeconômica = pobreza* (38 elementos);
2. Registros em que a *situação escolar atual = retirada da escola*, preencheu-se a *saída escolar = sim* (105 elementos);

3. Nos registros com *apresentou saída escolar = sim*, preencheu-se o valor de *repetência = sim* (524 elementos);
4. Nos registros que em *situação escolar atual = retirada da escola*, preencheu-se o valor de *projeções educacionais = não* (206 elementos);
5. Nos registros com valor de *intervenção = tratamento de drogas*, preencheu-se o valor de *apresenta consumo de substâncias = sim* (56 elementos);
6. Nos registros com valor de *intervenção = justiça juvenil*, preencheu-se o valor de *cometido delitos = sim* (114 elementos);
7. Nos registros com valor de *intervenção = proteção de direitos*, preencheu-se o valor de *sofreu violência = sim* (272 elementos);
8. Nos registros com valor de *tipo de violência ≠ vazio* ou *não*, preencheu-se o valor de *sofreu violência = sim* (23 elementos);
9. Nos registros com valor de *sofreu violência = não*, preencheu-se o valor de *tipo de violência = não* (1.389 elementos);
10. Nos registros com valor de *sofreu violência = não*, preencheu-se o valor de *principal agressor = não* (1.612 elementos).

Na sequência, foram realizadas algumas transformações no conjunto de dados, principalmente modificações de tipos e agregação de valores nos atributos.

1. Na variável *nacionalidade*, somente 2,4% dos valores não corresponde a chileno. Portanto, foi trocado o nome do atributo para *chileno* e transformou-se o atributo em binário (*sim* se a criança é chilena e *não* se tem outra nacionalidade).
2. O mesmo procedimento foi feito com a variável *etnia*, em que somente tinha 7,4% de valores com alguma categoria *indígena*. Assim, trocou-se o nome do atributo para *indígena* e transformou-o em binário (*sim* se a criança é de alguma etnia indígena e *não* se não é).
3. Os atributos *vinculação com redes comunitárias* e *vinculação com redes institucionais* foram unificados em apenas um de natureza binária chamado *vinculação com redes* (*sim* se a criança já participou de redes e *não* se nunca participou ou tem participado), o que considerou-se fazer já que não tem relevância o fato de ser uma rede comunitária ou institucional.
4. Criou-se um novo atributo do *tipo de violência*, porém com valores agregados. Foram consideradas as quatro categorias de violência infantil definidas no [Capítulo 2](#):

violência física e psicológica, abuso sexual e negligência como possíveis valores; atribuindo cada um dos tipos de abuso originais do conjunto alguma dessas categorias agrupadas.

Ademais, foram eliminados os atributos que não têm valor para a análise da predição da violência infantil através de ML. O primeiro foi o ano de ingresso, pois saber quando a criança entrou no centro não entrega informações relevantes, além de que apenas tinham-se três diferentes anos; o segundo é a linha de intervenção, pois não entrega informação da violência explicitamente e informações sobre consumo de substâncias e atos criminosos foram recuperadas na primeira etapa do pré-processamento; o terceiro foi o nome do projeto, pois descreve principalmente a cidade onde se executa os diferentes atendimentos às crianças; e a vinculação com redes comunitárias e institucionais foram eliminadas depois de unificá-las no atributo *vinculação com redes*. Posteriormente eliminaram-se instâncias com valores vazios. Foram descartados os registros sem valor em tipo de violência (1.154 no total), pois não permitiriam fazer a categorização da violência. Ademais, foram eliminadas todas as instâncias que tivessem vazios mais de 3 atributos, sendo essas um total de 668. Também considerou-se que as variáveis idade e sexo não podem ter elementos em branco, eliminando 9 instâncias.

Após esse pré-processamento, a quantidade de atributos e registros resultantes foi de 27 e 7.940, respectivamente, o que corresponde ao 81,26% da amostra inicial. É importante mencionar que essa totalidade de instâncias varia em relação ao momento de fazer a análise a depender do que se busca predizer. Por exemplo, quando se interessa determinar o tipo de violência, são excluídas todas as instâncias que não apresentam violência, o que diminui o tamanho da amostra. Outra consideração foi o desequilíbrio das quantidades de valores das variáveis preditoras (o que em alguns casos é significativo), motivo que implicou realizar as tarefas de amostragem, tanto *oversampling* como *undersampling*, para balancear a quantidade de classes, sempre baseado na que tinha menos registros. Essa última técnica aplicou-se através do método SMOTE (*Synthetic Minority Over-sampling TEchnique*) (CHAWLA et al., 2002) incorporado na ferramenta WEKA. Os resultados podem ser revisados na Tabela 5 da seção 5.1.

4.2.2 Conjuntos de Dados Utilizado para Método Não Supervisionado

Para esse conjunto de dados foram aplicadas duas tarefas de pré-processamento, uma de eliminação de atributos e registros, e outra de agregação de valores.

Sobre a eliminação de dados, primeiro foram descartadas todas as variáveis que não eram referentes às vitimizações, isso porque não têm o formato que permite aplicar regras de associação sobre esses elementos, deixando 64 atributos no total (32 vitimizações avaliadas em vida e ano). Na sequência, eliminaram-se todos os registros que não

apresentavam alguma resposta sobre as vitimizações sofridas em vida e ano, restando no total 14.700 instâncias.

A segunda tarefa foi criar duas agregações diferentes das vitimizações baseadas nos tipos de violência, isso pelo grande número de atributos que se tinha. Primeiro, foram consideradas as mesmas seis categorias que a enquete inicial define para classificar as diferentes violências estabelecidas na [subseção 4.1.2](#). A segunda agregação feita considerou os tipos de violência gerais definidos no [Capítulo 2](#): violência física e psicológica, abuso sexual, e negligência. A rotulação feita dos atributos foi:

- Violência Física: CC2, CC5, CC6, CC7, Cuidadores2, Pares1, Pares2, e Pares5;
- Violência Psicológica: CC1, CC3, CC4, Cuidadores1, Pares3, Pares4, Sexuais2, Indiretas6, Indiretas7, Digitais1, e Digitais2;
- Abuso Sexual: Sexuais1, Sexuais3, Sexuais4, Sexuais5, Sexuais6, e Sexuais7; e
- Negligência: Cuidadores3, Cuidadores4, Indiretas1, Indiretas2, Indiretas3, Indiretas4, e Indiretas5.

Assim, os novos conjuntos de dados têm 6 e 4 atributos respectivamente, sendo todos valores binários que possibilitam a utilização do algoritmo *Apriori*, mantendo para todos os conjuntos a mesma quantidade de registros (14.700).

4.3 Seleção de Atributos

Essa fase implica apenas em selecionar alguns atributos da coleção inicial para otimizar tanto o procedimento de treinamento dos modelos, como os resultados dos mesmos. Nesse caso, considerando a quantidade de atributos que conformavam ambos os conjuntos do estudo, estimou-se fazer a etapa de seleção de atributos principalmente com a finalidade de melhorar o desempenho dos modelos e os resultados em geral. Porém, fazer esse processo sobre as variáveis da coleção de polivitimização implicaria em dar uma maior importância a alguns tipos de violência que outros, o que não permitiria uma análise adequada. Portanto, a seleção de atributos apenas foi aplicada sobre o conjunto todo do centro de proteção de crianças. Fizeram-se três tipos diferentes de seleção de variáveis, através de um método de filtro, outro *wrapper*, e um último baseado na informação entregue pelo estado da arte.

- Método de Filtro – Utilizou-se o algoritmo de correlação de atributos do WEKA (*CorrelationAttributeEval*), que determina a dependência de todas as variáveis com o elemento que interessa ser predito, escolhendo as 5 mais relacionadas para fazer a análise.

- Método *Wrapper* – Foi usado o algoritmo *WrapperSubsetEval* do WEKA, escolhendo o método J48 (equivalente ao C4.5) para obter os subconjuntos de dados. Ademais, os atributos foram escolhidos através de um processo *forward selection*, que implica que começou-se com um conjunto vazio de variáveis e em cada iteração do algoritmo adicionou-se um elemento até que o desempenho de classificação do modelo não melhorasse.
- Baseado no Estado da Arte – Consideraram-se a maioria das variáveis obtidas, não utilizando apenas as que nunca foram mencionadas nos artigos que compõem o estado da arte, que são: *nacionalidade (chileno)*, *zona*, *região*, *último ano aprovado* e *projeções educacionais*. Para esse caso sempre foram considerados os mesmos atributos para gerar o modelo nas diferentes execuções realizadas do algoritmo C4.5.

É importante lembrar que se têm três atributos preditores nesse conjunto de dados (se a criança sofreu violência, tipo de violência e principal agressor). Portanto, foram aplicados os três tipos de seleção de atributo para cada uma dessas variáveis, ademais de uma execução sem seleção de atributos como detalha a [Figura 6](#). Conseqüentemente, tiveram-se quatro saídas diferentes nessa parte do trabalho.

4.4 Geração de Modelos

Essa etapa corresponde ao treinamento dos modelos que permitem fazer as previsões da violência infantil. A criação desses foi feita de maneira separada para cada um dos conjuntos de dados. Ainda que ambos estejam vinculados ao mesmo problema, a conformação de um em relação ao outro é muito diferente, isso principalmente pelos tipos de dados e abordagem que cada coleção tem (problema supervisionado e não supervisionado respectivamente).

4.4.1 Algoritmo C4.5 – Classificação

Como foi discutido no [Capítulo 2](#), para criar o modelo preditor do conjunto de dados do centro de proteção utilizou-se o algoritmo C4.5, o que permite encontrar regras que satisfazem os atributos para classificar um fenômeno em particular. Essa técnica na ferramenta WEKA representa-se pelo módulo J48.

O algoritmo foi configurado de três formas diferentes para testar o seu desempenho: com a configuração padrão do WEKA, limitando o tamanho da árvore resultante, e sem limitá-lo. Os elementos considerados para as diferentes configurações são: a quantidade mínima de instâncias por nó folha (*minNumObj*) e se a árvore é podada ou não (*unpruned*). Essas configurações são sumarizadas na sequência:

- Configuração Padrão – Não foram alterados os valores iniciais do algoritmo que a ferramenta tem por defeito, sendo esses: dois elementos para *minNumObj* e com a poda da árvore habilitada;
- Limitando o Tamanho da Árvore – Essa configuração foi aplicada para obter uma árvore pequena e simples de interpretar, tentando manter o desempenho de uma árvore maior. Os valores foram de 100 para *minNumObj* e com a poda habilitada;
- Sem Limitações no Tamanho da Árvore – Tentando melhorar o desempenho do modelo, foi aplicada essa configuração (uma maior árvore terá um maior ajuste sobre os dados), com *minNumObj* igual a 2 e sem a poda da árvore habilitada.

4.4.2 Algoritmo *Apriori* – Regras de Associação

Da mesma forma que para a análise supervisionada, para gerar as regras de associação utilizou-se a ferramenta WEKA 3.8, mas aplicando o algoritmo *Apriori*. Sobre as execuções realizadas, todas foram feitas com a mesma configuração, com valores mínimos para o suporte e a confiança de 0,05 e 0,75, respectivamente. Com isso, buscou-se conseguir um número relativamente elevado de regras para cada uma das 6 execuções do algoritmo (3 diferentes agrupamentos para os dados correspondentes às vitimizações vida e ano), visto que um suporte baixo permite considerar regras das quais seus antecedentes não têm muita presença no conjunto de dados total (nesse caso, regras com 5% de presença ou mais são válidas) e uma confiança elevada permite validar a taxa de acompanhamento do consequente (como mínimo terá que ter 75% de presença). No caso que esses valores gerem uma quantidade excessiva de regras, sempre será possível selecionar apenas as que apresentem melhores resultados nas métricas de avaliação.

4.5 Validação de Resultados

Nessa etapa da metodologia busca-se avaliar os modelos treinados anteriormente com o objetivo de estabelecer se são ou não adequados para a tarefa de predição da violência infantil. Considerando o trabalho separado que foi feito para cada um dos conjuntos de dados, a validação dos resultados também tem que ser feita em duas partes, principalmente porque modelos de classificação e de regras de associação não são avaliados da mesma forma.

4.5.1 Modelo de Classificação

Conforme exposto no [Capítulo 2](#), existem variadas técnicas para validar modelos gerados por técnicas supervisionadas, as que são aplicáveis ao algoritmo C4.5 de árvores de decisão (ROKACH; MAIMON, 2008). Utilizou-se *k-fold CV*, o que permite criar diferentes

subconjuntos de treinamento e teste a partir da coleção inicial, para assim gerar uma árvore para cada um e validar o desempenho da predição desses modelos (KOHAVI et al., 1995). Para esta parte da análise usou-se um k de 10 (quantidade de subconjuntos).

Para avaliar o modelo resultante do processo de CV, consideraram-se as métricas: predição, sensibilidade e especificidade. Todas elas foram geradas através dos valores TP, FP, TN, e FN que classificou o modelo. A partir desses elementos, são calculadas as AUC, ROC e PR, permitindo determinar, de uma maneira mais objetiva, a qualidade do modelo. Finalmente, estimou-se a MH das duas AUC para obter apenas um valor para medir a qualidade dos modelos de classificação.

4.5.2 Modelo de Regras de Associação

No processo de avaliação de uma regra de associação podem ser consideradas as duas métricas que permitem criá-las, o suporte e a confiança. Como foi discutido no Capítulo 2, o suporte corresponde à taxa de presença do conjunto antecedente da regra sobre todos os registros, procurando assim um valor elevado. Da mesma forma busca-se uma confiança significativa, o que reflete em que o consequente estará presente baseado no antecedente com uma alta probabilidade. Portanto, é esperado que ambos valores sejam próximos a 1. Ademais de usar o suporte e a confiança para avaliar as regras, utilizaram-se as métricas definidas na subseção 2.2.4.2: *lift*, *leverage*, e *conviction*. Por meio delas buscou-se apenas manter as que através de uma avaliação matemática tiveram desempenhos ótimos, e não conservar um grande número de regras de associação que puderam não ser representativas do conjunto de dados.

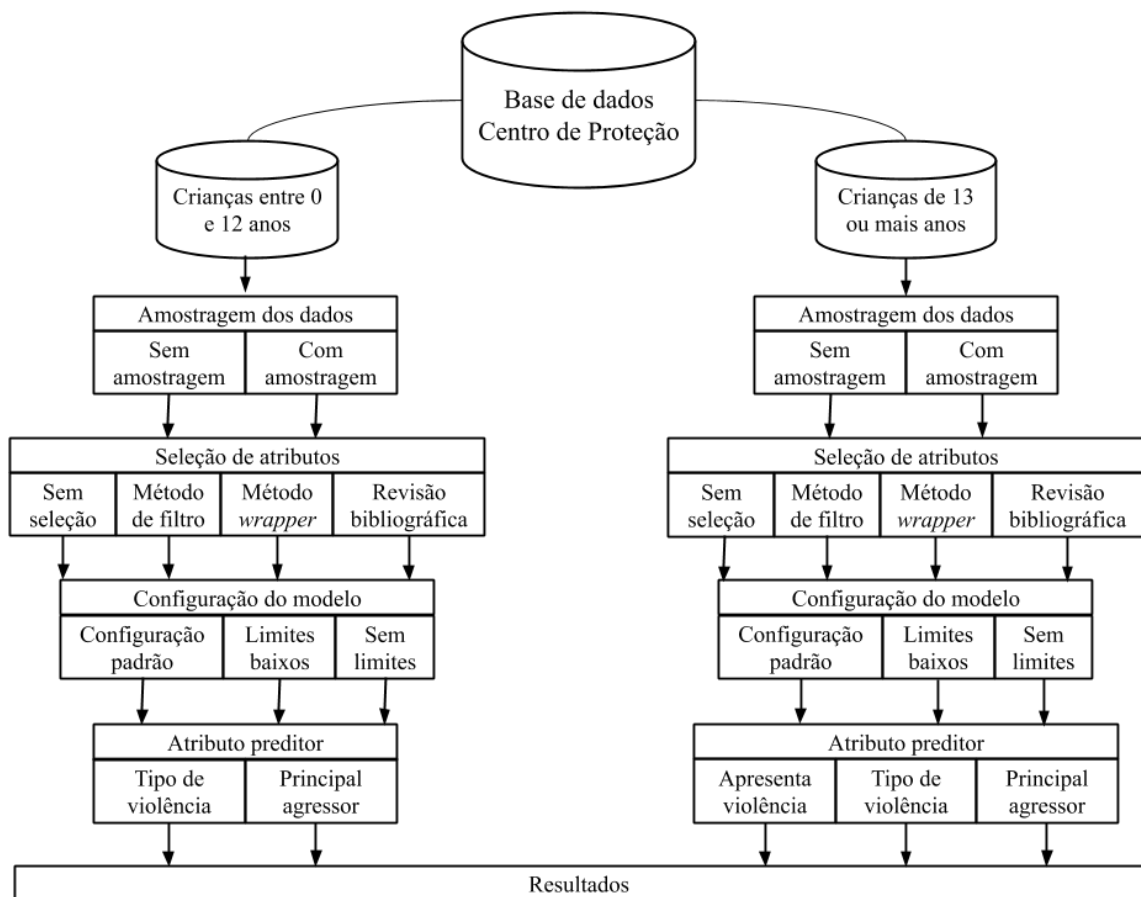
5 Resultados e Discussões

Esse capítulo apresenta os resultados conseguidos para os diferentes ajustes descritos na [subseção 4.4.1](#) e na [subseção 4.4.2](#) dos algoritmos *C4.5* e *Apriori* respectivamente. Para expor os resultados e a análise feita, as duas seções que seguem irão detalhar esses aspectos para cada uma das coleções de dados utilizadas.

5.1 Centro de Proteção de Crianças

Como foi detalhado no capítulo anterior, consideraram-se diferentes elementos para determinar as variações das execuções do algoritmo J48 do WEKA sobre os dados do centro de proteção. A [Figura 7](#) agrupa todos esses detalhando também o fluxo para encontrar os diferentes *inputs* da ferramenta.

Figura 7 – Fluxograma da metodologia proposta à tarefa de classificação, o que dividi-se em dois caminhos com conjuntos de dados diferentes.



Fonte: Elaborado pelo autor.

Cabe lembrar que a coleção inicial de dados foi dividida em duas, uma com os re-

gistros de crianças de 12 ou menos anos e outra com os maiores que essa idade. Para ambos os subconjuntos foram feitas as tarefas de: escolher um atributo preditor, tendo a segunda coleção um a mais (todas as crianças com 12 anos ou menos sim apresentavam violência); amostragem de dados tentando balancear a quantidade de classes; selecionar as variáveis relevantes baseado em diferentes considerações; e aplicar as diferentes configurações para executar o software. Portanto, foram 48 execuções para o primeiro subconjunto e 72 para o segundo (120 diferentes resultados no total). A [Tabela 5](#) apresenta esses resultados, dos quais são detalhados os valores obtidos para: a porcentagem de acertos (%), a MH das AUC ROC e PR, e o tamanho da árvore em quantidade de nós (nós folha/nós totais). O número no final da descrição de cada conjunto (na coluna atributo preditor) corresponde à quantidade de classes desse mesmo, e o que aparece na coluna amostragem ao total de registros resultantes depois de aplicar as técnicas de *oversampling* e *undersampling*.

Nota-se que o tamanho da maioria dos modelos é consideravelmente grande, portanto apresentar todos em formato gráfico não é pertinente nesse documento. Porém, para exemplificar os resultados a [Figura 8](#) descreve uma árvore menor (11 nós no total, com taxa de acertos de 81,3%, e MH de 0,79) na predição de presença de violência com a configuração: sem amostragem, método de filtro na seleção de atributos e sem limites no seu tamanho. Ademais, na [seção A.1](#) do [Apêndice A](#) pode ser visto o menor modelo (em formato de regra) obtido para cada um dos 5 atributos preditores analisados.

Para compreender como essa árvore faz a classificação de um caso, precisa-se seguir o fluxo que tem. A primeira regra especifica se a criança teve socialização em rua ou não. Caso o valor é *sim*, automaticamente o modelo classifica o registro como um possível caso de apresentar vulneração com uma taxa de 91,5% (o primeiro número que aparece entre parêntesis corresponde aos casos classificados corretamente e o segundo aos errados). Caso contrário, precisa ser aplicada a próxima regra, se tem cometido delito ou não, e novamente as respostas podem ser *sim* ou *não*. Se não tem cometido delito é classificado como *sim* (90,0%). Para o fluxo da esquerda (*tem cometido delitos = sim*), avalia-se o número de programas nos que participou a criança: se esse valor é maior a 0 a resposta do modelo é positiva (77,0%), caso contrario precisa ser avaliado se a criança teve vinculação com redes de apoio. Quando esse atributo toma o valor *não*, o caso é classificado negativamente (76,9%), e se é *sim* vai ser avaliado o valor de superlotação. Quando a criança não possui essa característica, o caso não apresenta violência (71,1%), e quando sua moradia sim possui superlotação classifica-se com *sim* (67,1%).

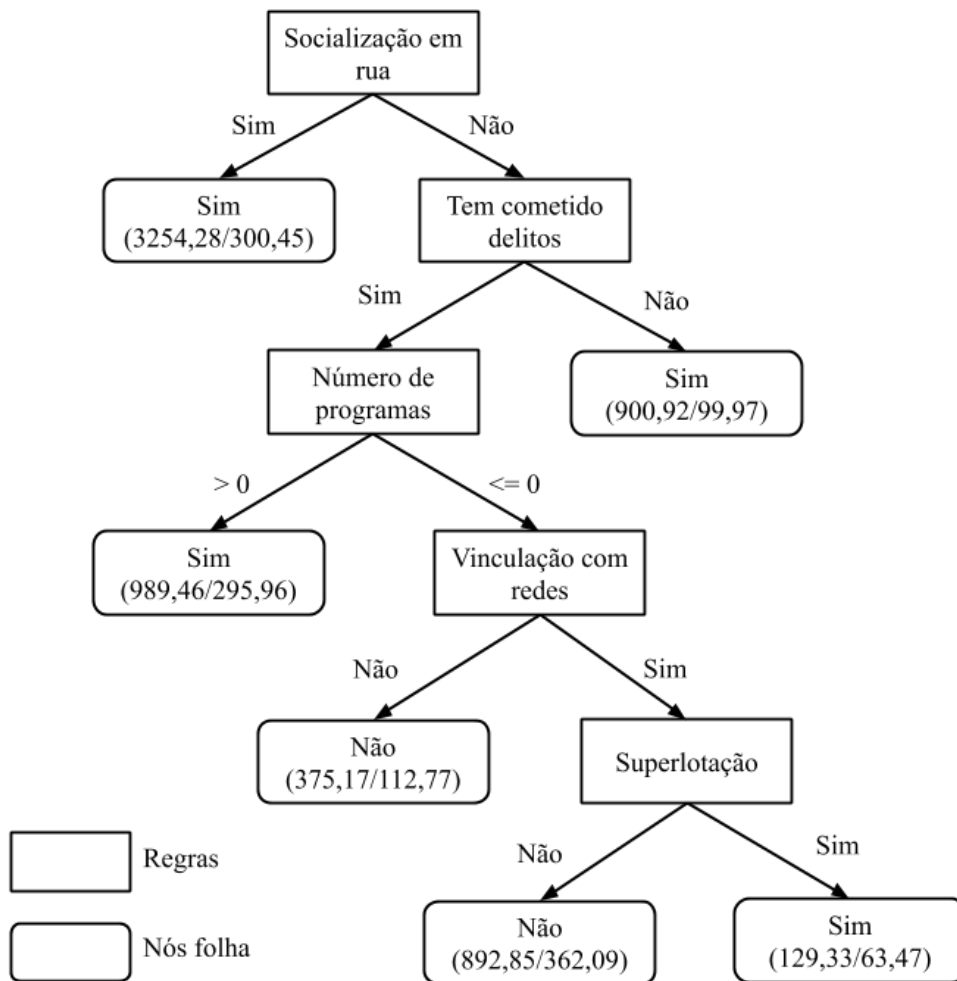
Retomando os resultados apresentados na [Tabela 5](#), dos cinco diferentes atributos preditores que a compõem, os dois primeiros correspondem aos registros de crianças que têm 12 ou menos anos (tipo de violência e principal agressor) e os três seguintes aos de crianças de mais de 12 anos (presença de violência, tipo de violência e principal agressor). Portanto, não é possível comparar os resultados obtidos entre esses atributos. Porém,

Tabela 5 – Resultados obtidos ao aplicar três configurações distintas do algoritmo J48 aos diferentes conjuntos de dados, especificando o número de classes e registros de cada um e destacando os melhores resultados para cada atributo preditor.

Atributo Preditor	Conjuntos de Dados		Configuração e Métricas de Avaliação											
			Seleção de Atributos			Padrão			Tamanho da Árvore Limitado			Tamanho da Árvore Ilimitado		
			%	MH	Nós	%	MH	Nós	%	MH	Nós	%	MH	Nós
Tipo de Violência com Registros de Crianças de 12 ou menos anos (4)	Sem Amostragem (1.375)	Sem Seleção	53,98	0,606	275/369	45,60	0,512	16/19	51,49	0,59	929/1143	51,49	0,59	929/1143
		Método de Filtro	41,67	0,474	52/77	39,85	0,435	8/10	39,85	0,485	84/128	39,85	0,485	84/128
	Com Amostragem (1.816)	Método Wrapper	55,13	0,627	97/148	46,04	0,516	16/19	54,84	0,63	157/239	54,84	0,63	157/239
		Revisão Bibliográfica	44,36	0,504	174/274	40,22	0,458	17/21	42,55	0,5	426/672	42,55	0,5	426/672
		Sem Seleção	59,69	0,669	430/568	45,37	0,515	18/24	61,01	0,664	887/1119	61,01	0,664	887/1119
		Método de Filtro	43,67	0,522	38/52	39,59	0,475	11/15	44,38	0,531	88/120	44,38	0,531	88/120
Principal Agressor com Registros de Crianças de 12 ou menos anos (6)	Sem Amostragem (1.205)	Método Wrapper	60,35	0,663	343/504	45,37	0,515	18/24	60,90	0,676	647/895	60,90	0,676	647/895
		Revisão Bibliográfica	51,60	0,567	304/472	42,51	0,481	17/22	50,83	0,573	474/733	50,83	0,573	474/733
	Com Amostragem (1.188)	Sem Seleção	38,01	0,44	356/470	33,36	0,365	9/10	36,68	0,429	816/1023	36,68	0,429	816/1023
		Método de Filtro	37,59	0,425	54/75	34,11	0,359	9/10	38,67	0,437	77/111	38,67	0,437	77/111
		Método Wrapper	41,83	0,483	118/181	34,11	0,359	9/10	40,75	0,493	181/284	40,75	0,493	181/284
		Revisão Bibliográfica	33,78	0,408	219/359	34,19	0,364	15/18	32,12	0,4	384/629	32,12	0,4	384/629
Apresenta Violência com Registros de Crianças de 12 anos (2)	Sem Amostragem (6.442)	Sem Seleção	45,45	0,501	406/536	33,84	0,389	17/20	45,37	0,501	757/958	45,37	0,501	757/958
		Método de Filtro	35,61	0,418	47/61	33,00	0,353	10/12	33,75	0,422	63/87	33,75	0,422	63/87
	Com Amostragem (6.212)	Método Wrapper	47,31	0,516	422/542	33,84	0,389	17/20	46,55	0,511	787/978	46,55	0,511	787/978
		Revisão Bibliográfica	39,65	0,446	215/371	29,55	0,354	18/23	38,81	0,438	316/545	38,81	0,438	316/545
		Sem Seleção	86,52	0,862	258/348	84,64	0,856	30/40	85,22	0,85	2070/2484	85,22	0,85	2070/2484
		Método de Filtro	81,46	0,786	4/7	81,46	0,786	4/7	81,29	0,791	6/11	81,29	0,791	6/11
Tipo de Violência com Registros de Crianças de 12 anos (4)	Sem Amostragem (4.999)	Método Wrapper	86,76	0,867	126/171	84,61	0,855	26/33	86,20	0,894	847/1104	86,20	0,894	847/1104
		Revisão Bibliográfica	82,38	0,826	154/231	82,09	0,831	15/25	80,79	0,82	1305/1799	80,79	0,82	1305/1799
	Com Amostragem (4.008)	Sem Seleção	86,91	0,892	294/393	81,07	0,879	28/38	85,93	0,883	2074/2480	85,93	0,883	2074/2480
		Método de Filtro	79,41	0,809	10/19	78,06	0,805	6/11	79,44	0,811	16/31	79,44	0,811	16/31
		Método Wrapper	86,93	0,894	263/350	81,23	0,881	35/46	86,27	0,882	1875/2232	86,27	0,882	1875/2232
		Revisão Bibliográfica	81,73	0,83	217/329	79,04	0,826	17/27	81,09	0,851	1277/1811	81,09	0,851	1277/1811
Principal Agressor com Registros de Crianças de 12 anos (7)	Sem Amostragem (4.640)	Sem Seleção	74,90	0,63	205/267	74,54	0,612	13/16	69,61	0,66	2437/2976	69,61	0,66	2437/2976
		Método de Filtro	73,42	0,601	7/13	72,94	0,602	4/7	73,48	0,632	20/39	73,48	0,632	20/39
	Com Amostragem (4.410)	Método Wrapper	74,80	0,607	38/55	74,22	0,612	13/16	71,46	0,639	841/1134	71,46	0,639	841/1134
		Revisão Bibliográfica	73,20	0,617	100/153	73,58	0,602	4/7	67,55	0,637	1395/2014	67,55	0,637	1395/2014
		Sem Seleção	63,22	0,688	885/1164	51,55	0,608	59/68	63,30	0,693	2350/2879	63,30	0,693	2350/2879
		Método de Filtro	43,76	0,521	51/101	42,71	0,472	9/17	43,96	0,526	72/143	43,96	0,526	72/143

Fonte: Elaborado pelo autor.

Figura 8 – Árvore de decisão obtida ao aplicar o algoritmo J48 na predição de presença de violência sobre o conjunto de dados de crianças maiores de 12 anos.



Fonte: Elaborado pelo autor.

pode-se contrastar os resultados conseguidos para cada um desses elementos. Assim, para escolher os melhores considerou-se a MH das curvas mais elevadas para cada grupo de execuções, as que corresponde a:

- Tipo de violência em crianças de 12 e menos anos: com amostragem, seleção de atributos por método *wrapper* e sem limitar o tamanho da árvore;
- Principal agressor em crianças de 12 e menos anos: com amostragem, seleção de atributos por método *wrapper* e valores padrão na execução do algoritmo;
- Presença de violência em crianças de mais de 12 anos: com amostragem, seleção de atributos por método *wrapper* e valores padrão na execução do algoritmo;
- Tipo de violência em crianças de mais de 12 anos: com amostragem, seleção de atributos por método *wrapper* e sem limitar o tamanho da árvore; e

- Principal agressor em crianças de mais de 12 anos: com amostragem, seleção de atributos por método *wrapper* e sem limitar o tamanho da árvore.

Baseado nas execuções anteriores, e de forma geral, pode-se estabelecer que em todos os casos o método de seleção de atributos com melhores resultados foi o *wrapper*. Relacionado ao processo de amostragem, esse não foi o mais eficaz apenas na tarefa de predizer se as crianças maiores a 12 anos apresentam violência ou não, execução que teve o melhor resultado de todos os testes, chegando quase a um 90% de precisão na taxa de acertos e 0,9 de MH.

Acerca dos diferentes atributos preditores, foram obtidos resultados variados que estão determinados principalmente pela quantidade de valores possíveis que têm cada um. Por exemplo, quando se classificou crianças maiores de 12 anos apresentarem ou não vulneração (apenas duas possibilidades), atingiu-se uma taxa de acertos de 86,2% e com MH de aproximadamente 0,9, o que corresponde a um ótimo resultado. De forma contrária, quando buscou-se classificar o principal agressor no conjunto de crianças de 12 ou menos anos (seis valores possíveis), a melhor execução apenas atingiu 47,3% de acertos com MH de 0,52, resultando em um modelo pouco confiável.

Em relação às configurações utilizadas, a maioria das vezes (68% das situações) os melhores resultados foram obtidos sem limitar o tamanho da árvore. Esse fenômeno acontece principalmente pelo *overfitting* do modelo aos dados, o que não é recomendável principalmente porque não é um padrão que pode ser aplicado em dados que não foram utilizados no treinamento. Ademais, quando a quantidade de nós da árvore é numerosa, tem-se o problema de que o modelo dificilmente pode ser lido por um humano caso a quantidade de regras seja muito elevada, motivo pelo qual seria necessária uma ferramenta computacional para realizar essa tarefa. Assim, ao limitar o tamanho da árvore através do mínimo de elementos nos nós folha e também pela poda da árvore, encontraram-se ainda modelos menores e com desempenhos aceitáveis. Esse comportamento pôde ser observado no caso de predizer a presença de violência sem amostragem e usando o método de filtro, o qual atingiu mais de um 81% de acertos (e MH de 0,786) com uma árvore de apenas 7 nós no total (o melhor resultado para essa categoria tem 1.104 nós). Portanto, não é possível estabelecer que sempre um modelo com muitas (ou poucas) regras é o melhor, sempre dependerá do contexto do problema.

Considerando os diferentes tipos de seleção de atributos realizadas, em 76,7% das execuções o método *wrapper* obteve os melhores resultados. Essa técnica empregou o mesmo algoritmo J48 para encontrar os subconjuntos, sempre começando com o conjunto vazio de atributos, e para cada uma das variantes das coleções de dados entregou diferentes variáveis para gerar as árvores. Portanto, não é possível definir apenas um grupo de elementos fixo para o problema da violência infantil. Sobre as outras formas de seleção

atributos: o método de filtro nunca foi o melhor, contudo é uma boa técnica quando se quer listar os elementos que se tem em base na sua relação com o valor preditor; e com a informação obtida na revisão bibliográfica, também não conseguiu-se nenhum melhor resultado, o que indica que o contexto define quais são os variáveis mais relevantes e não o problema como tal. Assim, foi possível estabelecer que os métodos de seleção de atributos *wrapper* são uma alternativa factível para a redução da dimensionalidade de um conjunto de dados nessa problemática.

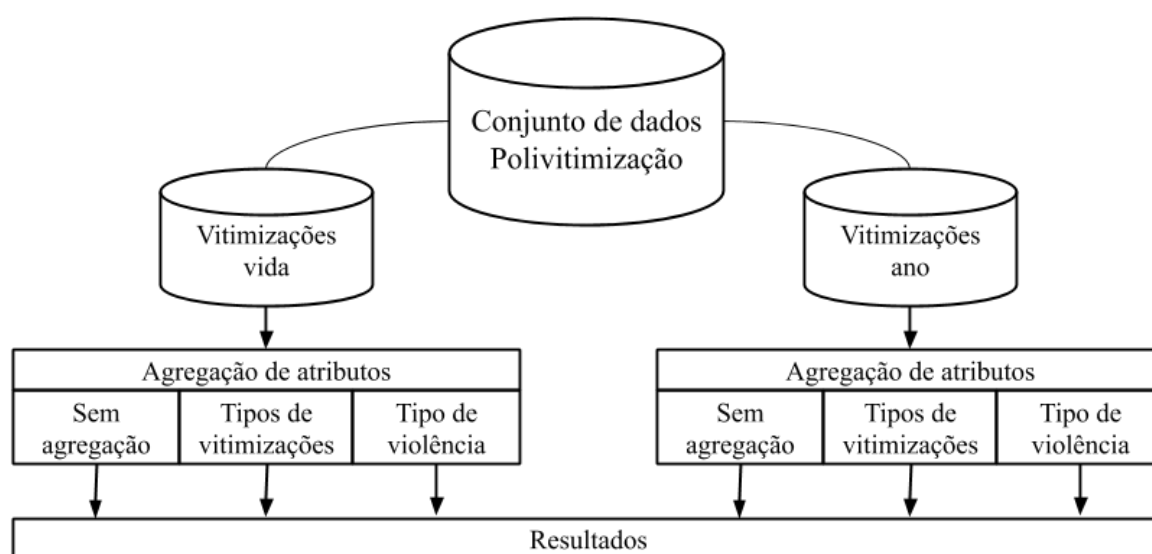
A amostragem de dados para balancear os conjuntos foi realizada para todos os atributos preditores. Na maioria dos casos conseguiram-se modelos com melhores desempenhos que quando não foi feito o balanceamento. Ademais, trabalhar com coleções balanceadas permite ter uma maior confiança, visto que o classificador não terá um viés sobre algum valor da variável preditora. Por exemplo, quando utilizou-se o atributo preditor principal agressor para crianças maiores de 12 anos, os modelos gerados com os conjuntos que não tiveram amostragem nunca superaram os 50% de acertos, porém obtiveram AUC ROC de no mínimo 0,55, efeito explicado com a presença de um dos valores predominantes na variável preditora (nesse caso os valores *mãe* e *pai*).

Considerando novamente os melhores resultados para cada um dos 5 diferente elementos preditos, estabelece-se que o modelo para determinar se uma criança maior a 12 anos está sofrendo de violência ou não é o único com um desempenho ótimo. Nesse caso, a MH varia aproximadamente entre 0,80 e 0,90, concluindo-se com esses valores que: o modelo faz uma discriminação correta entre as classificações possíveis (AUC ROC elevado), possui uma precisão e sensibilidade elevadas, e não possui *overfitting* ao momento de realizar a predição. Sobre os outros resultados, quando tentou-se determinar o tipo de violência (para as duas segmentações) e o principal agressor em crianças de 12 ou mais anos foram semelhantes. Para esses casos obtive-se a AUC ROC sempre maior a 0,8 (sendo um resultado desejado) mas a AUC PR com um valor baixo (aproximadamente de 0,6), determinando que o modelo não é um modelo adequado. Ademais, as taxas de classificações corretas são próximas 60%, o que não faz ao modelo confiável na tarefa de predizer esses 3 casos. Por último, para o resultado na predição do principal agressor em crianças de 12 ou menos anos, todas as métricas são deficientes, portanto não é um modelo confiável em nenhum caso. Isso pode ser explicado principalmente pela grande quantidade de valores que esse atributo pode tomar, sendo 6 no conjunto de crianças de 12 ou menos anos e 7 para o de 13 ou mais, dificultando a aprendizagem dos modelos ao ter tantas saídas possíveis.

5.2 Questionário sobre Polivitimização

Para o segundo conjunto de dados foram feitas variadas agregações dos tipos de vitimizações, tentando encontrar regras que tenham ótimos desempenhos e sejam de utilidade na abordagem do problema da violência infantil. Assim, a [Figura 9](#) descreve o fluxo seguido para os dois sub-conjuntos de dados (vitimizações vida e ano).

Figura 9 – Considerações e trabalhos feitos previa aplicação do algoritmo Apriori sobre os conjuntos de dados não supervisionados.



Fonte: Elaborado pelo autor.

Como foi descrito na [subseção 4.2.2](#) e apresenta-se na [Figura 9](#), fizeram-se duas agregações diferentes de dados, uma para os tipos (agrupações) de vitimizações definidas na [subseção 4.1.2](#) e outra para as quatro principais categorias de violência (psicológica, física, abuso sexual, e negligência), além de conservar os atributos sem modificações. Considerando a natureza do algoritmo, nesse caso não foram aplicadas diferentes configurações (valores fixos); utilizou-se um valor baixo para o suporte tentando obter uma quantidade significativa de regras e uma confiança elevada para filtrar as que têm uma presença do consequente sobre o antecedente não considerável, sendo esses valores 0,05 e 0,75 respectivamente. Assim, foram feitas 6 execuções diferentes do algoritmo Apriori, das quais o número de regras obtido em cada uma apresenta-se na [Tabela 6](#). A diferença nas quantidades de regras explicam-se principalmente pelo número de atributos que tem o conjunto e se é considerável a presença de valores positivos (sofreu algum tipo de violência ou vitimização) sobre a amostra total, o que aumenta as combinações de elementos para gerar novas regras. Em base dessa premissa, é esperado que o conjunto de vitimizações vida sem agregação tenha uma grande quantidade de regras entendendo que são 32 atributos e a história de vida toda da criança onde poderia ter sofrido alguma vitimização; ou que na vitimização ano quando a agregação é por tipo de violência tenha apenas 14 regras,

considerando que são 4 atributos diferentes e somente 1 ano de historia para avaliar a presença de violência.

Tabela 6 – Quantidade total de regras conseguidas nas diferentes execuções do algoritmo Apriori sobre os conjuntos de dados não supervisionados, utilizando um suporte de 0,05 e confiança de 0,75.

Conjunto de Dados	Agregação dos Atributos		
	Sem Agregação	Tipo de Vitimização	Tipo de Violência
Vitimizações Vida	6963	182	29
Vitimizações Ano	21	56	14

Fonte: Elaborado pelo autor.

Para avaliar os resultados, observaram-se os valores *lift*, *conviction*, e *leverage*, isso para cada uma das regras geradas. Pelo fato de que não existem valores padrão para essas métricas na avaliação, e o grande número de regras obtidas em alguns casos, é que apenas consideraram-se apresentar as que têm um valor maior à média em cada métrica para todas as execuções, resultando as quantidades apresentadas na [Tabela 7](#). Comparando com a [Tabela 6](#), é importante a diminuição de regras por conjunto, o que é esperado depois de ter aplicado esse filtro. Contudo, é relevante destacar que ainda cada execução tem pelo menos uma regra válida, as que têm ótimos valores para as métricas de avaliação.

Tabela 7 – Quantidade total de regras conseguidas com o algoritmo Apriori onde os seus valores das métricas *lift*, *conviction*, e *leverage* são maiores às medias de cada uma das execuções.

Conjunto de Dados	Agregação dos Atributos		
	Sem Agregação	Tipo de Vitimização	Tipo de Violência
Vitimizações Vida	704	3	1
Vitimizações Ano	7	3	1

Fonte: Elaborado pelo autor.

As regras conseguidas após aplicar esse filtro estão descritas na [Tabela 8](#), detalhando a qual conjunto de dados corresponde e os respectivos valores obtidos para as métricas de avaliação. Para o conjunto de vitimizações vida sem agregação, apenas incluíram-se as primeiras 20 regras ordenadas de forma descendente por seus valores de confiança; um número maior da regras de todas as execuções apresentam-se na [seção A.2](#) do [Apêndice A](#).

Antes de discutir sobre as regras obtidas, é importante avaliá-las através das diferentes métricas. Como já foi estabelecido, os valores utilizados como mínimo para o suporte e a confiança foram de 0,05 e 0,75 respectivamente, o que não limita que poderiam ser mais altos. Sobre o suporte, considera-se que a maioria das regras apresentadas na [Tabela 8](#) são baixos no geral; aumenta quando se tem algum tipo de agregação mas ainda continua não sendo um número elevado (a exceção do tipo de violência para a vitimiza-

Tabela 8 – Regras de associação geradas com o algoritmo Apriori para cada um dos diferentes conjuntos de dados diferenciados por vitimizações em vida e ano e qual é agregação dos atributos que apresenta detalhando o suporte e a confiança, ademais das métricas utilizadas para avaliá-las.

Conjunto de Dados	Tipo de Agregação	Regra de Associação		Valores das Métricas				
		Antecedente	Consequente	Confiança	Lift	Leverage	Conviction	
Vitimizações Vida	Sem Agregação	Indiretas4, Pares1, CC3, Cuidadores1, Cuidadores2	CC6	0,96	1,81	0,03	10,74	
		Indiretas4, Pares1, Cuidadores1, CC4, Cuidadores2	CC6	0,96	1,81	0,03	10,55	
		Indiretas2, Indiretas3, Indiretas4, Pares1, Cuidadores1, Cuidadores2	CC6	0,95	1,81	0,03	10,20	
		Indiretas2, Indiretas4, Pares1, Cuidadores1, Cuidadores2	CC6	0,95	1,80	0,03	9,97	
		Indiretas2, Indiretas4, Pares1, CC4, Cuidadores2	CC6	0,95	1,80	0,03	9,51	
		Indiretas2, Pares1, CC3, Cuidadores1, Cuidadores2	CC6	0,95	1,80	0,03	9,51	
		Indiretas2, Indiretas3, Indiretas4, Pares1, CC3, Cuidadores2	CC6	0,95	1,80	0,03	9,44	
		Indiretas2, Indiretas4, Pares1, CC3, Cuidadores2	CC6	0,95	1,80	0,03	9,42	
		Indiretas3, Indiretas4, Pares1, Cuidadores1, Cuidadores2	CC6	0,95	1,80	0,03	9,34	
		Indiretas3, Indiretas4, Pares1, CC3, Cuidadores2	CC6	0,95	1,80	0,03	9,32	
		Indiretas4, Pares1, CC4, Cuidadores2	CC6	0,95	1,80	0,03	9,18	
		Indiretas3, Indiretas4, Pares1, CC4, Cuidadores2	CC6	0,95	1,80	0,03	9,14	
		Indiretas4, Pares1, Cuidadores1, Cuidadores2	CC6	0,95	1,79	0,04	9,01	
		Indiretas4, Pares1, CC3, Cuidadores2	CC6	0,95	1,79	0,04	8,99	
		Indiretas4, Pares1, CCI, Cuidadores2	CC6	0,95	1,79	0,03	8,82	
		Indiretas2, Pares1, Cuidadores1, CC4, Cuidadores2	CC6	0,95	1,79	0,03	8,67	
Indiretas2, Indiretas4, Pares1, CCI, Cuidadores2	CC6	0,95	1,79	0,03	8,67			
Indiretas3, Indiretas4, Pares1, CCI, Cuidadores2	CC6	0,95	1,79	0,03	8,67			
Indiretas2, Indiretas3, Pares1, CC4, Cuidadores2	CC6	0,95	1,79	0,03	8,41			
Indiretas4, Pares1, CC3, Cuidadores1, CC4	CC6	0,94	1,79	0,03	8,38			
Indiretas, CC, Cuidadores, Digitais	Pares	0,90	1,48	0,06	4,04			
CC, Digitais, Sexuais	Indiretas, Pares	0,89	1,59	0,05	4,11			
Cuidadores, Digitais, Sexuais	Indiretas, CC, Pares	0,89	1,70	0,05	4,31			
VF, AS	Negligencia, VP	0,93	1,31	0,05	4,32			
Vitimização Ano	Sem Agregação	Indiretas2, Indiretas3, Pares1	CC6	0,85	2,99	0,03	4,69	
		Indiretas4, Pares1	CC6	0,84	2,97	0,04	4,54	
		Indiretas3, Pares1	CC6	0,82	2,89	0,05	4,00	
		Indiretas2, Pares1	CC6	0,82	2,87	0,05	3,89	
		Indiretas3, Indiretas5	Indiretas4	0,81	2,63	0,03	3,66	
		Cuidadores1, Pares1	CC6	0,81	2,86	0,03	3,78	
		Pares1	CC6	0,78	2,75	0,07	3,26	
		Indiretas, Pares, Digitais	CC	0,90	1,77	0,04	4,91	
		Indiretas, Cuidadores, Pares	CC	0,89	1,75	0,05	4,42	
		Pares, Digitais	CC	0,89	1,75	0,04	4,40	
VF, AS	Negligencia, VP	0,84	1,80	0,05	3,33			

Fonte: Elaborado pelo autor.

ção vida, regra que tem um suporte de 3.675 registros). Porém, as confianças das regras possuem valores adequados, estabelecendo que têm uma presença relevante na amostra. Para os registros de vitimizações vida a maioria têm um 0,93, o que implica que mais de 90% dos registros satisfazem a regra para um conjunto determinado pelo suporte. Apesar de que para as vitimizações ano esses valores diminuem, ainda são consideráveis; a maior parte cumpre com ter uma porcentagem de presença de mais de 80% sobre os registros totais conformados pelo correspondente suporte.

Sobre as três métricas de avaliação, pode-se destacar inicialmente que todas as regras satisfazem os valores mínimos para considerar que existe uma dependência entre o antecedente e o consequente; *lift* e *conviction* sempre são maiores a 1, e *leverage* sempre maior a 0. *lift* e *leverage* são métricas semelhantes, ambas comparam a presença do antecedente e o consequentes quando aparecem juntos e separados; o primeiro a través de uma taxa e o segundo com uma diferença (fato pelo qual os valores são próximos ao valor 0). Isso reafirma a existência da relação de dependência entre antecedente e consequente. *Conviction* permite associar o antecedente com a ausência do consequente com a respectiva direcionalidade que a regra possui ($X \Rightarrow Y$), o que não fazem as outras métricas; para todas as regras descritas na [Tabela 8](#) os valores novamente descrevem uma relação de dependência.

A respeito das regras encontradas, é simples de perceber que são bastante semelhantes umas às outras, o que explica-se principalmente pela predominância de alguns atributos sobre o resto. Por exemplo, para os dois conjuntos de vitimizações, a maioria das regras nos conjuntos sem agregação possuem como consequente o elemento CC6 (ataques físicos sem objetos), o que poderia ser entendido como uma vitimização muito mais frequente que o resto. Nos antecedentes das mesmas sim se tem mais vitimizações, mas ainda continua sendo um conjunto pequeno de elementos considerando os 32 atributos possíveis. Além disso, não tem-se nenhum atributo das categorias de vitimizações digitais e sexuais nas duas partes que compõem as regras. Sobre os tipos de vitimização, pode-se dizer que para as vitimizações vida as três regras sim são diversas nas suas formas, não assim com as geradas para as vitimizações ano, onde o consequente sempre é o mesmo. Para a agregação por tipo de violência, ambas regras são iguais (o que valida a sua qualidade), o que novamente pode ser explicado pela quantidade de atributos que tem-se para os dois conjuntos; com apenas 4 variáveis agregadas é esperado que a taxa de valores positivos aumente e as regras possíveis totais diminuam em número, sendo mais provável que alguma possa-se repetir. Esse último fato também valida as semelhanças dos dois conjuntos de dados (vida e ano) além da temporalidade que têm.

Em relação as mesmas regras apresentadas na [Tabela 8](#) e avaliando a utilidade que possam ter, talvez não sejam resultados que entreguem informação valiosa para um profissional que trabalha com crianças violentadas ou em contextos semelhante, conside-

rando a complexidade que têm na sua forma e as poucas variações de umas com outras. Tendo como exemplo disso, poderia ser excessivamente difícil e pouco relevante detetar todas as regras que não tiveram agregações no conjunto de vitimizações ano para apenas concluir que estão relacionadas com um ataque físico sem objeto (CC6). O anterior obriga a avaliar o tipo de resultado obtido com profissionais da área, além de definir se existem regras que a priori poderiam ser mais relevantes que outras, assim ter como objetivo avaliar esse tipo de sentenças que podem ser mais uteis na hora de estudar algum caso de violência infantil.

6 Conclusões

Esse último capítulo inclui as conclusões de todo o trabalho realizado, mencionando os aspectos descobertos que são relevantes ao estudo, além de quais são as atividades que devem ser realizadas em trabalhos futuros.

6.1 Trabalho Realizado

O principal objetivo desse estudo foi contribuir em um acercamento para dar solução a uma problemática social através da ciência da computação, especificamente detectar padrões que descrevem a violência infantil a partir de técnicas de ML. Isso poderia ajudar consideravelmente os trabalhadores que vinculam-se com as crianças que experimentam essas situações de sofrimento, permitindo a eles tomar melhores decisões na abordagem desses casos.

Como foi apresentado no [Capítulo 1](#), a violência infantil é uma problemática que está no mundo todo e que não discrimina por cultura, etnia, situação econômica ou qualquer característica social. São muitas as pesquisas desenvolvidas que contribuem ao entendimento ou possível solução desse problema, mas até agora não se têm diretrizes padrão para abordá-lo, o que implica que não existe consenso em determinar quais são os principais fatores que acrescentam (ou diminuem) a presença de violência infantil. Por outro lado, ML representa um conjunto de ferramentas computacionais que têm variados usos, sendo esses principalmente a descrição e predição (ou classificação) de algum fenômeno através de grandes quantidades de dados. Isso pode ser aplicado em qualquer área do conhecimento, e como foi exposto neste documento, também já fora usado para prever a violência infantil anteriormente.

No [Capítulo 3](#) foram revisados diferentes artigos baseados em técnicas estatísticas ou que fazem uso de ML para abordar o problema da predição da violência infantil, possuindo a maioria características importantes a destacar. Variados algoritmos foram utilizados obtendo diferentes valores de precisão ao testá-los. Conseqüentemente, não é possível estabelecer apenas uma técnica a utilizar para esse problema. Ademais, todos os conjuntos de dados são distintos, implicando que cada uma das pesquisas estudadas foi aplicada em um contexto social e cultural diferente. Portanto, a tarefa de padronização do problema da violência infantil é muito mais difícil de se concretizar.

A metodologia utilizada nesse trabalho possui etapas que comumente são empregadas em projetos de ML, sendo essas: entendimento do problema e obtenção dos dados, pré-processamento dos dados, treinamento e validação dos modelos e análises dos

resultados. Nesse caso, a que apresentou a maior dificuldade foi a obtenção dos dados supervisionados, isso por três motivos principais: (1) não existe um número grande de organizações que trabalhem com crianças no contexto de violência; (2) essas poucas instituições normalmente não geram nem armazenam constantemente registros padronizados dos casos que assistem; e (3) o fato de conceder dados de crianças para algum ente externo à organização às vezes pode ser um impedimento legal para realizar uma pesquisa desse tipo e, por conseguinte, nem sempre têm a disposição de entregar os dados coletados. Contudo, é importante destacar que existem iniciativas que permitem a qualquer pessoa ter um fácil acesso a conjuntos de dados semelhantes, como a que fez a entidade de prevenção do crime do governo do Chile ao aplicar uma enquete sobre polivitimização em crianças desse país, aportando uma coleção de registros não supervisionada.

Sobre a análise dos dados supervisionados, o algoritmo C4.5 aplicado nesse trabalho obteve variados resultados, os quais são refletidos através dos atributos a serem preditos. Quando utilizou-se a variável *se a criança sofreu violência ou não* (unicamente as maiores a 12 anos), todas as configurações empregadas conseguiram elevados desempenhos, chegando a predizer corretamente mais do 86% dos casos. Para os modelos da classificação do tipo de violência, para as crianças de 13 ou mais anos os resultados foram acetáveis, chegando a obter uma taxa de acertos próxima de 75%, enquanto para as crianças menores de 13 anos, os desempenhos dos modelos diminuíram, chegando a ter no máximo 60% de acertos, o que não corresponde a um modelo confiável. Sobre a predição do principal agressor, também conseguiram-se resultados pouco confiáveis, entre 64% e 47% de acertos para crianças de 13 ou mais anos e menores a 13 anos, respectivamente.

Considerando os dados não supervisionados e o algoritmo Apriori utilizados, para cada uma das execuções encontraram-se regras ou padrões que, tendo em conta os critérios sobre as métricas de avaliação usadas, são válidas computacionalmente. Porém, as execuções do algoritmo feitas sobre dados agregados (por vitimização ou tipo de violência) não tiveram muita variedade de elementos no antecedente e consequente das regras, isso principalmente pelo número pequeno de atributos que compõem esses conjuntos de dados. Quando não foi feita nenhuma agregação de atributos, o número de regras válidas foi o maior. Ainda assim, os valores que as compõem na maioria das vezes são os mesmos (indiretas 2, 3 e 4, Cuidadores1, Pares1, entre as mais frequentes). É por isso que os resultados encontrados nessa parte da pesquisa, além da validação computacional, precisam ser avaliados por profissionais que conheçam do contexto da violência infantil.

Cabe mencionar que a hipótese estabelecida no início dessa pesquisa foi pautada no seguinte questionamento: modelos baseados em regras gerados por meio de técnicas de ML com dados de instituições chilenas que trabalham com crianças abusadas podem entregar padrões que permitam estimar em que casos algum ato de violência infantil está acontecendo ou irá acontecer? Assim, notou-se que para o conjunto de dados supervisio-

nado, foram descobertos padrões que permitem estimar quando uma criança está sofrendo de violência e o seu tipo. Esse resultado é complementar àqueles obtidos com os registros não supervisionados, que permitem relacionar diferentes tipos (ou agrupamentos) de violência infantil.

Sobre os objetivos do trabalho, tanto o principal como os secundários foram atingidos. Por meio de métricas de validação conseguiram-se avaliar os modelos preditivos e regras criadas a partir dos dados de forma positiva, concluindo que podem ser utilizados em análises de casos de violência e detectaram os principais atributos que estão relacionados ao problema da violência infantil no contexto do Chile.

Considera-se importante realizar pesquisas desse tipo, onde são abordados problemas que realmente são significativos à sociedade a partir de uma perspectiva não convencional; nesse caso, por meio do uso de técnicas de ML.

6.2 Trabalhos Futuros

Após concluir a pesquisa ainda existem diferentes tarefas que não foram realizadas, pois não estavam definidas no escopo do projeto. Neste sentido, destacam-se dois trabalhos de grande importância. O primeiro está relacionado a avaliar os resultados obtidos não simplesmente a partir de uma perspectiva computacional, mas por uma visão dos indivíduos que trabalham diretamente com crianças que sofrem violência. Além de uma métrica para avaliar o desempenho de um modelo, a experiência dessas pessoas pode determinar se uma regra é válida ou não, ou se realmente os resultados entregam informação relevante no entendimento do problema. A segunda tarefa é o desenvolvimento de alguma ferramenta que permita utilizar os modelos gerados, com uma interface simples que possa ser usada facilmente por qualquer pessoa que trabalhe no contexto de violência infantil e que necessite de apoio para estudar esse tipo de casos. Isso propiciaria fazer uso de técnicas de ML que não necessariamente entreguem regras como resultado, possibilitando obter modelos com melhores desempenhos na tarefa de prever a violência infantil.

APÊNDICE A – Resultados

Nessa seção do apêndice são apresentados os modelos obtidos ao gerar diferentes execuções dos algoritmos J48 e Apriori da ferramenta WEKA 3.8.

A.1 J48

A cerca das árvores conseguidas com os dados supervisionados, a seguir expõem-se as descritas na [seção 5.1](#) como as de menor tamanho para cada um dos 5 atributos preditores, utilizando a configuração padrão da ferramenta. Os valores finais que aparecem em cada um dos nós folha correspondem aos registros classificados corretamente (esquerda) e os errôneos (direita). Ademais, a barra vertical (|) define o nível da árvore ao que corresponde o atributo de cada linha.

- Modelo 1 – Tipo de violência em crianças de 12 ou menos anos, com configuração: amostragem sobre os dados, e seleção de atributos por método de filtro. Valores das métricas de validação: porcentagem de acertos = 43,67%, AUC ROC = 0,696, AUC PR = 0,418, e MH = 0,522; com 52 nós no total.

```

SEXO = MASCULINO
| ZONA = NORTE
| | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SIST. ESC. FORMAL: VIOLÊNCIA FÍSICA (544.55/339.63)
| | SITUAÇÃO ESCOLAR ATUAL = SEM ATIVIDADE ESCOLAR: NEGLIGÊNCIA (19.38/4.17)
| | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SISTEMA ESCOLAR PARA ADULTOS: VIOLÊNCIA FÍSICA (0.0)
| | SITUAÇÃO ESCOLAR ATUAL = VALIDAÇÃO DE ESTUDOS: VIOLÊNCIA FÍSICA (0.0)
| | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SIST. ESC. FECHADO: NEGLIGÊNCIA (4.08/0.04)
| | SITUAÇÃO ESCOLAR ATUAL = ESCOLARIZADA: VIOLÊNCIA FÍSICA (0.0)
| ZONA = SUL
| | REPETÊNCIA = NÃO: VIOLÊNCIA PSICOLÓGICA (260.82/147.77)
| | REPETÊNCIA = SIM: NEGLIGÊNCIA (41.18/11.09)
| ZONA = CENTRO
| | ADULTO RESPONSÁVEL = MÃE
| | | REPETÊNCIA = NÃO: VIOLÊNCIA FÍSICA (32.0/18.0)
| | | REPETÊNCIA = SIM: NEGLIGÊNCIA (9.0/2.0)
| | ADULTO RESPONSÁVEL = AVÓ: NEGLIGÊNCIA (5.0/2.0)
| | ADULTO RESPONSÁVEL = MÃE E PAI: VIOLÊNCIA FÍSICA (2.0)
| | ADULTO RESPONSÁVEL = PAI: VIOLÊNCIA FÍSICA (14.0/1.0)
| | ADULTO RESPONSÁVEL = OUTRO FAMILIAR: NEGLIGÊNCIA (3.0/1.0)
| | ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: VIOLÊNCIA FÍSICA (0.0)
| | ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: VIOLÊNCIA FÍSICA (0.0)
SEXO = FEMININO
| ADULTO RESPONSÁVEL = MÃE
| | ZONA = NORTE: ABUSO SEXUAL (367.41/180.41)

```

```

| | ZONA = SUL
| | | REPETÊNCIA = NÃO: ABUSO SEXUAL (176.82/115.59)
| | | REPETÊNCIA = SIM: NEGLIGÊNCIA (10.87/3.58)
| | ZONA = CENTRO
| | | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SISTEMA ESCOLAR FORMAL: ABUSO SEXUAL (33.34/18.67)
| | | SITUAÇÃO ESCOLAR ATUAL = SEM ATIVIDADE ESCOLAR: NEGLIGÊNCIA (2.0)
| | | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SISTEMA ESCOLAR PARA ADULTOS: ABUSO SEXUAL (0.0)
| | | SITUAÇÃO ESCOLAR ATUAL = VALIDAÇÃO DE ESTUDOS: NEGLIGÊNCIA (1.67)
| | | SITUAÇÃO ESCOLAR ATUAL = CURSANDO SIST. ESCOLAR FECHADO: ABUSO SEXUAL (0.0)
| | | SITUAÇÃO ESCOLAR ATUAL = ESCOLARIZADA: ABUSO SEXUAL (0.0)
| ADULTO RESPONSÁVEL = AVÓ
| | ZONA = NORTE: NEGLIGÊNCIA (79.87/37.94)
| | ZONA = SUL: VIOLÊNCIA PSICOLÓGICA (36.53/17.53)
| | ZONA = CENTRO: VIOLÊNCIA FÍSICA (1.4/0.4)
| ADULTO RESPONSÁVEL = MÃE E PAI
| | ZONA = NORTE: ABUSO SEXUAL (24.88/13.88)
| | ZONA = SUL: NEGLIGÊNCIA (28.25/12.25)
| | ZONA = CENTRO: ABUSO SEXUAL (2.19/1.13)
| ADULTO RESPONSÁVEL = PAI
| | ZONA = NORTE
| | | REPETÊNCIA = NÃO: VIOLÊNCIA FÍSICA (33.72/12.45)
| | | REPETÊNCIA = SIM: ABUSO SEXUAL (2.06/1.06)
| | ZONA = SUL: NEGLIGÊNCIA (10.22/5.22)
| | ZONA = CENTRO: VIOLÊNCIA FÍSICA (3.17/0.17)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR: NEGLIGÊNCIA (39.95/15.59)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: NEGLIGÊNCIA (25.61/13.38)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: ABUSO SEXUAL (1.02/0.02)

```

- Modelo 2 – Principal agressor em crianças de 12 ou menos anos, com configuração: amostragem sobre os dados, e seleção de atributos por método de filtro. Valores das métricas de validação: porcentagem de acertos = 35,61%, AUC ROC = 0,669, AUC PR = 0,304, e MH = 0,418; com 61 nós no total.

```

REGIÃO = ATACAMA (III): PAI (2.0/1.0)
REGIÃO = ANTOFAGASTA (II)
| ADULTO RESPONSÁVEL = PAI: MÃE (43.76/23.21)
| ADULTO RESPONSÁVEL = MÃE
| | SEXO = MASCULINO: PAI (168.44/130.44)
| | SEXO = FEMININO: PADRASTO (284.46/190.74)
| ADULTO RESPONSÁVEL = MÃE E PAI: MÃE E PAI (26.46/16.46)
| ADULTO RESPONSÁVEL = AVÓ: MÃE (75.31/53.36)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: PADRASTO (14.25/10.23)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR
| | SEXO = MASCULINO: MÃE E PAI (12.06/6.06)
| | SEXO = FEMININO: OUTRO FAMILIAR (6.26/3.2)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: PADRASTO (0.0)
REGIÃO = LOS RIOS (XIV)
| ADULTO RESPONSÁVEL = PAI: OUTRO NÃO FAMILIAR (8.1/5.0)
| ADULTO RESPONSÁVEL = MÃE: OUTRO NÃO FAMILIAR (118.39/48.0)
| ADULTO RESPONSÁVEL = MÃE E PAI
| | ZONA NA QUAL MORA = RURAL: MÃE E PAI (3.0)
| | ZONA NA QUAL MORA = URBANA: OUTRO NÃO FAMILIAR (10.15/6.0)
| ADULTO RESPONSÁVEL = AVÓ
| | ZONA NA QUAL MORA = RURAL: MÃE E PAI (11.0/4.0)

```

```

| | ZONA NA QUAL MORA = URBANA
| | | SEXO = MASCULINO: MÃE (5.0/3.0)
| | | SEXO = FEMININO: OUTRO NÃO FAMILIAR (11.32/4.0)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: OUTRO FAMILIAR (3.04/1.04)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR: OUTRO NÃO FAMILIAR (0.0)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: OUTRO NÃO FAMILIAR (0.0)
REGIÃO = ARAUCANÍA (IX)
| ADULTO RESPONSÁVEL = PAI: MÃE (17.0/8.0)
| ADULTO RESPONSÁVEL = MÃE: MÃE E PAI (113.0/69.0)
| ADULTO RESPONSÁVEL = MÃE E PAI: MÃE E PAI (32.0/6.0)
| ADULTO RESPONSÁVEL = AVÓ
| | ZONA NA QUAL MORA = RURAL: OUTRO FAMILIAR (7.0/2.0)
| | ZONA NA QUAL MORA = URBANA: MÃE (10.0/5.0)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: MÃE E PAI (0.0)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR: MÃE (1.0)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: MÃE E PAI (0.0)
REGIÃO = METROPOLITANA (XIII)
| ADULTO RESPONSÁVEL = PAI: MÃE (3.04/1.04)
| ADULTO RESPONSÁVEL = MÃE: OUTRO FAMILIAR (62.81/43.81)
| ADULTO RESPONSÁVEL = MÃE E PAI: MÃE E PAI (3.04/1.04)
| ADULTO RESPONSÁVEL = AVÓ: MÃE E PAI (5.06/3.06)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: OUTRO FAMILIAR (0.0)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR: MÃE E PAI (4.05/2.05)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: OUTRO FAMILIAR (0.0)
REGIÃO = TARAPACÁ (I)
| SEXO = MASCULINO: MÃE (32.0/21.0)
| SEXO = FEMININO: OUTRO FAMILIAR (40.0/28.0)
REGIÃO = BIOBIO (VIII)
| ADULTO RESPONSÁVEL = PAI
| | ZONA NA QUAL MORA = RURAL: MÃE E PAI (2.0)
| | ZONA NA QUAL MORA = URBANA: PAI (4.0/2.0)
| ADULTO RESPONSÁVEL = MÃE: MÃE (32.0/11.0)
| ADULTO RESPONSÁVEL = MÃE E PAI: MÃE E PAI (3.0)
| ADULTO RESPONSÁVEL = AVÓ: MÃE E PAI (8.0/3.0)
| ADULTO RESPONSÁVEL = OUTRO NÃO FAMILIAR: MÃE E PAI (1.0)
| ADULTO RESPONSÁVEL = OUTRO FAMILIAR: MÃE (1.0)
| ADULTO RESPONSÁVEL = NÃO TEM REFERENTE: MÃE (1.0)
REGIÃO = VALPARAÍSO (V): MÃE E PAI (1.0)
REGIÃO = O'HIGGINS (VI): PAI (2.0/1.0)

```

- Modelo 3 – Presença de violência em crianças de mais de 12 anos, com configuração: sem amostragem sobre os dados, e seleção de atributos por método de filtro. Valores das métricas de validação: porcentagem de acertos = 81,46%, AUC ROC = 0,774, AUC PR = 0,798, e MH = 0,786; com 7 nós no total.

```

SOCIALIZAÇÃO EM RUA = NÃO
| TEM COMETIDO DELITOS = NÃO: SIM (900.92/99.97)
| TEM COMETIDO DELITOS = SIM
| | NUMERO DE PROGRAMAS DE PROTEÇÃO NOS QUAIS TEM PARTICIPADO <= 0: NÃO (1397.35/540.72)
| | NUMERO DE PROGRAMAS DE PROTEÇÃO NOS QUAIS TEM PARTICIPADO > 0: SIM (989.46/295.96)
SOCIALIZAÇÃO EM RUA = SIM: SIM (3254.28/300.45)

```

- Modelo 4 – Tipo de violência em crianças de mais de 12 anos, com configuração: sem amostragem sobre os dados, e seleção de atributos por método de filtro. Valores

das métricas de validação: porcentagem de acertos = 73,42%, AUC ROC = 0,604, AUC PR = 0,598, e MH = 0,601; com 13 nós no total.

```
SEXO = FEMININO
| APRESENTA CONSUMO DE SUBSTÂNCIAS = NÃO
| | TEM COMETIDO DELITOS = NÃO
| | | APRESENTOU SAÍDA ESCOLAR = NÃO: ABUSO SEXUAL (345.14/169.67)
| | | APRESENTOU SAÍDA ESCOLAR = SIM
| | | | IDADE <= 13: ABUSO SEXUAL (2.7/1.35)
| | | | IDADE > 13
| | | | | IDADE <= 16: NEGLIGÊNCIA (55.64/21.95)
| | | | | IDADE > 16: ABUSO SEXUAL (15.04/7.29)
| | TEM COMETIDO DELITOS = SIM: NEGLIGÊNCIA (50.84/17.63)
| APRESENTA CONSUMO DE SUBSTÂNCIAS = SIM: NEGLIGÊNCIA (708.62/239.82)
SEXO = MASCULINO: NEGLIGÊNCIA (3811.0/858.0)
```

- Modelo 5 – Principal agressor em crianças de mais de 12 anos, com configuração: sem amostragem sobre os dados, e seleção de atributos por método de filtro. Valores das métricas de validação: porcentagem de acertos = 47,07%, AUC ROC = 0,552, AUC PR = 0,329, e MH = 0,412; com 21 nós no total.

```
APRESENTA CONSUMO DE SUBSTÂNCIAS = NÃO
| TEM COMETIDO DELITOS = NÃO
| | SEXO = FEMININO
| | | APRESENTOU SAÍDA ESCOLAR = NÃO: OUTRO NÃO FAMILIAR (322.87/238.52)
| | | APRESENTOU SAÍDA ESCOLAR = SIM
| | | | IDADE <= 13: OUTRO FAMILIAR (2.71/1.53)
| | | | IDADE > 13
| | | | | IDADE <= 15
| | | | | IDADE <= 14: MÃE (14.47/8.15)
| | | | | IDADE > 14: MÃE E PAI (16.8/8.8)
| | | | | IDADE > 15: MÃE (35.86/23.0)
| | SEXO = MASCULINO
| | | APRESENTOU SAÍDA ESCOLAR = NÃO
| | | | IDADE <= 15: MÃE (87.65/60.33)
| | | | IDADE > 15: PAI (26.6/14.8)
| | | APRESENTOU SAÍDA ESCOLAR = SIM
| | | | IDADE <= 13: OUTRO NÃO FAMILIAR (2.38/0.38)
| | | | IDADE > 13: MÃE E PAI (30.94/12.97)
| TEM COMETIDO DELITOS = SIM: MÃE E PAI (137.0/65.96)
APRESENTA CONSUMO DE SUBSTÂNCIAS = SIM: MÃE E PAI (3962.72/1983.91)
```

A.2 Apriori

As regras obtidas através do algoritmo Apriori para cada uma das 6 execuções (todas com suporte de 0,05 e confiança de 0,75 como valores mínimos) são apresentadas a seguir no mesmo formato que são entregadas pela ferramenta WEKA 3.8, o que está conformado por: elementos que compõem o antecedente entre parêntesis ([]) seguidos da quantidade, símbolos ==>, elementos que conformam o conseqüente entre parêntesis

([]) seguidos da quantidade, confiança (conf), *lift*, *leverage* (lev), e *conviction* (conv). Considerando que alguns conjuntos de dados proporcionaram muitas regras, limitou-se a quantidade a apresentar nesse documento a 25 por execução, ordenadas em forma ascendente sobre os valores da confiança.

- Regras para conjunto de dados das vitimizações em vida sem agregação de atributos (25 regras de um total de 6.963).

1. [Indiretas2, Indiretas4, Pares1, CC3, Cuidadores1, Cuidadores2]: 792 ==> [CC6]: 761
<conf:(0.96)> lift:(1.82) lev:(0.02) conv:(11.66)
2. [Indiretas4, Pares1, CC1, Cuidadores1, Cuidadores2]: 845 ==> [CC6]: 811
<conf:(0.96)> lift:(1.81) lev:(0.02) conv:(11.38)
3. [Indiretas3, Indiretas4, Pares1, CC3, Cuidadores1, Cuidadores2]: 789 ==> [CC6]: 757
<conf:(0.96)> lift:(1.81) lev:(0.02) conv:(11.26)
4. [Indiretas4, Pares1, CC3, Cuidadores1, Cuidadores2]: 935 ==> [CC6]: 895
<conf:(0.96)> lift:(1.81) lev:(0.03) conv:(10.74)
5. [Indiretas4, Pares1, Cuidadores1, CC4, Cuidadores2]: 873 ==> [CC6]: 835
<conf:(0.96)> lift:(1.81) lev:(0.03) conv:(10.55)
6. [Indiretas2, Indiretas3, Indiretas4, Pares1, Cuidador1, Cuidador2]: 909 ==> [CC6]: 868
<conf:(0.95)> lift:(1.81) lev:(0.03) conv:(10.2)
7. [Indiretas4, Pares1, CC1, CC3, Cuidadores2]: 833 ==> [CC6]: 795
<conf:(0.95)> lift:(1.8) lev:(0.02) conv:(10.06)
8. [Indiretas2, Indiretas4, Pares1, Cuidadores1, Cuidadores2]: 1079 ==> [CC6]: 1029
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.97)
9. [Indiretas4, Pares1, CC3, CC4, Cuidadores2]: 815 ==> [CC6]: 777
<conf:(0.95)> lift:(1.8) lev:(0.02) conv:(9.85)
10. [Indiretas4, Pares1, Cuidadores1, Pares3, Cuidadores2]: 834 ==> [CC6]: 794
<conf:(0.95)> lift:(1.8) lev:(0.02) conv:(9.58)
11. [Indiretas2, Indiretas4, Pares1, CC4, Cuidadores2]: 949 ==> [CC6]: 903
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.51)
12. [Indiretas2, Pares1, CC3, Cuidadores1, Cuidadores2]: 949 ==> [CC6]: 903
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.51)
13. [Indiretas2, Indiretas3, Indiretas4, Pares1, CC3, Cuidadores2]: 902 ==> [CC6]: 858
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.44)
14. [Indiretas2, Indiretas4, Pares1, CC3, Cuidadores2]: 1060 ==> [CC6]: 1008
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.42)
15. [Indiretas2, Indiretas3, Indiretas4, Pares1, CC4, Cuidadores2]: 811 ==> [CC6]: 771
<conf:(0.95)> lift:(1.8) lev:(0.02) conv:(9.32)
16. [Indiretas3, Indiretas4, Pares1, Cuidadores1, Cuidadores2]: 1071 ==> [CC6]: 1018
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.34)
17. [Indiretas3, Indiretas4, Pares1, CC3, Cuidadores2]: 1048 ==> [CC6]: 996
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.32)
18. [Indiretas2, Pares1, CC1, Cuidadores1, Cuidadores2]: 839 ==> [CC6]: 797
<conf:(0.95)> lift:(1.8) lev:(0.02) conv:(9.19)
19. [Indiretas4, Pares1, CC4, Cuidadores2]: 1111 ==> [CC6]: 1055
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.18)
20. [Indiretas3, Indiretas4, Pares1, CC4, Cuidadores2]: 931 ==> [CC6]: 884
<conf:(0.95)> lift:(1.8) lev:(0.03) conv:(9.14)
21. [Indiretas2, Indiretas3, Pares1, CC3, Cuidadores1, Cuidadores2]: 807 ==> [CC6]: 766
<conf:(0.95)> lift:(1.79) lev:(0.02) conv:(9.05)
22. [Indiretas2, Pares1, CC1, CC4, Cuidadores2]: 784 ==> [CC6]: 744
<conf:(0.95)> lift:(1.79) lev:(0.02) conv:(9.01)
23. [Indiretas2, Pares1, CC1, CC3, Cuidadores2]: 818 ==> [CC6]: 776
<conf:(0.95)> lift:(1.79) lev:(0.02) conv:(8.96)

24. [Indiretas4, Pares1, Cuidadores1, Cuidadores2]: 1301 ==> [CC6]: 1234
<conf:(0.95)> lift:(1.79) lev:(0.04) conv:(9.01)
25. [Indiretas4, Pares1, CC3, Cuidadores2]: 1259 ==> [CC6]: 1194
<conf:(0.95)> lift:(1.79) lev:(0.04) conv:(8.99)

- Regras para conjunto de dados das vitimizações em vida utilizando o tipo de vitimização para a agregação de atributos (25 regras de um total de 182).

1. [CC, Pares, Cuidadores, Digitais, Sexuais]: 1653 ==> [Indiretas]: 1627
<conf:(0.98)> lift:(1.19) lev:(0.02) conv:(10.6)
2. [CC, Cuidadores, Digitais, Sexuais]: 1775 ==> [Indiretas]: 1746
<conf:(0.98)> lift:(1.19) lev:(0.02) conv:(10.25)
3. [Pares, Cuidadores, Digitais, Sexuais]: 1689 ==> [Indiretas]: 1658
<conf:(0.98)> lift:(1.19) lev:(0.02) conv:(9.14)
4. [Cuidadores, Digitais, Sexuais]: 1828 ==> [Indiretas]: 1794
<conf:(0.98)> lift:(1.19) lev:(0.02) conv:(9.05)
5. [Indiretas, Pares, Cuidadores, Digitais, Sexuais]: 1658 ==> [CC]: 1627
<conf:(0.98)> lift:(1.28) lev:(0.02) conv:(11.99)
6. [CC, Pares, Digitais, Sexuais]: 1943 ==> [Indiretas]: 1904
<conf:(0.98)> lift:(1.19) lev:(0.02) conv:(8.41)
7. [CC, Pares, Cuidadores, Digitais]: 2881 ==> [Indiretas]: 2823
<conf:(0.98)> lift:(1.19) lev:(0.03) conv:(8.46)
8. [CC, Digitais, Sexuais]: 2129 ==> [Indiretas]: 2084
<conf:(0.98)> lift:(1.18) lev:(0.02) conv:(8.02)
9. [Pares, Cuidadores, Digitais, Sexuais]: 1689 ==> [CC]: 1653
<conf:(0.98)> lift:(1.27) lev:(0.02) conv:(10.56)
10. [Indiretas, Pares, Cuidadores, Sexuais]: 2377 ==> [CC]: 2324
<conf:(0.98)> lift:(1.27) lev:(0.03) conv:(10.19)
11. [CC, Cuidadores, Digitais]: 3194 ==> [Indiretas]: 3122
<conf:(0.98)> lift:(1.18) lev:(0.03) conv:(7.58)
12. [Pares, Cuidadores, Digitais]: 2979 ==> [Indiretas]: 2911
<conf:(0.98)> lift:(1.18) lev:(0.03) conv:(7.48)
13. [Pares, Digitais, Sexuais]: 2004 ==> [Indiretas]: 1958
<conf:(0.98)> lift:(1.18) lev:(0.02) conv:(7.38)
14. [Digitais, Sexuais]: 2229 ==> [Indiretas]: 2175
<conf:(0.98)> lift:(1.18) lev:(0.02) conv:(7.02)
15. [CC, Pares, Cuidadores, Sexuais]: 2385 ==> [Indiretas]: 2324
<conf:(0.97)> lift:(1.18) lev:(0.02) conv:(6.66)
16. [Cuidadores, Digitais]: 3362 ==> [Indiretas]: 3274
<conf:(0.97)> lift:(1.18) lev:(0.03) conv:(6.54)
17. [Pares, Cuidadores, Sexuais]: 2450 ==> [CC]: 2385
<conf:(0.97)> lift:(1.27) lev:(0.03) conv:(8.59)
18. [Indiretas, Cuidadores, Digitais, Sexuais]: 1794 ==> [CC]: 1746
<conf:(0.97)> lift:(1.27) lev:(0.02) conv:(8.47)
19. [Indiretas, Pares, Digitais, Sexuais]: 1958 ==> [CC]: 1904
<conf:(0.97)> lift:(1.27) lev:(0.03) conv:(8.24)
20. [CC, Pares, Digitais]: 3694 ==> [Indiretas]: 3590
<conf:(0.97)> lift:(1.18) lev:(0.04) conv:(6.09)
21. [Cuidadores, Digitais, Sexuais]: 1828 ==> [CC]: 1775
<conf:(0.97)> lift:(1.26) lev:(0.03) conv:(7.83)
22. [Pares, Cuidadores, Sexuais]: 2450 ==> [Indiretas]: 2377
<conf:(0.97)> lift:(1.17) lev:(0.02) conv:(5.73)
23. [Indiretas, Pares, Cuidadores, Digitais]: 2911 ==> [CC]: 2823
<conf:(0.97)> lift:(1.26) lev:(0.04) conv:(7.57)
24. [Pares, Digitais, Sexuais]: 2004 ==> [CC]: 1943

<conf:(0.97)> lift:(1.26) lev:(0.03) conv:(7.48)
 25. [CC, Cuidadores, Sexuais]: 2673 ==> [Indiretas]: 2588
 <conf:(0.97)> lift:(1.17) lev:(0.03) conv:(5.38)

- Regras para conjunto de dados das vitimizações em vida utilizando o tipo de violência para a agregação de atributos (25 regras de um total de 29).

1. [Negligência, Violência Física, Abuso Sexual]: 3494 ==> [Violência Psicológica]: 3387
 <conf:(0.97)> lift:(1.22) lev:(0.04) conv:(6.68)
2. [Violência Física, Abuso Sexual]: 3629 ==> [Violência Psicológica]: 3509
 <conf:(0.97)> lift:(1.22) lev:(0.04) conv:(6.19)
3. [Violência Psicológica, Violência Física, Abuso Sexual]: 3509 ==> [Negligência]: 3387
 <conf:(0.97)> lift:(1.17) lev:(0.03) conv:(4.93)
4. [Violência Física, Abuso Sexual]: 3629 ==> [Negligência]: 3494
 <conf:(0.96)> lift:(1.16) lev:(0.03) conv:(4.61)
5. [Violência Psicológica, Abuso Sexual]: 3966 ==> [Negligência]: 3792
 <conf:(0.96)> lift:(1.16) lev:(0.03) conv:(3.92)
6. [Negligência, Abuso Sexual]: 4002 ==> [Violência Psicológica]: 3792
 <conf:(0.95)> lift:(1.19) lev:(0.04) conv:(3.92)
7. [Abuso Sexual]: 4227 ==> [Negligência]: 4002
 <conf:(0.95)> lift:(1.14) lev:(0.03) conv:(3.23)
8. [Abuso Sexual]: 4227 ==> [Violência Psicológica]: 3966
 <conf:(0.94)> lift:(1.18) lev:(0.04) conv:(3.33)
9. [Violência Física, Abuso Sexual]: 3629 ==> [Negligência, Violência Psicológica]: 3387
 <conf:(0.93)> lift:(1.31) lev:(0.05) conv:(4.32)
10. [Negligência, Violência Física]: 9286 ==> [Violência Psicológica]: 8591
 <conf:(0.93)> lift:(1.17) lev:(0.08) conv:(2.76)
11. [Violência Psicológica, Violência Física]: 9288 ==> [Negligência]: 8591
 <conf:(0.92)> lift:(1.12) lev:(0.06) conv:(2.3)
12. [Violência Física]: 10193 ==> [Violência Psicológica]: 9288
 <conf:(0.91)> lift:(1.15) lev:(0.08) conv:(2.32)
13. [Violência Física]: 10193 ==> [Negligência]: 9286
 <conf:(0.91)> lift:(1.1) lev:(0.06) conv:(1.94)
14. [Abuso Sexual]: 4227 ==> [Negligência, Violência Psicológica]: 3792
 <conf:(0.9)> lift:(1.26) lev:(0.05) conv:(2.8)
15. [Violência Psicológica]: 11664 ==> [Negligência]: 10449
 <conf:(0.9)> lift:(1.08) lev:(0.05) conv:(1.66)
16. [Negligência, Violência Psicológica, Abuso Sexual]: 3792 ==> [Violência Física]: 3387
 <conf:(0.89)> lift:(1.29) lev:(0.05) conv:(2.86)
17. [Violência Psicológica, Abuso Sexual]: 3966 ==> [Violência Física]: 3509
 <conf:(0.88)> lift:(1.28) lev:(0.05) conv:(2.65)
18. [Negligência, Abuso Sexual]: 4002 ==> [Violência Física]: 3494
 <conf:(0.87)> lift:(1.26) lev:(0.05) conv:(2.41)
19. [Negligência]: 12159 ==> [Violência Psicológica]: 10449
 <conf:(0.86)> lift:(1.08) lev:(0.05) conv:(1.47)
20. [Abuso Sexual]: 4227 ==> [Violência Física]: 3629
 <conf:(0.86)> lift:(1.24) lev:(0.05) conv:(2.16)
21. [Violência Psicológica, Abuso Sexual]: 3966 ==> [Negligência, Violência Física]: 3387
 <conf:(0.85)> lift:(1.35) lev:(0.06) conv:(2.52)
22. [Negligência, Abuso Sexual]: 4002 ==> [Violência Psicológica, Violência Física]: 3387
 <conf:(0.85)> lift:(1.34) lev:(0.06) conv:(2.39)
23. [Violência Física]: 10193 ==> [Negligência, Violência Psicológica]: 8591
 <conf:(0.84)> lift:(1.19) lev:(0.09) conv:(1.84)
24. [Abuso Sexual]: 4227 ==> [Violência Psicológica, Violência Física]: 3509
 <conf:(0.83)> lift:(1.31) lev:(0.06) conv:(2.16)

25. [Abuso Sexual]: 4227 ==> [Negligência, Violência Física]: 3494
 <conf:(0.83)> lift:(1.31) lev:(0.06) conv:(2.12)

- Regras para conjunto de dados das vitimizações no último ano sem agregação de atributos.

1. [Indiretas3, Indiretas5]: 973 ==> [Indiretas2]: 847
 <conf:(0.87)> lift:(2) lev:(0.03) conv:(4.33)
2. [Indiretas4, Indiretas5]: 1095 ==> [Indiretas2]: 938
 <conf:(0.86)> lift:(1.97) lev:(0.03) conv:(3.92)
3. [Indiretas2, Indiretas3, Pares1]: 884 ==> [CC6]: 750
 <conf:(0.85)> lift:(2.99) lev:(0.03) conv:(4.69)
4. [Indiretas4, Pares1]: 1134 ==> [CC6]: 956
 <conf:(0.84)> lift:(2.97) lev:(0.04) conv:(4.54)
5. [Indiretas3, Pares1]: 1284 ==> [CC6]: 1055
 <conf:(0.82)> lift:(2.89) lev:(0.05) conv:(4)
6. [Indiretas5]: 1484 ==> [Indiretas2]: 1216
 <conf:(0.82)> lift:(1.89) lev:(0.04) conv:(3.12)
7. [Indiretas2, Pares1]: 1305 ==> [CC6]: 1066
 <conf:(0.82)> lift:(2.87) lev:(0.05) conv:(3.89)
8. [Indiretas3, Indiretas5]: 973 ==> [Indiretas4]: 790
 <conf:(0.81)> lift:(2.63) lev:(0.03) conv:(3.66)
9. [Cuidadores1, Pares1]: 957 ==> [CC6]: 777
 <conf:(0.81)> lift:(2.86) lev:(0.03) conv:(3.78)
10. [Indiretas4, Digitais1]: 1061 ==> [Indiretas3]: 834
 <conf:(0.79)> lift:(1.87) lev:(0.03) conv:(2.69)
11. [Pares1]: 2184 ==> [CC6]: 1705
 <conf:(0.78)> lift:(2.75) lev:(0.07) conv:(3.26)
12. [Indiretas3, Indiretas4, Cuidadores1]: 1337 ==> [Indiretas2]: 1041
 <conf:(0.78)> lift:(1.79) lev:(0.03) conv:(2.55)
13. [Indiretas2, Indiretas4, Cuidadores1]: 1340 ==> [Indiretas3]: 1041
 <conf:(0.78)> lift:(1.84) lev:(0.03) conv:(2.58)
14. [Indiretas2, Indiretas5]: 1216 ==> [Indiretas4]: 938
 <conf:(0.77)> lift:(2.5) lev:(0.04) conv:(3.01)
15. [Indiretas4, CC4]: 1080 ==> [Indiretas2]: 832
 <conf:(0.77)> lift:(1.77) lev:(0.02) conv:(2.45)
16. [Indiretas4, CC6, Cuidadores1]: 1007 ==> [Indiretas3]: 768
 <conf:(0.76)> lift:(1.81) lev:(0.02) conv:(2.43)
17. [Indiretas4, CC6, Cuidadores1]: 1007 ==> [Indiretas2]: 767
 <conf:(0.76)> lift:(1.75) lev:(0.02) conv:(2.36)
18. [Indiretas3, Indiretas4, CC6]: 1451 ==> [Indiretas2]: 1104
 <conf:(0.76)> lift:(1.75) lev:(0.03) conv:(2.36)
19. [Indiretas4, CC3]: 1211 ==> [Indiretas3]: 918
 <conf:(0.76)> lift:(1.8) lev:(0.03) conv:(2.38)
20. [Indiretas2, Indiretas4, CC6]: 1461 ==> [Indiretas3]: 1104
 <conf:(0.76)> lift:(1.79) lev:(0.03) conv:(2.36)
21. [Indiretas4, Digitais1]: 1061 ==> [Indiretas2]: 800
 <conf:(0.75)> lift:(1.73) lev:(0.02) conv:(2.29)

- Regras para conjunto de dados das vitimizações no último ano utilizando o tipo de vitimização para a agregação de atributos (25 regras de um total de 56).

1. [Cuidadores, Digitais, Sexuais]: 805 ==> [Indiretas]: 768
 <conf:(0.95)> lift:(1.42) lev:(0.02) conv:(6.91)

2. [CC, Cuidadores, Pares, Sexuais]: 799 ==> [Indiretas]: 757
<conf:(0.95)> lift:(1.41) lev:(0.01) conv:(6.06)
3. [Cuidadores, Pares, Sexuais]: 862 ==> [Indiretas]: 813
<conf:(0.94)> lift:(1.4) lev:(0.02) conv:(5.62)
4. [CC, Digitais, Sexuais]: 966 ==> [Indiretas]: 911
<conf:(0.94)> lift:(1.4) lev:(0.02) conv:(5.63)
5. [CC, Cuidadores, Pares, Digitais]: 1007 ==> [Indiretas]: 947
<conf:(0.94)> lift:(1.4) lev:(0.02) conv:(5.38)
6. [Digitais, Sexuais]: 1152 ==> [Indiretas]: 1077
<conf:(0.93)> lift:(1.39) lev:(0.02) conv:(4.94)
7. [CC, Cuidadores, Digitais]: 1473 ==> [Indiretas]: 1377
<conf:(0.93)> lift:(1.39) lev:(0.03) conv:(4.95)
8. [Cuidadores, Pares, Digitais]: 1092 ==> [Indiretas]: 1018
<conf:(0.93)> lift:(1.38) lev:(0.02) conv:(4.75)
9. [Indiretas, Cuidadores, Pares, Sexuais]: 813 ==> [CC]: 757
<conf:(0.93)> lift:(1.83) lev:(0.02) conv:(7)
10. [Indiretas, Cuidadores, Pares, Digitais]: 1018 ==> [CC]: 947
<conf:(0.93)> lift:(1.83) lev:(0.03) conv:(6.94)
11. [CC, Pares, Sexuais]: 1099 ==> [Indiretas]: 1019
<conf:(0.93)> lift:(1.38) lev:(0.02) conv:(4.42)
12. [Cuidadores, Pares, Sexuais]: 862 ==> [CC]: 799
<conf:(0.93)> lift:(1.82) lev:(0.02) conv:(6.61)
13. [Cuidadores, Digitais]: 1757 ==> [Indiretas]: 1627
<conf:(0.93)> lift:(1.37) lev:(0.03) conv:(4.37)
14. [CC, Cuidadores, Sexuais]: 1169 ==> [Indiretas]: 1081
<conf:(0.92)> lift:(1.37) lev:(0.02) conv:(4.28)
15. [Cuidadores, Pares, Digitais]: 1092 ==> [CC]: 1007
<conf:(0.92)> lift:(1.81) lev:(0.03) conv:(6.23)
16. [Pares, Sexuais]: 1218 ==> [Indiretas]: 1121
<conf:(0.92)> lift:(1.37) lev:(0.02) conv:(4.05)
17. [CC, Pares, Digitais]: 1441 ==> [Indiretas]: 1326
<conf:(0.92)> lift:(1.37) lev:(0.02) conv:(4.05)
18. [Cuidadores, Sexuais]: 1393 ==> [Indiretas]: 1274
<conf:(0.91)> lift:(1.36) lev:(0.02) conv:(3.79)
19. [Indiretas, Pares, Sexuais]: 1121 ==> [CC]: 1019
<conf:(0.91)> lift:(1.78) lev:(0.03) conv:(5.34)
20. [CC, Digitais]: 2292 ==> [Indiretas]: 2082
<conf:(0.91)> lift:(1.35) lev:(0.04) conv:(3.54)
21. [Pares, Digitais]: 1621 ==> [Indiretas]: 1472
<conf:(0.91)> lift:(1.35) lev:(0.03) conv:(3.52)
22. [CC, Sexuais]: 1808 ==> [Indiretas]: 1641
<conf:(0.91)> lift:(1.35) lev:(0.03) conv:(3.51)
23. [CC, Cuidadores, Pares]: 2029 ==> [Indiretas]: 1834
<conf:(0.9)> lift:(1.34) lev:(0.03) conv:(3.38)
24. [Pares, Sexuais]: 1218 ==> [CC]: 1099
<conf:(0.9)> lift:(1.77) lev:(0.03) conv:(4.98)
25. [Indiretas, Pares, Digitais]: 1472 ==> [CC]: 1326
<conf:(0.9)> lift:(1.77) lev:(0.04) conv:(4.91)

- Regras para conjunto de dados das vitimizações no último ano utilizando o tipo de violência para a agregação de atributos.

1. [Violência Psicológica, Violência Física, Abuso Sexual]: 1646 ==> [Negligência]: 1516
<conf:(0.92)> lift:(1.35) lev:(0.03) conv:(4.02)
2. [Negligência, Violência Física, Abuso Sexual]: 1649 ==> [Violência Psicológica]: 1516

- <conf:(0.92)> lift:(1.61) lev:(0.04) conv:(5.26)
3. [Violência Física, Abuso Sexual]: 1805 ==> [Negligência]: 1649
<conf:(0.91)> lift:(1.34) lev:(0.03) conv:(3.68)
 4. [Violência Física, Abuso Sexual]: 1805 ==> [Violência Psicológica]: 1646
<conf:(0.91)> lift:(1.59) lev:(0.04) conv:(4.82)
 5. [Violência Psicológica, Abuso Sexual]: 2360 ==> [Negligência]: 2133
<conf:(0.9)> lift:(1.33) lev:(0.04) conv:(3.32)
 6. [Abuso Sexual]: 2908 ==> [Negligência]: 2560
<conf:(0.88)> lift:(1.3) lev:(0.04) conv:(2.67)
 7. [Violência Psicológica, Violência Física]: 4991 ==> [Negligência]: 4300
<conf:(0.86)> lift:(1.27) lev:(0.06) conv:(2.31)
 8. [Negligência, Violência Física]: 5074 ==> [Violência Psicológica]: 4300
<conf:(0.85)> lift:(1.48) lev:(0.09) conv:(2.8)
 9. [Violência Física, Abuso Sexual]: 1805 ==> [Negligência, Violência Psicológica]: 1516
<conf:(0.84)> lift:(1.8) lev:(0.05) conv:(3.33)
 10. [Negligência, Abuso Sexual]: 2560 ==> [Violência Psicológica]: 2133
<conf:(0.83)> lift:(1.46) lev:(0.05) conv:(2.56)
 11. [Violência Física]: 6111 ==> [Negligência]: 5074
<conf:(0.83)> lift:(1.22) lev:(0.06) conv:(1.89)
 12. [Violência Física]: 6111 ==> [Violência Psicológica]: 4991
<conf:(0.82)> lift:(1.43) lev:(0.1) conv:(2.33)
 13. [Violência Psicológica]: 8415 ==> [Negligência]: 6842
<conf:(0.81)> lift:(1.2) lev:(0.08) conv:(1.71)
 14. [Abuso Sexual]: 2908 ==> [Violência Psicológica]: 2360
<conf:(0.81)> lift:(1.42) lev:(0.05) conv:(2.26)

Referências

- AFIFI, T. O.; MACMILLAN, H. L. Resilience following child maltreatment: A review of protective factors. *The Canadian Journal of Psychiatry*, SAGE Publications Sage CA: Los Angeles, CA, v. 56, n. 5, p. 266–272, 2011. Citado na página 25.
- AGRAWAL, R.; IMIELIŃSKI, T.; SWAMI, A. Mining association rules between sets of items in large databases. In: ACM. *Acm sigmod record*. [S.l.], 1993. v. 22, n. 2, p. 207–216. Citado 5 vezes nas páginas 20, 28, 37, 38 e 60.
- AGRAWAL, R.; SRIKANT, R. et al. Fast algorithms for mining association rules. In: *Proc. 20th int. conf. very large data bases, VLDB*. [S.l.: s.n.], 1994. v. 1215, p. 487–499. Citado na página 28.
- AIKEN, L. S.; WEST, S. G.; RENO, R. R. *Multiple regression: Testing and interpreting interactions*. [S.l.]: Sage, 1991. Citado na página 44.
- ALLISON, P. D. *Missing data*. [S.l.]: Sage publications, 2001. Citado na página 29.
- ALPAYDIN, E. *Introduction to machine learning*. [S.l.]: MIT press, 2009. Citado na página 27.
- ALTEMEIER, W. A. et al. Prediction of child abuse: A prospective study of feasibility. *Child Abuse & Neglect*, Elsevier, v. 8, n. 4, p. 393–400, 1984. Citado 3 vezes nas páginas 42, 48 e 57.
- ALTMAN, D. G.; BLAND, J. M. Diagnostic tests. 1: Sensitivity and specificity. *BMJ: British Medical Journal*, BMJ Publishing Group, v. 308, n. 6943, p. 1552, 1994. Citado na página 35.
- ALVES, L. et al. Crime prediction through urban metrics and statistical learning. *Physica A*, p. 435–443, 2018. Citado na página 19.
- AMRIT, C. et al. Identifying child abuse through text mining and machine learning. *Expert systems with applications*, Elsevier, v. 88, p. 402–418, 2017. Citado 2 vezes nas páginas 49 e 58.
- B. THORNBERRY T., S. C. K. In the wake of childhood maltreatment. 1997. Citado na página 24.
- BABBIE, E. R. *The basics of social research*. [S.l.]: Cengage Learning, 2013. Citado na página 54.
- BAIR, E. Semi-supervised clustering methods. *Wiley Interdisciplinary Reviews: Computational Statistics*, Wiley Online Library, v. 5, n. 5, p. 349–361, 2013. Citado na página 28.
- BEGLE, A. M.; DUMAS, J. E.; HANSON, R. F. Predicting child abuse potential: An empirical investigation of two theoretical frameworks. *Journal of Clinical Child & Adolescent Psychology*, Taylor & Francis, v. 39, n. 2, p. 208–219, 2010. Citado 4 vezes nas páginas 9, 45, 46 e 58.

- BELLMAN, R. *Dynamic programming*. [S.l.]: Courier Corporation, 2013. Citado na página 30.
- BENESTY, J. et al. Pearson correlation coefficient. In: *Noise reduction in speech processing*. [S.l.]: Springer, 2009. p. 1–4. Citado na página 30.
- BIRLESON, P. The validity of depressive disorder in childhood and the development of a self-rating scale: a research report. *Journal of Child Psychology and Psychiatry*, Wiley Online Library, v. 22, n. 1, p. 73–88, 1981. Citado na página 66.
- BISHOP, C. M. *Pattern Recognition and Machine Learning*. [S.l.]: Springer, 2006. Citado na página 26.
- BLATTBERG, R. C.; KIM, B.-D.; NESLIN, S. A. Evaluation of statistical models. In: *Database Marketing*. [S.l.]: Springer, 2008. p. 309–310. Citado na página 34.
- BOGOMOLOV, A. et al. Once upon a crime: Towards crime prediction from demographics and mobile data. *16th International Conference on Multimodal Interaction*, p. 427–434, 2014. Citado na página 19.
- BORGAN, Ø. et al. *Handbook of Statistical Methods for Case-Control Studies*. [S.l.]: CRC Press, 2018. Citado na página 54.
- BOYD, K.; ENG, K. H.; PAGE, C. D. Area under the precision-recall curve: point estimates and confidence intervals. In: SPRINGER. *Joint European conference on machine learning and knowledge discovery in databases*. [S.l.], 2013. p. 451–466. Citado na página 36.
- BREIMAN, L. *Classification and regression trees*. [S.l.]: Routledge, 2017. Citado na página 33.
- BRIN, S. et al. Dynamic itemset counting and implication rules for market basket data. *Acm Sigmod Record*, ACM, v. 26, n. 2, p. 255–264, 1997. Citado na página 40.
- BROWNE, M. et al. *Informe Final Análisis Multivariable de Estudio Polivictimización*. [S.l.], 2018. Citado na página 66.
- BURRELL, B.; THOMPSON, B.; SEXTON, D. Predicting child abuse potential across family types. *Child abuse & neglect*, Elsevier, v. 18, n. 12, p. 1039–1049, 1994. Citado 2 vezes nas páginas 42 e 57.
- BUTCHART, A.; HARVEY, A. *Prevención del maltrato infantil: qué hacer, y cómo obtener evidencias*. [S.l.], 2009. Citado na página 25.
- CAPLAN, J. M.; KENNEDY, L. W. Risk terrain modeling compendium. *Rutgers Center on Public Security, Newark*, 2011. Citado na página 51.
- CHANDRASHEKAR, G.; SAHIN, F. A survey on feature selection methods. *Computers & Electrical Engineering*, Elsevier, v. 40, n. 1, p. 16–28, 2014. Citado 2 vezes nas páginas 30 e 31.
- CHAPELLE, O.; SCHOLKOPF, B.; ZIEN, A. Semi-supervised learning (chappelle, o. et al., eds.; 2006)[book reviews]. *IEEE Transactions on Neural Networks*, IEEE, v. 20, n. 3, p. 542–542, 2009. Citado na página 28.

- CHAWLA, N. V. et al. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, v. 16, p. 321–357, 2002. Citado na página 68.
- CHEN, T.; GUESTRIN, C. Xgboost: A scalable tree boosting system. In: ACM. *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*. [S.l.], 2016. p. 785–794. Citado na página 55.
- CHILD WELFARE INFORMATION GATEWAY. *Understanding the Effects of Maltreatment on Early Brain Development*. [S.l.], 2001. Citado na página 24.
- CHILD WELFARE INFORMATION GATEWAY. *Risk and Protective Factors for Child Abuse and Neglect*. [S.l.], 2004. 2–3 p. Citado na página 25.
- CHILD WELFARE INFORMATION GATEWAY. *Long-Term Consequences of Child Abuse and Neglect*. [S.l.], 2013. 3–7 p. Citado na página 24.
- CHILD WELFARE INFORMATION GATEWAY. *Protective Factors Approaches in Child Welfare*. [S.l.], 2014. 1 p. Citado na página 25.
- CHOULDECHOVA, A. et al. A case study of algorithm-assisted decision making in child maltreatment hotline screening decisions. In: *Conference on Fairness, Accountability and Transparency*. [S.l.: s.n.], 2018. p. 134–148. Citado 4 vezes nas páginas 11, 55, 56 e 58.
- COELHO, L. P.; RICHERT, W. *Building machine learning systems with Python*. [S.l.]: Packt Publishing Ltd, 2015. Citado na página 31.
- COLMAN, R. A.; WIDOM, C. S. Childhood abuse and neglect and adult intimate relationships: A prospective study. *Child abuse & neglect*, Elsevier, v. 28, n. 11, p. 1133–1151, 2004. Citado na página 24.
- COOK, A. et al. Complex trauma in children and adolescents. *Psychiatric annals*, SLACK Incorporated, v. 35, n. 5, p. 390–398, 2005. Citado na página 26.
- COX, D. R. Regression models and life-tables. *Journal of the Royal Statistical Society: Series B (Methodological)*, Wiley Online Library, v. 34, n. 2, p. 187–202, 1972. Citado na página 44.
- COX, D. R. *Analysis of survival data*. [S.l.]: Routledge, 2018. Citado na página 44.
- CRECECONTIGO. *El Maltrato Infantil*. 2015. Citado na página 24.
- DALEY, D. et al. Risk terrain modeling predicts child maltreatment. *Child abuse & neglect*, Elsevier, v. 62, p. 29–38, 2016. Citado 4 vezes nas páginas 9, 51, 52 e 58.
- DAVIS, J.; GOADRICH, M. The relationship between precision-recall and roc curves. In: ACM. *Proceedings of the 23rd international conference on Machine learning*. [S.l.], 2006. p. 233–240. Citado na página 36.
- DEPANFILIS, D.; ZURAVIN, S. J. Predicting child maltreatment recurrences during treatment. *Child Abuse & Neglect*, Elsevier, v. 23, n. 8, p. 729–743, 1999. Citado 2 vezes nas páginas 44 e 57.
- DETOURS, V. et al. Integration and cross-validation of high-throughput gene expression data: comparing heterogeneous data sets. *FEBS letters*, Wiley Online Library, v. 546, n. 1, p. 98–102, 2003. Citado na página 29.

- DONG, M. et al. The interrelatedness of multiple forms of childhood abuse, neglect, and household dysfunction. *Child abuse & neglect*, Elsevier, v. 28, n. 7, p. 771–784, 2004. Citado na página 26.
- DUBOWITZ, H. et al. Identifying children at high risk for a child maltreatment report. *Child abuse & neglect*, Elsevier, v. 35, n. 2, p. 96–104, 2011. Citado 2 vezes nas páginas 46 e 58.
- DURLAK, J. A. Common risk and protective factors in successful prevention programs. *American journal of orthopsychiatry*, Wiley Online Library, v. 68, n. 4, p. 512–520, 1998. Citado na página 25.
- FARREN, D. Factores asociados al maltrato infantil en adolescentes escolares de la comuna de recoleta. *IV Congreso Nacional de Investigación sobre Violencia y Delincuencia*, p. 169–193, 2007. Citado na página 18.
- FAWCETT, T. An introduction to roc analysis. *Pattern recognition letters*, Elsevier, v. 27, n. 8, p. 861–874, 2006. Citado na página 36.
- FELITTI, V. J. et al. Relationship of childhood abuse and household dysfunction to many of the leading causes of death in adults: The adverse childhood experiences (ace) study. *American journal of preventive medicine*, Elsevier, v. 14, n. 4, p. 245–258, 1998. Citado na página 24.
- FINKELHOR, D. et al. The juvenile victimization questionnaire: reliability, validity, and national norms. *Child abuse & neglect*, Elsevier, v. 29, n. 4, p. 383–412, 2005. Citado 2 vezes nas páginas 26 e 65.
- FINKELHOR, D.; ORMROD, R. K.; TURNER, H. A. Poly-victimization: A neglected component in child victimization. *Child abuse & neglect*, Elsevier, v. 31, n. 1, p. 7–26, 2007. Citado 2 vezes nas páginas 25 e 26.
- FINKELHOR, D. et al. Polyvictimization: Children’s exposure to multiple types of violence, crime, and abuse. *National survey of children’s exposure to violence*, United States Department of Justice, 2011. Citado na página 26.
- FLAHERTY, C. W.; PATTERSON, D. A. Predicting child physical abuse recurrence: comparison of a neural network to logistic regression. *Journal of Technology in Human Services*, Taylor & Francis, v. 21, n. 4, p. 93–111, 2003. Citado 2 vezes nas páginas 53 e 57.
- FRANK, R. et al. Frequency of child abuse seen in a pediatric surgery unit. *Pediatric surgery international*, Springer, v. 7, n. 6, p. 454–458, 1992. Citado 2 vezes nas páginas 47 e 57.
- GAN, G.; MA, C.; WU, J. *Data clustering: theory, algorithms, and applications*. [S.l.]: Siam, 2007. Citado na página 27.
- GARCÍA, S.; LUENGO, J.; HERRERA, F. *Data preprocessing in data mining*. [S.l.]: Springer, 2015. Citado na página 29.
- GAYATHRI, K. et al. Hierarchical activity recognition or dementia care using markov logic network. *Personal and Ubiquitous Computing*, p. 271–285, 2015. Citado na página 19.

- GELMAN, A.; HILL, J. *Data analysis using regression and multilevel/hierarchical models*. [S.l.]: Cambridge university press, 2006. Citado na página 42.
- GOES, A.; STEINER, M. Proposta de metodologia para a criação de etiquetas de classificação - estudo de caso: Desempenho escolar. *Gestão e Produção*, p. 177–191, 2016. Citado na página 19.
- GÓMEZ, L. *Menores Víctimas y Testigos de Violencia Familiar*. Tese (Doutorado) — Facultad de Derecho, Universidad Zaragoza, 2011. Citado na página 18.
- GUYON, I.; ELISSEEFF, A. An introduction to variable and feature selection. *Journal of machine learning research*, v. 3, n. Mar, p. 1157–1182, 2003. Citado na página 31.
- GUYON, I. et al. *Feature extraction: foundations and applications*. [S.l.]: Springer, 2008. Citado na página 30.
- HALKIDI, M.; BATISTAKIS, Y.; VAZIRGIANNIS, M. On clustering validation techniques. *Journal of intelligent information systems*, Springer, v. 17, n. 2-3, p. 107–145, 2001. Citado na página 39.
- HASTIE, T.; TIBSHIRANI, R.; FRIEDMAN, J. Unsupervised learning. In: *The elements of statistical learning*. [S.l.]: Springer, 2009. p. 485–585. Citado 2 vezes nas páginas 27 e 37.
- HAWKINS, D. M. The problem of overfitting. *Journal of chemical information and computer sciences*, ACS Publications, v. 44, n. 1, p. 1–12, 2004. Citado na página 30.
- HE, H.; GARCIA, E. A. Learning from imbalanced data. *IEEE Transactions on Knowledge & Data Engineering*, Ieee, n. 9, p. 1263–1284, 2008. Citado na página 29.
- HORIKAWA, H. et al. Development of a prediction model for child maltreatment recurrence in japan: A historical cohort study using data from a child guidance center. *Child abuse & neglect*, Elsevier, v. 59, p. 55–65, 2016. Citado 2 vezes nas páginas 47 e 58.
- HOSMER, D. W.; LEMESHOW, S.; STURDIVANT, R. X. *Applied logistic regression*. [S.l.]: John Wiley & Sons, 2013. Citado 2 vezes nas páginas 52 e 55.
- JAKSCH, T.; ORTNER, R.; AUER, P. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, v. 11, n. Apr, p. 1563–1600, 2010. Citado na página 28.
- JAMES, G. et al. *An introduction to statistical learning*. [S.l.]: Springer, 2013. Citado na página 32.
- KASS, G. V. An exploratory technique for investigating large quantities of categorical data. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, Wiley Online Library, v. 29, n. 2, p. 119–127, 1980. Citado na página 33.
- KIM, W. et al. A taxonomy of dirty data. *Data mining and knowledge discovery*, Springer, v. 7, n. 1, p. 81–99, 2003. Citado na página 29.
- KOHAVI, R. et al. A study of cross-validation and bootstrap for accuracy estimation and model selection. In: MONTREAL, CANADA. *Ijcai*. [S.l.], 1995. v. 14, n. 2, p. 1137–1145. Citado 3 vezes nas páginas 34, 47 e 72.

- KOHONEN, T. The self-organizing map. *Proceedings of the IEEE*, IEEE, v. 78, n. 9, p. 1464–1480, 1990. Citado na página 28.
- KOTSIANTIS, S.; KANELLOPOULOS, D. Association rules mining: A recent overview. *GESTS International Transactions on Computer Science and Engineering*, v. 32, n. 1, p. 71–82, 2006. Citado na página 38.
- LECANNELIER, F. *El trauma oculto en la infancia. Guía científicamente informada para padres, educadores y profesionales*. [S.l.]: Penguin Random House Grupo Editorial, 2018. Citado na página 24.
- LEE, J. A.; VERLEYSEN, M. *Nonlinear dimensionality reduction*. [S.l.]: Springer Science & Business Media, 2007. Citado na página 30.
- LIN, T. Y. Attribute transformations for data mining i: Theoretical explorations. *International journal of intelligent systems*, Wiley Online Library, v. 17, n. 2, p. 213–222, 2002. Citado na página 29.
- LINDSAY, M. et al. Neural predictors of initialiting alcohol use during adolescence. *The American Journal of Psychiatry*, p. 172–185, 2017. Citado na página 19.
- LITTLE, J.; RIXON, A. Computer learning and risk assessment in child protection. *Child Abuse Review: Journal of the British Association for the Study and Prevention of Child Abuse and Neglect*, Wiley Online Library, v. 7, n. 3, p. 165–177, 1998. Citado 4 vezes nas páginas 9, 48, 49 e 57.
- LIU, F. et al. Dual teaching: A practical semi-supervised wrapper method. *arXiv preprint arXiv:1611.03981*, 2016. Citado na página 28.
- LIU, H.; MOTODA, H. *Feature selection for knowledge discovery and data mining*. [S.l.]: Springer Science & Business Media, 2012. Citado na página 31.
- MAGLOGIANNIS ILIAS, G. *Emerging artificial intelligence applications in computer engineering: real word ai systems with applications in ehealth, hci, information retrieval and pervasive technologies*. [S.l.]: Ios Press, 2007. Citado na página 26.
- MALDONADO, S.; WEBER, R. Modelos de selección de atributos para support vector machines. *Revista Ingeniería de Sistemas*, v. 26, p. 49–70, 2012. Citado na página 31.
- MARQUARDT, D. W.; SNEE, R. D. Ridge regression in practice. *The American Statistician*, Taylor & Francis Group, v. 29, n. 1, p. 3–20, 1975. Citado na página 31.
- MARSHALL, D. B.; ENGLISH, D. J. Neural network modeling of risk assessment in child protective services. *Psychological Methods*, American Psychological Association, v. 5, n. 1, p. 102, 2000. Citado 2 vezes nas páginas 52 e 57.
- MASTEN, A. S.; GARMEZY, N. Risk, vulnerability, and protective factors in developmental psychopathology. In: *Advances in clinical child psychology*. [S.l.]: Springer, 1985. p. 39–40. Citado na página 25.
- MCDONALD, T.; MARKS, J. A review of risk factors assessed in child protective services. *Social Service Review*, University of Chicago Press, v. 65, n. 1, p. 112–132, 1991. Citado na página 25.

- MCHUGH, M. L. The chi-square test of independence. *Biochemia medica: Biochemia medica*, Medinska naklada, v. 23, n. 2, p. 143–149, 2013. Citado 2 vezes nas páginas 30 e 45.
- MILNER, J. Características sociales y psicológicas de los maltratadores físicos del niño. *II Congreso Estatal sobre Infancia Maltratada*, 1993. Citado na página 18.
- MOHD, P. et al. A neural network students' performance prediction model (nnsppm). *IEEE International Conference on Smart Instrumentation, Measurement and Applications (ICSIMA)*, 2013. Citado na página 19.
- MOND, B.; PECARIĆ, J. A mixed arithmetic-mean-harmonic-mean matrix inequality. *Linear algebra and its applications*, Elsevier, v. 237, p. 449–454, 1996. Citado na página 36.
- MONTGOMERY, D. C.; PECK, E. A.; VINING, G. G. *Introduction to linear regression analysis*. [S.l.]: John Wiley & Sons, 2012. Citado na página 45.
- NISBET, R.; ELDER, J.; MINER, G. *Handbook of statistical analysis and data mining applications*. [S.l.]: Academic Press, 2009. Citado na página 28.
- OLVERA-LÓPEZ, J. A. et al. A review of instance selection methods. *Artificial Intelligence Review*, Springer, v. 34, n. 2, p. 133–143, 2010. Citado na página 31.
- ORDONEZ, C. Association rule discovery with the train and test approach for heart disease prediction. *IEEE Transactions on Information Technology in Biomedicine*, IEEE, v. 10, n. 2, p. 334–343, 2006. Citado na página 39.
- PAMPEL, F. C. *Logistic regression: A primer*. [S.l.]: Sage, 2000. Citado na página 27.
- PIATETSKY-SHAPIO, G. Discovery, analysis, and presentation of strong rules. *Knowledge discovery in databases*, Menlo Park, CA: AAI/MIT, p. 229–238, 1991. Citado na página 40.
- PYLE, D. *Data preparation for data mining*. [S.l.]: morgan kaufmann, 1999. Citado na página 29.
- QUINLAN, J. R. Discovering rules by induction from large collections of examples. *Expert systems in the micro electronics age*, Edinburgh University Press, 1979. Citado na página 33.
- QUINLAN, J. R. Induction of decision trees. *Machine learning*, Springer, v. 1, n. 1, p. 81–106, 1986. Citado na página 33.
- QUINLAN, J. R. *C4. 5: programs for machine learning*. [S.l.]: Morgan Kaufmann, 1993. Citado 3 vezes nas páginas 33, 34 e 60.
- RODRIGUEZ, C. M.; GREEN, A. J. Parenting stress and anger expression as predictors of child abuse potential. *Child Abuse & Neglect*, Pergamon, v. 21, n. 4, p. 367–377, 1997. Citado 2 vezes nas páginas 43 e 57.
- ROJAS, R. *Neural networks: a systematic introduction*. [S.l.]: Springer Science & Business Media, 2013. Citado na página 27.

- ROKACH, L.; MAIMON, O. Z. *Data mining with decision trees: theory and applications*. [S.l.]: World scientific, 2008. Citado 4 vezes nas páginas 20, 27, 32 e 71.
- ROSENBERG, M. *Society and the adolescent self-image*. [S.l.]: Princeton university press, 1965. Citado na página 66.
- RUDAS, T. *Odds ratios in the analysis of contingency tables*. [S.l.]: Sage, 1998. Citado na página 54.
- RUMELHART, D. E.; HINTON, G. E.; WILLIAMS, R. J. *Learning internal representations by error propagation*. [S.l.], 1985. Citado na página 52.
- RUSSELL, J. Predictive analytics and child protection: Constraints and opportunities. *Child Abuse and Neglect*, p. 182–189, 2015. Citado na página 19.
- RUSSELL, S. J.; NORVIG, P. *Artificial intelligence: a modern approach*. [S.l.]: Malaysia; Pearson Education Limited, 2016. Citado 2 vezes nas páginas 27 e 39.
- RUTHERFORD, A. *Introducing ANOVA and ANCOVA: a GLM approach*. [S.l.]: Sage, 2001. Citado na página 30.
- SCHWARTZ, D. R.; KAUFMAN, A. B.; SCHWARTZ, I. M. Computational intelligence techniques for risk assessment and decision support. *Children and Youth Services Review*, Elsevier, v. 26, n. 11, p. 1081–1095, 2004. Citado 2 vezes nas páginas 48 e 57.
- SCHWARTZ, I. M. et al. Predictive and prescriptive analytics, machine learning and child welfare risk assessment: The broward county experience. *Children and Youth Services Review*, Elsevier, v. 81, p. 309–320, 2017. Citado 2 vezes nas páginas 50 e 58.
- SEBER, G. A.; LEE, A. J. *Linear regression analysis*. [S.l.]: John Wiley & Sons, 2012. Citado na página 52.
- SENAME. *SENAME por tus Derechos, Maltrato Infantil y Adolescente*. 2016. Citado na página 23.
- SHANNON, C. E. A mathematical theory of communication. *Bell system technical journal*, Wiley Online Library, v. 27, n. 3, p. 379–423, 1948. Citado na página 33.
- SHAVLIK, J. et al. *Readings in machine learning*. [S.l.]: Morgan Kaufmann, 1990. Citado na página 27.
- SILVERMAN, A. B.; REINHERZ, H. Z.; GIACONIA, R. M. The long-term sequelae of child and adolescent abuse: A longitudinal community study. *Child abuse & neglect*, Elsevier, v. 20, n. 8, p. 709–723, 1996. Citado na página 24.
- SLEDJESKI, E. M. et al. The use of risk assessment to predict recurrent maltreatment: A classification and regression tree analysis (cart). *Prevention science*, Springer, v. 9, n. 1, p. 28–37, 2008. Citado 2 vezes nas páginas 53 e 57.
- SØGAARD, A. Semi-supervised learning and domain adaptation in natural language processing. *Synthesis Lectures on Human Language Technologies*, Morgan & Claypool Publishers, v. 6, n. 2, p. 1–103, 2013. Citado na página 28.

- STITH, S. M. et al. Risk factors in child maltreatment: A meta-analytic review of the literature. *Aggression and violent behavior*, Elsevier, v. 14, n. 1, p. 13–29, 2009. Citado na página 25.
- SUBSECRETARÍA DE PREVENCIÓN DEL DELITO. *Primera Encuesta Nacional de Polivictimización en Niñas, Niños y Adolescentes*. [S.l.], 2018. Citado 2 vezes nas páginas 20 e 65.
- SUTTON, R. S.; BARTO, A. G. *Introduction to reinforcement learning*. [S.l.]: MIT press Cambridge, 1998. Citado na página 29.
- THURSTON, H.; MIYAMOTO, S. The use of model based recursive partitioning as an analytic tool in child welfare. *Child abuse & neglect*, Elsevier, v. 79, p. 293–301, 2018. Citado 3 vezes nas páginas 54, 56 e 58.
- TIBSHIRANI, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, Wiley Online Library, v. 58, n. 1, p. 267–288, 1996. Citado na página 31.
- UNICEF. *Maltrato Infantil en Chile*. [S.l.], 2005. 2–3 p. Citado na página 23.
- UNICEF. *Convención sobre los Derechos del Niño*. 2014. Citado na página 17.
- UNICEF. *Una Situación Habitual*. [S.l.], 2017. 3–6 p. Citado na página 17.
- VAITHIANATHAN, R. et al. Children in the public benefit system at risk of maltreatment: Identification via predictive modeling. *American journal of preventive medicine*, Elsevier, v. 45, n. 3, p. 354–359, 2013. Citado 2 vezes nas páginas 51 e 58.
- VAN DER KOLK, B. A. Developmental trauma disorder: Toward a rational diagnosis for children with complex trauma histories. *Psychiatric annals*, SLACK Incorporated, v. 35, n. 5, p. 401–408, 2005. Citado na página 25.
- WANG, L. *Support vector machines: theory and applications*. [S.l.]: Springer Science & Business Media, 2005. Citado na página 27.
- WITTEN, I. H. et al. *Data Mining: Practical Machine Learning Tools and Techniques*. [S.l.]: Morgan Kaufmann, 2016. 7–8 p. Citado 2 vezes nas páginas 19 e 26.
- WU, J. *Advances in K-means clustering: a data mining thinking*. [S.l.]: Springer Science & Business Media, 2012. Citado na página 28.
- YAP, B. W. et al. An application of oversampling, undersampling, bagging and boosting in handling imbalanced datasets. In: SPRINGER. *Proceedings of the First International Conference on Advanced Data and Information Engineering (DaEng-2013)*. [S.l.], 2014. p. 13–22. Citado na página 29.
- ZHANG, C.; ZHANG, S. *Association rule mining: models and algorithms*. [S.l.]: Springer-Verlag, 2002. Citado 3 vezes nas páginas 38, 39 e 40.
- ZOU, H.; HASTIE, T. Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, Wiley Online Library, v. 67, n. 2, p. 301–320, 2005. Citado na página 31.