

UNIVERSIDADE FEDERAL DO RIO GRANDE DO SUL
INSTITUTO DE FILOSOFIA E CIÊNCIAS HUMANAS
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA POLÍTICA

TIAGO VIER

**O USO DA INTELIGÊNCIA ARTIFICIAL NAS CIÊNCIAS SOCIAIS:
o caso do patriotismo dos brasileiros.**

Porto Alegre
2020

TIAGO VIER

**O USO DA INTELIGÊNCIA ARTIFICIAL NAS CIÊNCIAS SOCIAIS:
o caso do patriotismo dos brasileiros.**

Tese submetida ao Programa de Pós-Graduação em Ciência Política da Universidade Federal do Rio Grande do Sul como requisito parcial para obtenção do título de Doutor em Ciência Política.

Orientador: Prof. Dr. Henrique Carlos de Oliveira de Castro

Linha de Pesquisa: Cultura Política

Porto Alegre
2020

CIP - Catalogação na Publicação

Vier, Tiago

O uso da Inteligência Artificial nas ciências
sociais: o caso do patriotismo dos brasileiros. /
Tiago Vier. - 2020
189 f.

Orientador: Henrique Carlos de Oliveira de Castro.

Tese (Doutorado) - Universidade Federal do Rio
Grande do Sul, Instituto de Filosofia e Ciências
Humanas, Programa de Pós-Graduação em Ciência
Política, Porto Alegre, BR-RS, 2020.

1. Inteligência Artificial. 2. Aprendizado de
máquina. 3. Cultura Política. 4. Pesquisa Mundial de
Valores. 5. Patriotismo. I. Castro, Henrique Carlos de
Oliveira de, orient. II. Título.

TIAGO VIER

**O USO DA INTELIGÊNCIA ARTIFICIAL NAS CIÊNCIAS SOCIAIS:
o caso do patriotismo dos brasileiros.**

Tese submetida ao Programa de Pós-Graduação em Ciência Política da Universidade Federal do Rio Grande do Sul como requisito parcial para obtenção do título de Doutor em Ciência Política.

Orientador: Prof. Dr. Henrique Carlos de Oliveira de Castro

Aprovado em: Porto Alegre, 27 de maio de 2020.

BANCA EXAMINADORA:

Orientador: Prof. Dr. Henrique Carlos de Oliveira de Castro
Faculdade de Ciências Econômicas (UFRGS)

Prof. Dr. Daniel Jaime Capistrano de Oliveira
School of Education (UCD-University College of Dublin)

Prof. Dr. Dante Augusto Couto Barone
Instituto de Informática (UFRGS)

Prof. Dr. Luís Gustavo Mello Grohmann
Instituto de Filosofia e Ciências Humanas (UFRGS)

Porto Alegre
2020

À Professora Sonia Ranincheski.

AGRADECIMENTOS

Os meus agradecimentos vão, em primeiro lugar, ao meu orientador, Henrique Carlos de Castro. Esta tese só foi possível porque ele aceitou o desafio de orientar um estudante com outra formação e percurso profissional e, com a paciência e a atenção necessárias, me apoiou totalmente durante o período de tese. Agradeço imensamente pelas conversas, variadas, divertidas e inspiradoras, e pelas orientações durante esta já saudosa jornada.

Agradeço também à banca de avaliação desta tese, composta pelos Professores Dante Augusto Couto Barone, Daniel Capistrano e Luís Gustavo Mello Grohmann, pela rigorosa leitura do meu trabalho e pelos preciosos comentários e críticas construtivas que me foram transmitidas durante a defesa. Também agradeço aos professores do PPG em Ciência Política com os quais pude interagir e, em especial, da Linha de Pesquisa em Cultura Política, Rodrigo Stumpf González e Sonia Ranincheski, pela atenção e pelos ensinamentos dispensados durante o curso.

Em seguida, devo agradecer também à sociedade brasileira, que me deu a oportunidade de estudar em uma universidade pública, gratuita e de qualidade como a Universidade Federal do Rio Grande do Sul. Agradeço também à CAPES, pelo apoio ao projeto que permitiu a realização da Sétima Onda da Pesquisa Mundial de Valores no Brasil e no âmbito do qual pude me beneficiar de um período de seis meses no exterior como bolsista. Essa etapa foi essencial para o meu desenvolvimento acadêmico. A internacionalização permanece, a meu ver, uma etapa essencial na formação de acadêmicos e pesquisadores capazes de produzir conhecimento original, de qualidade e passível de impacto na sociedade. Pelo acolhimento e pelo apoio durante a estada na Unil em Lausanne durante o período da bolsa, agradeço ao Professor Christian Thöni e à Dr^a Deborah Kistler.

Outro conjunto de agradecimentos é devido aos meus colegas do grupo de pesquisa em Cultura Política, da Pesquisa Mundial de Valores Brasil e do Laboratório de Pesquisa Pirandello com quem compartilhei muitos momentos de aprendizado, trabalho e diversão nestes últimos quatro anos. Infelizmente, não posso citar todos nominalmente, mas agradeço, em especial, aos colegas que tiveram alguma participação nos trabalhos que desenvolvi durante esse período.

Não posso deixar de agradecer à minha família e, em especial, meus pais, por terem despertado o gosto pela leitura e pelo aprendizado e pelo apoio em todas as empreitadas.

Por fim, um agradecimento especial à Carolina Brito, minha parceira de tudo, pelo apoio e pelo incentivo para continuar com minha formação acadêmica e profissional. Além da leitura crítica e da revisão do texto da tese, devo a ela muito do meu amadurecimento como pesquisador e acadêmico iniciante. Esta tese não teria sentido sem a sua presença intelectual e afetiva constante.

Il est triste que souvent pour être un bon patriote on soit l'ennemi du reste des hommes.(...) Être bon patriote, c'est souhaiter que la ville s'enrichisse par le commerce, et soit puissante par les armes. Il est clair qu'un pays ne peut gagner sans qu'un autre perde, et qu'il ne peut vaincre sans faire des malheureux. Tel est donc la condition humaine, que souhaiter la grandeur de son pays c'est souhaiter du mal à ses voisins. Celui qui voudrait que la patrie ne fût jamais ni plus grande, ni plus petite, ni plus riche, ni plus pauvre, ferait le citoyen de l'univers

Voltaire, Dictionnaire philosophique, portatif.
1764

RESUMO

Este trabalho propõe a integração da inteligência artificial (IA) nas ciências sociais por meio de ferramentas de aprendizado de máquina. Para avaliar a utilidade desse ferramental para a ciência política, foi escolhido um tema clássico da área: o patriotismo. Em termos metodológicos, a tese explora uma forma de integrar as ferramentas de aprendizado de máquina no processo de produção do conhecimento em um momento indutivo, no qual auxilia o cientista a produzir hipóteses sobre o fenômeno estudado. No caso deste trabalho, usando dados das sete ondas da Pesquisa Mundial de Valores (WVS) para 114 países, coletados entre 1981 e 2019, o aprendizado de máquina permitiu gerar hipóteses explicativas sobre o fenômeno do patriotismo e sua expressão na forma de orgulho nacional. Entre essas hipóteses, figura a forte ligação com valores tradicionais, como a religião e a família, e a crença em formas de organização social sustentadas por instituições de natureza autoritária e hierárquica. A possibilidade de gerar hipóteses usando essas ferramentas é um resultado favorável à integração da IA nas ciências sociais. As dificuldades e os desafios relacionados ao tratamento dos dados e as inferências a partir de modelos preditivos apontam para a necessidade de constituição de equipes multidisciplinares para implementação de projetos de pesquisa mais ambiciosos.¹

Palavras Chave: **1. Inteligência Artificial. 2. Aprendizado de máquina. 3. Cultura Política. 4. Pesquisa Mundial de Valores. 5. Patriotismo.**

¹VIER, Tiago. **O uso da Inteligência Artificial nas ciências sociais: o caso do patriotismo dos brasileiros.** Porto Alegre, 2020. 189 f. Tese (Doutorado em Ciência Política) - Programa de Pós-Graduação em Ciência Política, Instituto de Filosofia e Ciências Humanas. Universidade Federal do Rio Grande do Sul, Porto Alegre, 2020

ABSTRACT

This dissertation proposes the integration of artificial intelligence (AI) in the social sciences through the use of machine learning techniques. A classical theme from the political sciences, patriotism, was chosen to evaluate the utility of these tools and methods. In methodological terms, the dissertation explores the integration of machine learning in the process of knowledge production in the social sciences. Integrated in what is called here an inductive moment, machine learning might support the scientist in generating hypotheses about the phenomena under study. In this work, using data from seven waves of the World Values Survey (WVS) for 114 countries collected between 1981 and 2019, machine learning allowed to generate explanatory hypotheses for patriotism as pride for one's nationality. The most important insights are related to the linkages found between pride and traditional values, like religion and family, and the belief in forms of social organization that are sustained by institutions of authoritarian and hierarchical nature. The very possibility of generating hypotheses using these techniques is in itself a result that supports the integration of AI techniques in the social sciences. The challenges and difficulties related to data engineering and making explanatory inferences from predictive models indicates the need of constituting multidisciplinary teams in order to implement more ambitious research projects.²

Keywords: **1. Artificial Intelligence. 2. Machine Learning. 3. Political Culture. 4. World Values Survey. 5. Patriotism.**

²VIER, Tiago. **O uso da Inteligência Artificial nas ciências sociais: o caso do patriotismo dos brasileiros**. Porto Alegre, 2020. 189 f. Tese (Doutorado em Ciência Política) - Programa de Pós-Graduação em Ciência Política, Instituto de Filosofia e Ciências Humanas. Universidade Federal do Rio Grande do Sul, Porto Alegre, 2020

LISTA DE TABELAS

| | | |
|----|--|-----|
| 1 | Orgulho de ser Brasileiro 1991-2018 | 16 |
| 2 | Disparidades entre os caminhos preditivos e explicativos | 60 |
| 3 | Métricas utilizadas no aprendizado de máquina | 70 |
| 4 | Variáveis usadas na regressão logística | 79 |
| 5 | Resultados da comparação dos algoritmos | 83 |
| 6 | Performance do Extreme Boosting | 85 |
| 7 | Resultados das previsões para o WVS7 e Brasil | 91 |
| 8 | Amostra de Trabalho extraída do WVS por sociedade e onda | 123 |
| 9 | Algoritmos testados no benchmark | 140 |
| 10 | Artigos empíricos sobre patriotismo e orgulho nacional | 142 |

LISTA DE FIGURAS

| | | |
|----|--|-----|
| 1 | Patriotismo e Nacionalismo na base Google Books, 1800-2008 | 18 |
| 2 | Esquema de pesquisa usada no trabalho | 22 |
| 3 | Patriotismos e suas influências. | 51 |
| 4 | Surveys realizadas pelo WVS e EVS entre 1981 e 2019 | 54 |
| 5 | Separação da base de dados do WVS | 56 |
| 6 | Procedimento adotado na modelagem preditiva | 68 |
| 7 | Resultados da regressão logística em razão de chances | 81 |
| 8 | Resultados dos algoritmos testados | 84 |
| 9 | Curva ROC treinamento e validação | 86 |
| 10 | Curva ROC comparada para a amostras do WVS7 e Brasil | 88 |
| 11 | Matrizes de Confusão para as amostras do WVS7 e Brasil | 90 |
| 12 | Importância das variáveis no modelo XgBoost | 95 |
| 13 | Importância das variáveis usando o método da permutação | 96 |
| 14 | Diferença entre os métodos de importancia global | 97 |
| 15 | Impacto das variáveis agrupadas | 98 |
| 16 | Importância local das variáveis para seis países. | 101 |
| 17 | Importância local das variáveis para o Brasil | 102 |
| 18 | Não-respostas na base de trabalho | 128 |
| 19 | Correlação entre as predições dos algoritmos de melhor desempenho. | 137 |

*O presente trabalho foi realizado com apoio da
Coordenação de Aperfeiçoamento de Pessoal de
Nível Superior – Brasil (CAPES) – Código N°
88887.186170/2018-00.*

SUMÁRIO

| | |
|---|-----------|
| INTRODUÇÃO | 15 |
| A ESTRUTURA DA TESE | 20 |
| 1 INTELIGÊNCIA ARTIFICIAL E APRENDIZADO DE MÁQUINA | 24 |
| 1.1 APRENDIZADO DE MÁQUINA | 28 |
| 1.2 APLICAÇÕES DA IA NAS CIÊNCIAS SOCIAIS | 30 |
| 2 PATRIOTISMO | 33 |
| 2.1 DEFINIÇÃO HISTÓRICA DE PATRIOTISMO | 34 |
| 2.1.1 O patriotismo no Brasil | 37 |
| 2.2 A PESQUISA EMPÍRICA SOBRE PATRIOTISMO | 40 |
| 2.2.1 Estudos empíricos sobre o patriotismo dos brasileiros | 48 |
| 3 METODOLOGIA | 53 |
| 3.1 A PESQUISA MUNDIAL DE VALORES | 53 |
| 3.1.1 Amostras de treino, validação e predição | 55 |
| 3.1.2 Seleção e transformação das variáveis | 56 |
| 3.1.3 Tratamento de não-respostas | 58 |
| 3.2 DEDUÇÃO E INDUÇÃO NA MODELAGEM ESTATÍSTICA | 60 |
| 3.3 MOMENTO DEDUTIVO | 61 |
| 3.4 MOMENTO INDUTIVO | 63 |
| 3.5 ANÁLISE, INFERÊNCIA E INTERPRETAÇÃO DOS RESULTADOS | 71 |
| 3.6 AMBIENTE DE DESENVOLVIMENTO E PACOTES | 75 |

| | | |
|----------|---|------------|
| 4 | RESULTADOS | 77 |
| 4.1 | RESULTADOS DA MODELAGEM EXPLICATIVA | 77 |
| 4.2 | RESULTADOS DO APRENDIZADO DE MÁQUINA | 82 |
| 4.2.1 | Estudo comparativo | 82 |
| 4.2.2 | Aprendizado com o XGBoost | 85 |
| 4.3 | PREDIÇÃO NAS AMOSTRAS WVS7 E BRASIL | 87 |
| 4.3.1 | Probabilidades | 87 |
| 4.3.2 | Predições | 89 |
| 4.4 | INFERÊNCIA | 92 |
| 4.4.1 | Importância global | 93 |
| 4.4.2 | Importância local | 98 |
| 5 | DISCUSSÃO | 103 |
| 5.1 | IMPLICAÇÃO PARA AS TEORIAS DO PATRIOTISMO | 104 |
| 5.2 | IMPLICAÇÕES METODOLÓGICAS | 110 |
| | CONCLUSÃO | 113 |
| | APÊNDICE | 122 |
| A | APÊNDICE A - AMOSTRA DE TRABALHO | 123 |
| B | APÊNDICE B - NÃO-RESPOSTAS E IMPUTAÇÃO | 127 |
| C | APÊNDICE C - ESTUDO COMPARATIVO | 134 |
| D | APÊNDICE D - ARTIGOS EMPÍRICOS SOBRE PATRIOTISMO | 142 |
| | BIBLIOGRAFIA | 153 |
| | ÍNDICE REMISSIVO | 189 |

INTRODUÇÃO

A inteligência artificial (IA) vive um novo momento de euforia. Técnicas de aprendizado de máquina têm sido aplicadas com sucesso em diversas áreas: extração de padrões de imagens e vídeos (YANG et al., 2015), tradução de línguas mortas (LUO; CAO; BARZILAY, 2019), decisões judiciais (BERK, 2012), mitigação de mudanças climáticas (ATHEY, 2017), detecção de fraudes (BOLTON; HAND, 2002), identificação de doenças (SOUZA et al., 2020), e elaboração de programas de educação sob medida (LUCKIN et al., 2016). Além de jogarem GO ou xadrez melhor que humanos (SILVER et al., 2016), já existem sistemas inteligentes capazes de conduzir veículos autônomos (ENZWEILER, 2015).

Em paralelo a esses admiráveis movimentos na ciência da computação, as ciências sociais passam por um momento de importante questionamento. Por um lado, parecem ter perdido a sua utilidade social devido à pulverização dos temas de pesquisas e da reticência ou incapacidade em explicar fenômenos da atualidade (CALHOUN; WIEVIORKA, 2013). Por outro, estão em crise com seus próprios métodos e técnicas. Em particular, veem-se um descrédito com relação à validade e à replicabilidade das conclusões científicas devido ao uso desinformado de métodos quantitativos tradicionais e o desafio em tratar os chamados “big-data” (KING, 2014; CLARK; GOLDBER, 2015; WASSERSTEIN; LAZAR, 2016).

Como as ciências sociais e os cientistas sociais podem beneficiar-se dos avanços na área da inteligência artificial para construir conhecimento e explicar fenômenos sociais? O que a inteligência artificial e o aprendizado de máquina podem fazer pela ciência social e como? Mais precisamente, como pode ser aplicado o conjunto de métodos e técnicas de aprendizado de máquina para explicar os fenômenos ligados ao homem e suas interações sociais, e, mais especificamente, a prática humana relacionada ao poder e ao Estado — a política (BOBBIO, 1993)? Esse é o problema central do trabalho e a pergunta central que tentamos responder.

Espera-se que a inteligência artificial integrada à prática das ciências sociais na forma do aprendizado de máquina possa aumentar, consideravelmente, a capacidade de produzir conhecimento social. Para verificar se esta tese é plausível e explorar maneiras de usar a IA nas ciências sociais escolhemos abordar um tema caro à ciência política e em evidência no debate público atual: o patriotismo.

A escolha desse tema foi motivada, primeiramente, pela observação do declínio das declarações de orgulho de ser brasileiro na Pesquisa Mundial de Valores. As declarações de orgulho nacional são uma forma comum de observar o nível de patriotismo em uma sociedade, como veremos adiante. Conforme pode ser visto na Tabela 1, desde a primeira onda, no início dos anos 1990, a frequência das declarações de orgulho cai de forma importante no Brasil.

Tabela 1: Orgulho de ser Brasileiro 1991-2019. A tabela abaixo apresenta os dados das cinco ondas do WVS realizadas no Brasil entre 1991 e 2018. Os quatro níveis da variável original foram agrupados em dois níveis para computar as frequências em cada onda realizada

| | 1991 | 1997 | 2006 | 2014 | 2019 | Total |
|--|------|------|------|------|-------|-------|
| Alto Orgulho - Muito orgulhoso e Orgulhoso | | | | | | |
| Casos | 1534 | 953 | 1249 | 1144 | 1084 | 5964 |
| Frequência | 86.6 | 83.8 | 83.8 | 77.7 | 62.9 | |
| <i>Resíduo</i> | - | - | - | -0.9 | -18.0 | |
| Baixo Orgulho - Não muito orgulhoso e Não sou orgulhoso | | | | | | |
| Casos | 237 | 184 | 241 | 329 | 640 | 1631 |
| Frequência | 13.4 | 16.2 | 16.2 | 22.3 | 37.1 | |
| Total | 1771 | 1137 | 1490 | 1473 | 1724 | 7595 |

^a A variável original foi transformada em uma variável dicotômica da auto-apreciação individual do orgulho da sua nacionalidade na qual as respostas 'Não muito orgulhoso' e Nada Orgulhoso são re-codificados como Baixo Orgulho, e Muito Orgulho, Orgulhoso são codificados como Alto Orgulho.

^b Os resíduos são a diferença entre os eventos observados e esperados ajustados para os totais de cada onda e resposta. Em regra geral, um resíduo fora do intervalo -2-2 indica um número de casos esperados significativamente menor ou maior do que o esperado no caso de hipótese nula verdadeira

A escolha do patriotismo é também motivada pela importância acadêmica do tema. O patriotismo é um conceito clássico da ciência política, que está presente desde o surgimento da

disciplina. Transformado em nacionalismo, tornou-se para alguns o “fio condutor” da história, que atravessa a modernidade desde a Revolução Francesa até os tempos atuais (SMITH, 1998). Benedict Anderson (1991) escreveu, em 1991, que a era dos nacionalismos não havia terminado e que a condição nacional era ainda o valor de maior legitimidade na vida política dos nossos tempos³.

Dentro desse tema, a contribuição pretendida com a tese se enquadra na perspectiva teórica da cultura política na qual o patriotismo faz parte de um subconjunto dos fenômenos culturais, moldado por crenças, normas e valores políticos. Comportamento que é estudado desde as origens da própria ciência política: na *República*, Platão (1991, p. 86) se preocupava com a “disposição dos homens” para a república e com a convicção patriótica a que os cidadãos nunca deveriam renunciar. Maquiavel, e mais tarde Montesquieu, exaltaram o patriotismo dos Romanos como uma das formas de sustentação do Império. Nas revoluções liberais, a cultura política patriótica estava presente na maneira como Rousseau (2003, p. 83) descreveu a importância das leis gravadas nos corações dos cidadãos e da necessidade de consentimento e sentimentos patrióticos que isso implicaria para o sucesso do contrato social⁴. Tocqueville (2001, p. 338) reconheceu uma forma de patriotismo reflexivo nos costumes dos norte-americanos determinante para o sucesso da democracia. Inspirado pelos antigos, e em particular por este último, o “*The Civic Culture*”, estudo que inaugurou a área da cultura política, também considerou o patriotismo como parte de uma cultura “cívica” e determinante para a democracia (ALMOND, 1956, 1989, p. 5).

A relevância do tema tem relação, também, com as explicações sobre o homem e a sociedade de forma geral. C. W. Mills (1959) e J. Elster (1989, 1994) argumentam que uma das tarefas mais importantes das ciências sociais é explicar por que não estamos mais no “estado da natureza” e por que sociedades têm um determinado nível de ordenamento. Entender o nível de patriotismo ajuda a explicar o grau de cooperação e o quanto os indivíduos podem prever as ações de seus compatriotas, e, assim, o grau de ordenamento de determinada sociedade.

Visto de forma positiva, o patriotismo, nesse caso, aumenta a cooperação entre os cidadãos e reduz o conflito social. Quanto mais patrióticos, mais ordenados e distantes estaríamos do estado da natureza. Além de envolver os cidadãos em busca de justiça social para

³ *Indeed, nation-ness is the most universally legitimate value in the political life of our time.*

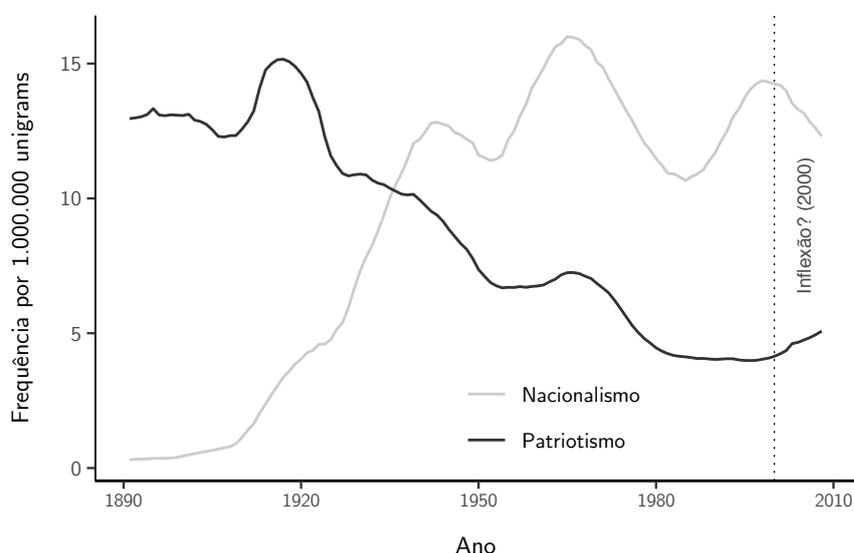
⁴ *“laws engraved in citizens’ hearths”... “local situation and temper of inhabitants”.*

seus compatriotas (BAR-TAL; STAUB, 1997; VIROLI, 1997), constitui-se em um reservatório de apoio difuso à comunidade política e às suas instituições em tempos de crise (EASTON, 1965).

Mas o patriotismo também apresenta riscos à coesão interna de uma sociedade. O impacto negativo do apego exacerbado a um grupo é o aumento do preconceito em relação a outros grupos, externos ou internos, que não compartilhem a mesma cultura, etnia ou linguagem. Retóricas chauvinistas são incentivadas, dificultando a adoção de visões mais cosmopolitas de mundo sustentadas por valores humanos comuns (NUSSBAUM; COHEN, 2002).

A escolha do tema parece em linha também com o seu ressurgimento no debate acadêmico. Depois de um período de hegemonia do estudo do nacionalismo, o interesse pelo patriotismo parece ter ressurgido a partir dos anos 2000. Desde então, vê-se um crescente interesse em livros e publicações sobre o tema, acompanhado de um declínio por outro termo indissociável, o nacionalismo (ver Figura 1).

Figura 1: Patriotismo e Nacionalismo na base Google Books, 1800-2008



Fonte: Elaboração a partir do Google NGram. Dados suavizados por regressão local com banda de 10 anos

A renovação no interesse acadêmico indica também que se trata de um tema importante para a sociedade. De fato, as narrativas patrióticas parecem ter ressurgido no debate público. A partir dos anos 2000, por razões distintas, em diferentes partes do mundo,

a necessidade de conter a globalização e valorizar a integridade dos estados nacionais passou a fazer parte do discurso político. Na Europa, diversas demandas por narrativas nacionais surgiram com força e geraram mudanças nos sistemas de educação nacionais para “ensinar” o patriotismo nas escolas e novas leis regulando a memória nacional⁵. Além disso, o patriotismo motivou novos movimentos geopolíticos, como a anexação da Crimeia e o próprio Brexit⁶.

Apelos ao patriotismo também se fizeram presentes em processos eleitorais em diversos países. A condição nacional passou a ser disputada por diferentes grupos políticos, tanto com cores do ultra-nacionalismo da extrema-direita⁷, quanto com cores mais democráticas e compatíveis com a globalização e valores compartilhados⁸.

No Brasil, o discurso patriótico e da soberania nacional foi resgatado como tema das eleições de 2018. O lema “Brasil acima de tudo e Deus acima de todos” e a bandeira nacional como símbolo de campanha aglutinaram os brasileiros descontentes com o desempenho econômico, com a corrupção e com a violência, em torno de valores ligados à pátria e à

⁵A partir do momento em que oito estados da antiga União Soviética foram integrados à União Europeia (UE) em 2004, a necessidade de se adaptar às legislações vigentes na UE, contra o Holocausto, por exemplo, gerou demandas internas por narrativas visando marcar o sofrimento da era comunista (SOROKA; KRAWATZEK, 2019). No mesmo período, o presidente Vladimir Putin deu início às tentativas de reabilitar o passado soviético, estratégia que culminou na modificação da lei da educação em 2012 em direção a promover patriotismo e identidade nacional (TSYRLINA-SPADY; LOVORN, 2015). No Japão, desde sua primeira gestão como Primeiro-Ministro, entre 2006 e 2007, o conservador Shinzo Abe também tem se esforçado para reorientar a cultura japonesa em direção à valorização do passado e ao orgulho da linhagem imperial, tornando igualmente compulsório o ensino do patriotismo nas escolas (BBC, 2007; SIEG, 2014). Segundo Soroka e Krawatzek (2019), nas duas últimas décadas, foram aprovadas mais de duzentas leis, resoluções e declarações regulando a memória histórica, a maioria vinda da Europa, e diversas delas tentando resgatar o passado e fomentar o amor à pátria.

⁶Em 2014, Vladimir Putin fomentou e permitiu a anexação da Crimeia usando justificativas relacionadas à identidade cultural dos Crimeios. A narrativa nacionalista também venceu com o Brexit, o “British Exit”, no referendo em 2016. Nos Estados Unidos, o gatilho para o resgate do patriotismo foi o ataque às torres gêmeas em 2001, que alterou a maneira como os americanos viam seu excepcionalismo e alterou a visão positiva da globalização. Essa mudança de orientação pode ter sido um dos fatores que culminaram com a eleição de Donald Trump em 2017 e com o incremento das políticas protecionistas do *America First*, nomeadamente irrompendo um conflito comercial com a China.

⁷Em 2004, na França, o *Front Nacional*, com discurso claramente nacionalista, e muitas vezes xenofóbico, obteve 25% dos votos nas eleições para o parlamento europeu. Na Alemanha e na Itália, o patriotismo foi incorporado no discurso da extrema-direita, de forma antidemocrática e xenofóbica, ou, ainda, islamofóbica como ilustram movimentos como o *Patriotic Europeans against Islamization of the Occident* (PEGIDA). Na Alemanha, o partido *Alternative für Deutschland* (AfD) propaga o “*wir sind das Volk*” (nós somos o povo). Na Itália, o *Lega*, de Matteo Salvini, usam como palavras de ordem a identidade e o orgulho com base em valores da família.

⁸Angela Merkel e outros líderes moderados tentam recuperar o patriotismo como conceito democrático e plural. Merkel, na premiação da Fundação Fulbright recebida em Berlim, em 28 de janeiro 2019, disse que o “*Patriotismo significa para mim pensar os próprios interesses sempre conjuntamente com o interesse dos outros. Por isso, nunca vou deixar de clamar pelo reforço da ordem multilateral em torno de valores e regras comuns*”. O líder do partido de centro-esquerda italiana *Partito Democratico* (PD), Nicola Zingaretti, fez um discurso recentemente pedindo pela redescoberta do patriotismo europeu, de forma a se opor ao surto nacionalista (GIUFFRIDA, 2019).

religião, “anti-globalistas” e contra partidos de esquerda. Em um momento no qual as declarações de alto orgulho estão em declínio, mesmo depois de longos anos de ação estatal para reforçar o patriotismo e os símbolos nacionais de forma geral, fornecer uma explicação para as bases desse comportamento parece ser especialmente pertinente para o caso brasileiro.

Assim como Mills chamou atenção nos anos 1950, Calhoun e Wieviorka (2013) afirmam, nos dias atuais, que a sociedade pode se beneficiar das ciências sociais para pensar o mundo atual e transformar crises em debate. Para eles, fenômenos sociais devem ser abordados não somente pela relevância acadêmica, mas também pela sua importância na atualidade e pelo seu interesse público. Mesmo que isso seja feito sem a distância histórica ideal que permite inferências mais robustas. A pretensão exposta nesses termos é, de certa forma, o dever de um acadêmico em formação.

A ESTRUTURA DA TESE

O primeiro capítulo da tese dialoga com as ciências da computação, primeiramente definindo e descrevendo brevemente a história da IA. Em seguida, são apresentados conceitos e técnicas de aprendizado de máquina, além de outros conceitos e técnicas relacionadas como *big-data* e mineração de dados. Uma vez definidos os conceitos, são apresentadas algumas referências do uso prévio dessas técnicas nas ciências sociais e algumas das novas perspectivas na forma da chamada “ciência social computacional” (ALVAREZ, 2016).

Para ilustrar a utilidade da IA para a ciência política, foi escolhido um fenômeno clássico da ciência política, o patriotismo. O segundo capítulo trata desse fenômeno. O capítulo inicia-se com uma breve apresentação do conceito e seu histórico na literatura da ciência política. Em seguida, são expostos o quadro analítico da cultura política e o corpo teórico com o qual dialogam as análises no trabalho e com o qual são observados os dados. Como descrito acima, patriotismo é considerado como parte integrante da cultura política de indivíduos e sociedades, que tem sustentação em um conjunto de valores, crenças e atitudes políticas. A parte final do capítulo é dedicada aos resultados empíricos que podem ser encontrados, principalmente, nas literaturas da cultura política e da psicologia social.

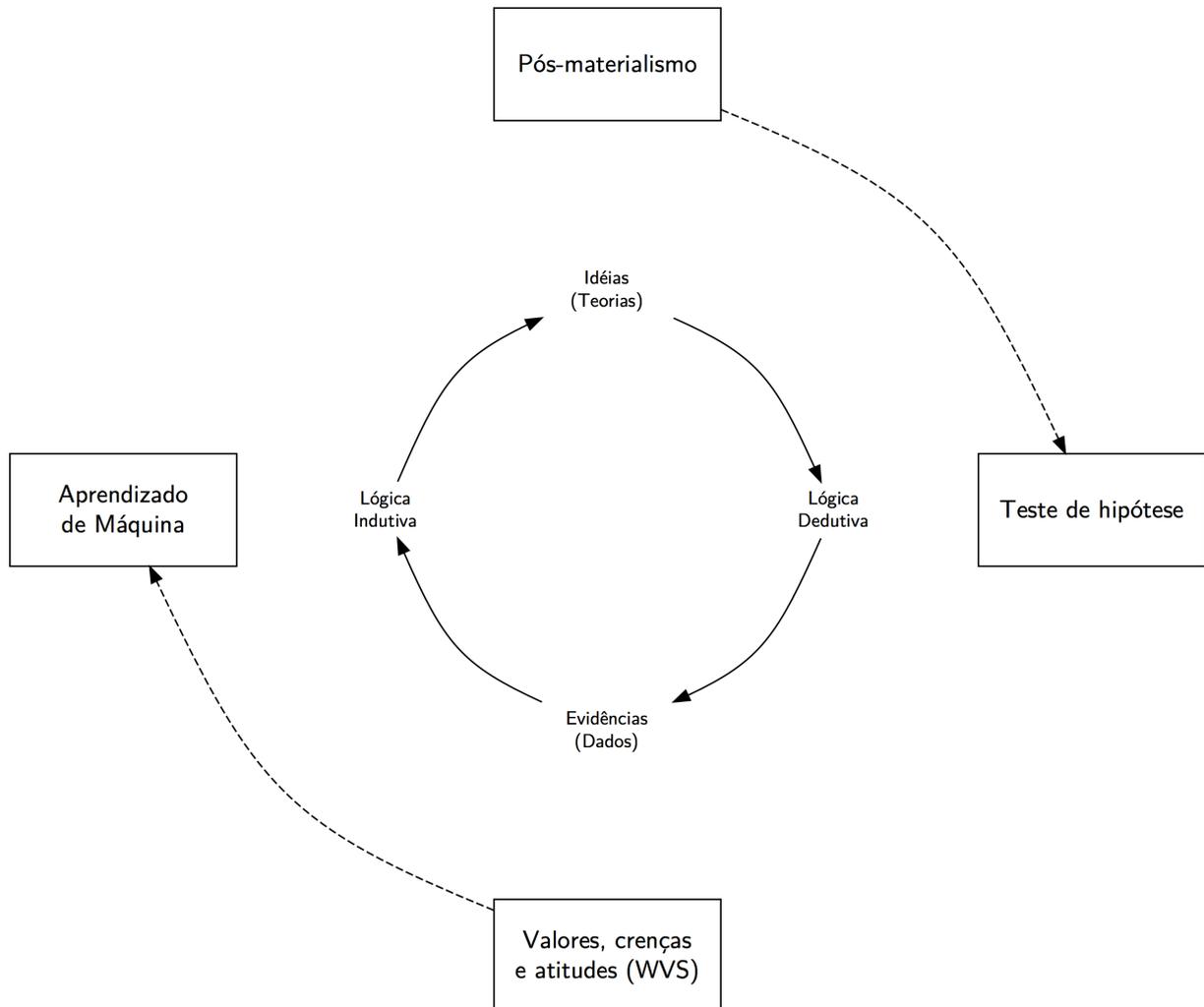
O terceiro capítulo do trabalho apresenta a metodologia que foi aplicada para

experimentalizar o uso do aprendizado de máquina e tentar encontrar uma ou mais hipóteses de explicação para o patriotismo. A intenção não é propor uma nova forma de trabalhar as ciências sociais, mas de integrar novas ferramentas e técnicas no processo de construção do conhecimento. Segundo Ragin e Amoroso (2010), esse processo implica um diálogo constante entre *empíria vs. teoria* ou de entre métodos dedutivos e indutivos. Espera-se que essas novas ferramentas possam ser um novo parceiro do pesquisador, não somente na seleção de variáveis em uma perspectiva quantitativa, mas ajudando-o no processo de síntese e criação de “*imagens*” dos objetos de pesquisa baseada em evidências empíricas.

Para chegar a esse objetivo, o capítulo três descreve os dados da Pesquisa Mundial de Valores (WVS) usados no trabalho e todas as técnicas que foram usadas nos dois momentos dedutivos e indutivos indicados na Figura 2. O momento indutivo é ilustrado no canto inferior esquerdo da figura, onde busca-se construir uma imagem do patriotismo a partir dos dados ou observações da Pesquisa Mundial de Valores. Em termos de modelagem estatística, trata-se de tomar o caminho preditivo com técnicas de aprendizado de máquina e tentar extrair ou inferir a “teoria” identificada pelos algoritmos usando as ferramentas de inferência disponíveis. Esses resultados são usados para dialogar com as teorias que tentam explicar o patriotismo e aprimorá-las. Em um processo contínuo, as novas teorias aprimoradas com a ajuda da IA podem ser testadas em um novo ciclo dedutivo, ilustrado no canto superior direito, usando técnicas de modelagem explicativa. Tendo em vista que o foco são as técnicas de aprendizado de máquina, o tratamento de dados e as principais técnicas aplicadas no aprendizado de máquina são apresentados com mais detalhes.

O capítulo quatro apresenta os resultados. Primeiramente, são apresentados os resultados da modelagem explicativa. Em seguida, são apresentados os resultados do estudo de *benchmarking* que permitiu fazer a escolha de um algoritmo e os resultados do aprendizado. Os resultados em termos preditivos foram comparados entre os modelos explicativos e dedutivos. Uma vez feita a comparação, foram exploradas técnicas que permitem fazer inferências sobre os resultados obtidos na modelagem preditiva. A primeira delas consiste na identificação da importância global das variáveis, e a segunda, na criação de modelos de substituição para obterem-se explicações locais.

Figura 2: Esquema de pesquisa usada no trabalho, com base em Ragin (2010), King et al. (1994). A figura indica o uso do aprendizado de máquina em uma abordagem indutiva (esquerda). Para efeitos de comparação e para demonstrar a complementaridade, também é feito um teste de hipótese clássico (direita).



Nos termos de Ragin (e Mill), a indução não se resume a fazer inferências no sentido da modelagem estatística. Uma boa indução consiste em selecionar as observações mais pertinentes e relevantes, e construir uma imagem do fenômeno que está sendo estudado. Para criar essa imagem, os resultados analíticos foram colocados em perspectiva no capítulo cinco. Neste capítulo, tenta-se selecionar e dar forma às inferências estatísticas de maneira a reconstruir uma “imagem” do patriotismo à luz do conhecimento e das teorias existentes. A “representação” do patriotismo reconstruída, em nível mundial e nacional, é o resultado da interação entre a imagem produzida usando os dados do WVS e a perspectiva analítica da cultura política.

No capítulo cinco, também são discutidos os pontos positivos e negativos da

experimentação com aprendizado de máquina em termos metodológicos. Em particular, abordamos as preocupações distintas da ciência da computação, e da modelagem preditiva em geral, e das ciências sociais entre explicação — o “por que” e o “como” de determinadas fenômenos — e predição — o “o que” e “quando” alguém vai agir de determinada maneira (BOELAERT; OLLION, 2018; WALLACH, 2016). Também é discutida a questão da parcimônia e da abordagem geralmente “*all but the kitchen sink*” da modelagem preditiva em comparação com a tradição das ciências sociais.

Na conclusão, analisa-se o resultado da integração do aprendizado de máquina para o caso do patriotismo e para as ciências sociais de forma mais ampla. São tratadas também as limitações da abordagem com relação aos dados estruturados do WVS e as perspectivas que se abrem como resultado da pesquisa.

1. INTELIGÊNCIA ARTIFICIAL E APRENDIZADO DE MÁQUINA

... a capacidade de assimilar as informações existentes e adaptar-se constantemente às novas informações e situações de modo a estruturar o seu universo...

Jean Piaget, "La Psychologie de l'intelligence"

Assim é a maneira como Piaget (2012, pp. 3-6) define a inteligência humana na sua perspectiva construtivista e cognitivista. Uma capacidade adaptativa essencial do ser humano, que envolve perceber, entender, fazer previsões sobre um mundo.

Entender como pensamos é a motivação dos pesquisadores há séculos, e parte da pesquisa em inteligência artificial deriva desse questionamento. Para Simon (1995), convencionou-se chamar de inteligência artificial os "fenômenos que podem ser observados quando computadores desempenham tarefas que, se desempenhadas por pessoas, seriam vistas como resultado da inteligência ou do raciocínio"¹.

Mas a IA tenta ir além da compreensão da inteligência humana na medida em que tenta substituí-la por inteiro pela construção de artefatos inteligentes (RUSSELL; NORVIG, 2016). As principais definições da inteligência artificial têm origem neste foco na engenharia de máquinas inteligentes, que se inicia, formalmente, na década de 1950, em paralelo ao desenvolvimento dos primeiros computadores, e motivada pela hipótese da máquina pensante de Turing (1950). A definição original e ainda atual foi cunhada por John McCarthy (1955) como "a ciência e engenharia de produzir máquinas inteligentes" no documento de preparação do que se considera ser a primeira conferência em inteligência artificial, a "*Dartmouth Summer Research Project on Artificial Intelligence*", realizada em 1956. A enciclopédia de filosofia de

¹ *the phenomena that appear when computers perform tasks that, if performed by people, would be regarded as requiring intelligence, or thinking.*

Stanford propõe uma definição similar, como sendo o “campo dedicado a construir animais e pessoas artificiais, ou que se pareçam com eles” (BRINGSJORD; GOVINDARAJULU, 2018).

Outras definições são relacionadas às capacidades de uma inteligência artificial. Nilsson (1998), por exemplo, considera IA o comportamento inteligente de artefatos, que envolve as capacidades de percepção, raciocínio, aprendizado, comunicação e ação. Na psicologia e na ciência cognitiva, a IA é definida na direção oposta como uma disciplina que visa entender os seres inteligentes (CUMMINS; POLLOCK, 1992). Essa perspectiva estava presente desde Dartmouth, principalmente com Herbert Simon, já mencionado acima.

Em torno da ideia unificadora de construção ou reprodução de um *agente* inteligente, talvez a classificação mais comum da IA seja entre a inteligência fraca e forte. A IA “forte” se refere a criar máquinas conscientes que possam pensar por si mesmas ou que *tenham* uma mente própria. A IA “fraca” busca criar máquinas que processam informações que *pareçam* ter o mesmo repertório mental dos humanos ou que se comportem de forma inteligente. Na IA forte, o agente é a própria máquina; na IA fraca, os humanos permanecem sempre como agentes². Essa distinção está clara na definição proposta por Intelligence e Council (1997), em que a IA é um “conjunto de cálculos que tornam possível assistir usuários a perceber, raciocinar e agir, e que pode ou não estar sob controle de uma máquina” (*i.e.* computadores ou sistemas robóticos).

A IA também pode ser caracterizada pelas duas abordagens principais que nasceram da conferência em Dartmouth e que podem ser caracterizadas em abordagens de cima para baixo e de baixo para cima. Na primeira, a IA depende do desenvolvimento de sistemas dedutivos seguindo a tese logicista e simbólica, desenvolvida, entre outros, por Allen Newell e Herbert Simon (1995; 1972; 1977). A hipótese central envolve o processamento de símbolos estruturados, em que um conjunto de regras é aplicado, gerando novas estruturas. Os símbolos são as representações de conceitos ou objetos que são processados usando lógica.

Na segunda abordagem, a IA é alcançada a partir de processos indutivos, baseados na probabilidade. Essa abordagem foi desenvolvida, principalmente, na corrente conexionista (no sentido da conexão entre neurônios) da pesquisa em neurocomputação (*e.g.* MINSKY, 1961). A abordagem de baixo para cima pode ser considerada sub-simbólica, mas não

²Não há espaço suficiente para tratar de uma distinção detalhada entre IA forte e fraca. Para uma discussão em diversos aspectos sobre o tema ver, por exemplo, Russell e Norvig (2016), Cap. 26

logicista. Considera a inteligência como emergente e fundamentada nos dados (*grounded*) no meio ambiente. Duas subabordagens podem ser identificadas aqui, uma delas puramente probabilística e outra conexionista. Um exemplo desta segunda são as redes neurais que se interessam pela habilidade em aprender inspiradas em modelos biológicos, tentando mimetizar o que encontramos em nível celular (NEWELL; SIMON, 1976). A premissa básica é que o comportamento inteligente deriva das conexões entre nós, denominadas neurônios em analogia ao sistema natural. Os primeiros algoritmos desse tipo apareceram nos anos 1940-50 (e.g. MCCULLOCH; PITTS, 1943; ROSENBLATT, 1958).

Alguns historiadores da IA identificam períodos de desenvolvimento que corresponderiam a “verões” e “invernos”, ou seja, momentos de euforia e sucesso, seguidos de outros de dificuldades e decepções (e.g. RUSSELL; NORVIG, 2016, p.24). Os anos que se seguiram à conferência de Dartmouth, até o final dos anos 1970, foram anos de entusiasmo e expectativas. Nesse período, apesar de os primeiros movimentos da IA terem sido feitos no conexionismo e em redes neurais, prevaleceu a IA com base em sistemas simbólicos, e foram os sistemas especialistas que conseguiram se destacar³. Esse movimento pode ser ilustrado pelos os sucessos de Newell e Simon (1976) com o *General Problem Solver*. Na década de 1970, a tentativa frustrada de tratar problemas mais complexos, para além da ilusão de que se poderia contar com um poder computacional cada vez maior e que isso seria suficiente para resolvê-los, levou a IA de volta à realidade. O período seguinte consolidou a expansão dos sistemas especialistas, mas também viu o renascimento da esperança nas redes neurais. No final da década de 1980, as promessas não cumpridas levaram a um novo inverno, o mais longo do IA, que durou até o final da primeira década do século XXI. A pesquisa com modelos conexionistas foi retomada nos anos 2010, e, a partir dos progressos computacionais, é a protagonista do novo verão que se vive atualmente (RUSSELL; NORVIG, 2016).

Na perspectiva de construir artefatos inteligentes ou que tentam reproduzir o que a mente humana pode fazer, a operacionalização é geralmente a partir da aplicação de testes de inteligência, sendo o Teste de Turing o mais importante dos testes de habilidade mental⁴,

³Sistemas especialistas, como o nome indica, incorporam o conhecimento de especialistas em termos de fatos e regras e ganham aparência de inteligência manipulando informações na forma de símbolos a partir desse conhecimento (GARSON, 1990).

⁴O “jogo da imitação” foi proposto por Turing (1950) para tentar responder à pergunta “Máquinas podem pensar?”. O objetivo do jogo para um entrevistador é de distinguir um homem e de uma máquina por meio de perguntas e respostas por escrito. Para testar a IA, substitui-se um dos entrevistados pela máquina e testa-se a sua capacidade de imitação.

para além de jogos como o xadrez e o GO (NEWELL, 2015, 1970). Isso explica a publicidade dada aos sucessos obtidos no jogo de xadrez com o *DeepBlue*, que em maio de 1997 bateu Gary Kasparov (ver KASPAROV; GREENGARD, 2017), e do *AlphaZero*, sagrado campeão em 2016 mais de 20 anos depois (SHEN, 2016). O *AlphaZero* é uma evolução do *AlphaGo* e usa métodos de aprendizado profundo (*deep learning*), que é a maneira como foi rebatizada nos anos 2010 a pesquisa em redes neurais profundas. Os sucessos com o aprendizado profundo marcam o fim do último inverno em que se encontrava a IA desde os anos 1990.

Desde então, os resultados desses desenvolvimentos são visíveis em diversos campos. Além de jogar GO ou xadrez melhor que humanos (SILVER et al., 2016), já existem sistemas inteligentes capazes de conduzir veículos autônomos (ENZWEILER, 2015), de extrair padrões de imagens e vídeos (YANG et al., 2015), de reconhecer a fala e a escrita e traduzir textos e áudios, de planificar e organizar agendas, ou, ainda, de detectar fraudes (BOLTON; HAND, 2002), de identificar mensagens indesejáveis (SPAM), de fazer planificação logística. Além disso, existem aplicações específicas para facilitar decisões judiciais (BERK, 2012), prever os efeitos das mudanças climáticas e possíveis ações de mitigação (ATHEY, 2017), identificar doenças (SOUZA et al., 2020) e propor programas de educação sob medida (LUCKIN et al., 2016).

Apesar desses imensos progressos, a perspectiva de construção de uma inteligência próxima da humana parece ainda distante. Um exemplo é o baixo desempenho desses sistemas em testes usados regularmente para testar a inteligência humana, como provas escolares. Metz (2019) descreve um desses casos e atribui a dificuldade das máquinas em tirar boas notas à complexidade de tratar simultaneamente gráficos, textos e o contexto geral das questões⁵. Um bom desempenho em uma prova como essas é tarefa muito mais complexa que os jogos de xadrez ou Go e talvez seja uma referência mais apropriada para avaliar o desempenho de máquinas inteligentes. A tentativa de criação de “cientistas” artificiais também entrou na pauta, mas ainda existem limitações importantes (e.g. Ross D. King et al. (2005); Kotthoff et al. (2019), Steinruecken et al. (2019)).

Uma dessas limitações estaria relacionada à capacidade de aprendizado, uma condição necessária para chegarmos ao desenvolvimento de uma inteligência próxima da humana. Em

⁵O algoritmo mencionado neste estudo teve desempenho similar ao de uma pessoa quando as questões mais complicadas usando gráficos foram retiradas da prova.

referência à definição proposta por Russell e Norvig (2016) e que se aplica particularmente a este trabalho, Payrovnaziri et al. (2020) definiram inteligência artificial diretamente ligada ao aprendizado. Segundo eles, a inteligência artificial pode ser definida como “sistemas que agem de forma humana através do aprendizado de máquina e, mais especificamente, usando análises preditivas”. Esse campo da IA tem visto um desenvolvimento acelerado nos últimos anos e é tratado a seguir.

1.1. APRENDIZADO DE MÁQUINA

Como visto acima, uma inteligência pressupõe uma capacidade de adaptação e aprendizado constantes e de forma autônoma (ALPAYDIN, 2010, 2017). Assim, o aprendizado de máquina (AM, ou ML para *machine learning*), ou ainda aprendizagem automática, é um dos “pré-requisitos” ou “capabilidades” da IA. Isso porque, sem poder aprender, dificilmente poderíamos afirmar que uma máquina fosse inteligente. Devido às suas diversas aplicações, o aprendizado de máquina também é chamado de uma “forma” de IA (BOELAERT; OLLION, 2018), ou considerada como a “parte” mais popular da IA (ALPAYDIN, 2017).

As bases do aprendizado de máquina também podem ser encontradas desde o surgimento das pesquisas em IA, por exemplo, com a ideia do “*Child Programme Idea*”, em que Turing (1950) propunha simular a inteligência de uma criança para deixá-la aprender sozinha em vez de criar um programa para simular a mente adulta.

Na prática, aprendizado de máquina se refere ao processo no qual algoritmos usam os dados disponíveis para criar generalizações baseados em modelos estatísticos. Outra definição do AM na ciência da computação é “a arte de programar computadores para otimizar um critério de desempenho usando dados construídos ou passados” (ALPAYDIN, 2010, 2017). Essa definição é orientada ao objetivo: o aprendizado de máquina se propõe a construir sistemas que melhorem sua performance em tarefas nas quais lhe são fornecidos exemplos de performance ideais para a tarefa em questão, ou então que melhorem sua performance pela repetição da experiência na tarefa (BRINGSJORD; GOVINDARAJULU, 2018; RUSSELL; NORVIG, 2016; MITCHELL, Tom M., 1997).

Apesar de os sistemas lógicos também aprenderem e os sistemas especialistas terem usado aprendizado automático para algumas de suas previsões, o AM se desenvolveu de forma mais importante dentro da abordagem de baixo para cima da IA. Nessa abordagem, o agente não é programado *a priori*, tanto pelas dificuldades de antecipar todas as situações possíveis, quanto pela impossibilidade de antecipar as mudanças ao longo do tempo. Pelo contrário, o aprendizado se refere à capacidade de detectar e extrapolar/generalizar padrões a partir de uma série de dados, sem que seja necessário usar regras predefinidas. Apesar disso, os simbolistas ainda contribuem bastante para a área, juntamente com os conexionistas, os evolucionistas, os Bayesianos e os analogistas, cada grupo com seus “algoritmos-mestres” (DOMINGOS, 2015).

Segundo Russell e Norvig (2016), o aprendizado de máquina pode ser classificado em função do retorno (*feedback*) disponível durante o aprendizado⁶. Os três tipos de aprendizado de máquina mais comuns correspondem aos três tipos de retorno possíveis. No caso do aprendizado supervisionado o agente observa a relação entre os dados de entrada e saída e tenta aproximar uma função que a descreva. Essa função é a melhor hipótese para a ocorrência da variável resposta e usada para prever sua ocorrência. O retorno que o algoritmo recebe enquanto “treina” diferentes explicações usando as variáveis explicativas é a própria variável resposta que tem um valor definido⁷.

No caso do aprendizado não supervisionado, o algoritmo não dispõe de nenhum retorno explícito. O aprendizado consiste em encontrar padrões e tentar categorizar os dados em função de conceitos genéricos de agrupamento. Ao contrário dos anteriores, essa categoria de algoritmos não necessita que seja identificada uma variável resposta que possa *supervisionar* as iterações entre variáveis explicativas. Enfim, o aprendizado reforçado se refere à capacidade de um agente mecânico inteligente de aprender a como evoluir em um ambiente, tomando decisões com base em seus erros e acertos. Feitas as previsões pelo modelo, os erros são punidos e os acertos recebem um retorno na forma de recompensas que reforçam o aprendizado feito anteriormente.

No caso deste trabalho, foi aplicado o aprendizado de máquina supervisionado. Para esses casos, pode-se fazer também uma distinção em função do tipo de problema a ser resolvido

⁶Segundo os autores, o tipo de aprendizado também depende do componente do agente que se tenta otimizar, do conhecimento e da representação dos dados disponíveis.

⁷Na literatura, as variáveis dependentes, ou resposta, e independentes, ou explicativas, podem ser chamadas, respectivamente, de “output”, “target”, “class” e “label”; e de “inputs”, “features” ou “attributes”.

(RUSSELL; NORVIG, 2016). Para os casos em que a função a ser descoberta tem é uma série limitada de valores, trata-se de um problema de *classificação*. Quando o retorno é um valor contínuo, trata-se de resolver um problema de *regressão*. O trabalho apresentado a seguir trata de resolver um problema de classificação de pessoas em função das declarações de orgulho da nacionalidade⁸

Antes de prosseguir, é necessário também definir os conceitos de algoritmo e modelo que são usados de forma recorrente no trabalho. Algoritmos são entendidos, comumente, como um conjunto de instruções usadas para transformar os dados de entrada em dados de saída e, assim, resolver um problema computacional (INTELLIGENCE; COUNCIL, 1997; ALPAYDIN, 2017). Um *algoritmo de aprendizado de máquina* se refere a um algoritmo que usa um outro algoritmo para modificá-lo durante o processo de aprendizado⁹. Este outro algoritmo, que é gerado durante o treinamento com os dados, é um modelo ajustado aos dados que é chamado aqui simplesmente de modelo¹⁰.

Como veremos com mais detalhes no capítulo seguinte, o aprendizado se trata do processo no qual se tentam otimizar os hiperparâmetros do primeiro e são selecionados os melhores parâmetros do segundo, em função do melhor ajuste aos dados disponíveis. Hiperparâmetros são as configurações externas dos algoritmos de aprendizado de máquina e que não podem ser aprendidos durante o treinamento com os dados. Parâmetros se referem ao ajuste dos modelos.

1.2. APLICAÇÕES DA IA NAS CIÊNCIAS SOCIAIS

A tentativa de integração da IA nas ciências sociais não é inédita. A cada novo momento de excitação em torno da IA, renovam-se as expectativas e surgem novas formas de combinar os métodos e as técnicas de pesquisa. Como visto acima, Herbert Simon e Allen Newell, precursores da IA, interessaram-se pela aplicação na psicologia cognitiva desde seus primórdios. Nos anos 1970, um dos focos foi a simulação de funções cognitivas individuais

⁸Para os modelos não supervisionados, em que não existe um retorno, trata-se de estimar densidades ou identificar sub-conjuntos (e.g. *clustering*, *vector quantization*).

⁹Na literatura, ainda são chamados de *learners*, ou estimadores.

¹⁰Vale fazer um pequeno comentário aqui: nas ciências naturais, “modelo” se refere, normalmente, a uma dedução lógica sobre algum fenômeno que pode ser testada procurando ou simulando dados. Um *fit*, ou ajuste, não se constitui em um modelo nesses termos, mas apenas em uma descrição dos dados.

e a simulação de grupos sociais. Nesse mesmo período, as técnicas de categorização, ou clusterização de informações, começaram a chamar a atenção.

Já nos anos 1990, Bainbridge et al. (1994, p. 408) propuseram um campo específico de pesquisa na sociologia a que chamaram de “Inteligência Artificial Social”, definida como “a aplicação de técnicas de inteligência mecânica a fenômenos sociais, [...] incluindo construção de teoria e análise de dados”¹¹. Bainbridge e seus colaboradores citam diversas experiências, entre elas a de Kimber (1991), que testou um algoritmo de redes neurais artificiais usando dados eleitorais, e Chablo (1996), que, na Antropologia, usou o ID3 (Iterative Dichotomiser 3), precursor do algoritmo C5.0 baseado em árvores de decisão. Também no final dos anos 1990, Brent (1989) propôs o uso de sistemas especialistas para auxiliar no desenho de pesquisas em ciências sociais.

A capacidade de computação melhorou muito desde as primeiras tentativas e fez com que as abordagens ficassem mais acessíveis. Em termos de processamento, a utilização das Unidades de Processamento Gráfico (GPU) ou de processadores especiais para redes neurais, além de algoritmos mais modernos acelerou a obtenção dos resultados. Para além das bases de dados sociais, cada vez mais em livre acesso, existem novas fontes de dados em “volume, variedades e velocidades” muito acima do que existiam no passado, os chamados “*big-data*” (FAVARETTO et al., 2020; OLLION; BOELAERT, 2015).

Na esteira desses desenvolvimentos, diversos usos recentes de técnicas de AM podem ser identificados. Mihaela e Stefan Robila (2019) identificaram diversas aplicações nas ciências do comportamento e sociais, incluindo o diagnóstico de algumas condições, predição, desenvolvimento psicológico e compreensão geral do comportamento. Na ciência política, mais especificamente, podem ser identificados por exemplo o uso disseminado em análises de texto e a multiplicação de métodos Bayesianos, muito usuais em algoritmos de AM. Existem também aplicações em dados eleitorais ou para identificar padrões de colaboração entre membros de uma mesma comunidade a partir de uma *survey*. Apesar do foco ser quase sempre em análises quantitativas, as aplicações de AM também existem para uso em dados qualitativos.¹²

¹¹ *Broadly defined, Artificial Social Intelligence (ASI) is the application of machine intelligence techniques to social phenomena. ASI includes both theory building and data analysis.*

¹² Para análise de texto ver, por exemplo, King, Pan e Roberts (2013) e Grimmer e Stewart (2013). Para métodos Bayesianos, ver Clinton, Jackman e Rivers (2004). Para aplicações com dados eleitorais, ver, por

Também já existem análises em cultura política usando os mesmos dados estruturados de *survey* do WVS que são usados no presente estudo (NASCIMENTO; BARONE; CASTRO, 2019). Muitas dessas técnicas de AM foram incorporadas nas perspectivas e ambições do que alguns autores chamam de Ciências Sociais Computacionais (LAZER et al., 2009).

Mais adiante, será apresentada a proposta metodológica, mas antes é preciso apresentar o estudo de caso ou a pergunta das ciências sociais que foi usada para testar as técnicas de aprendizado de máquina neste trabalho.

exemplo, Levin, Pomares e Alvarez (2016), Cantu e Saiegh (2010) e Beauchamp (2017). Settles e Dow (2013) aplicam aprendizado de máquina em dados de *survey* e o trabalho de Nan-Chen Chen et al. (2018) é um exemplo de aplicação em pesquisa qualitativa.

2. PATRIOTISMO

No capítulo anterior, o problema central da tese foi apresentado e discutido. Foram introduzidas as definições de inteligência artificial e aprendizado de máquina e formulada a hipótese exploratória de que as ferramentas deste campo da ciência da computação podem ser integradas no processo de construção do conhecimento nas ciências sociais. Espera-se que o aprendizado automático sirva para aprimorar ou construir novas teorias passíveis de teste empírico em um segundo momento hipotético-dedutivo.

O objetivo deste capítulo é explicitar o fenômeno escolhido das ciências sociais que serve como caso de estudo, o patriotismo. Patriotismo é um conceito disputado, podendo significar diversas coisas, como, por exemplo, uma forma de apego ao grupo social ou país de residência (BAR-TAL; STAUB, 1997), uma atitude (KOSTERMAN; FESHACH, 1989), uma crença (DOOB, 1976), um valor (INGLEHART, R., 1971), uma doutrina (KEDOURIE, 1960), uma ideologia (se entendida como sinônimo de nacionalismo), um movimento (BERNARDES, 2006), um princípio (WILSON, 1899), uma linguagem política (BRUBAKER, 2004), uma religião cívica (ROUSSEAU, 2003). Em todos os casos, pode ter conotações positivas ou negativas. Essa característica dificulta enormemente a delimitação do objeto de estudo, fazendo com que qualquer revisão seja necessariamente parcial.

Sem pretensões de ser exaustivo sobre um tema complexo e com uma literatura muito extensa, trata-se de apresentar alguns elementos tanto em termos teóricos quanto de pesquisa empírica usando *survey* que serão úteis para a construção da tese. O que parece ter utilidade é dispor de uma ou mais “representações”, no sentido proposto por Ragin do diálogo entre ideias e provas empíricas, do patriotismo no mundo e no Brasil que possam servir de guia na interpretação dos resultados do aprendizado automático. Intuitivamente, espera-se que essas técnicas permitam identificar diversas teses de forma simultânea e, assim, que sejam desconsideradas as disputas que envolvem as diferentes explicações possíveis para

um determinado fenômeno.

É isso que se faz a seguir. A primeira parte do capítulo apresenta o conceito de patriotismo em termos históricos. São apresentadas algumas das ideias e teorias que foram desenvolvidas ao longo da história e os fenômenos sociais que motivaram o estudo do tema no mundo e no Brasil. A segunda parte do capítulo apresenta o conceito operacional que foi usado e revisa os resultados de pesquisas empíricas sobre o tema. Antes de apresentar os resultados empíricos, propriamente ditos, descreve-se, brevemente, o quadro analítico compartilhado pela cultura política e pela psicologia social, de forma a explicitar a natureza das provas empíricas e como são observadas em pesquisas tipo *survey*. Ao final do capítulo, são apresentados os achados empíricos sobre o caso brasileiro e uma análise das lacunas que podem ser preenchidas por este estudo.

2.1. DEFINIÇÃO HISTÓRICA DE PATRIOTISMO

Na sua origem, patriota e compatriota eram termos usados para designar indivíduos que tinham ancestrais em um mesmo país ou território mais ou menos delimitado, pertencendo a um mesmo clã.¹ Entre o final do século XVII e meados do século XVIII, o conceito de patriotismo foi ampliado e se mantém praticamente inalterado até hoje.

Em sua forma mais simplificada, patriotismo significa um sentimento de afeição especial, ou de amor, ao país que um indivíduo reconhece como a sua pátria. Com base nos escritos de Nathanson (1993), Primoratz (2017) descreve esta ampliação do conceito na Enciclopédia de Filosofia de Stanford. Para ele, patriotismo não é somente o amor ao seu país, mas também “uma identificação pessoal com ele” e uma “preocupação especial pela prosperidade do país e pelo bem-estar dos seus compatriotas, o que implica a possibilidade de se sacrificar por ele”².

Apesar das poucas mudanças na sua definição, o patriotismo tomou diversas formas ao longo da história devido à dissolução, ao surgimento e à junção de novas “pátrias” que

¹ Isso fica claro na origem etimológica dos termos latino *patriota* e grego antigo *patris* que significam “do mesmo país”. O radical pátria, do latim *pater* e do grego *patris*, designa o país dos pais, a “terra natal” ou “terra paterna”.

² ... *love of one's country, identification with it, and special concern for its well-being and that of compatriots.*

encarnaram, sucessivamente, esse poder e cultura. Gregos e Romanos tinham amor tanto às pólis quanto à própria civilização grega. Durante a Idade Média, a pátria era personificada nos príncipes. A formação das cidades na Baixa Idade Média fez com que o patriotismo fosse também dirigido à liberdade que essas aglomerações lhes proporcionavam (VIROLI, 1997). Com a paz de Westfalia em 1648, a pátria se tornou menos amorfa, com um território mais claro e delimitado e o poder centralizado no príncipe (HOBSBAWM, 1992).

A partir do início do século XVIII, com as revoluções na Inglaterra, nos Estados Unidos da América e na França, pode-se notar uma nova transformação do conceito de pátria. A ordem aristocrática anterior, hierárquica e determinista, foi quebrada, criando um novo referencial identitário abstrato — o estado-nação —. A partir desse momento, terreno está pronto para que o patriotismo se transforme em nacionalismo, servindo de guia para a história desde então (SMITH, 1998).

Agora que tinham sido criadas as nações, deveriam ser criados os nacionais. Como disse Massimo d'Azeglio: "*We have made Italy, now we have to make Italians*" (citado em HOBSBAWM, 1992, p. 44). Os governos destes novos estados-nação tiveram papel fundamental na socialização nacional e na construção de um discurso hegemônico que fomentou a identificação nacional e legitimou as novas nações. Além das próprias guerras, os estados passaram a impor uma língua nacional, sistemas educacionais, serviço militar etc. Além disso, o estado passa a ser o grande fiador de uma cultura nacional imaginada pelos cidadãos (GELLNER, 2008; HOBSBAWM, 1992; TILLY, 1990; ANDERSON, 1991)

A partir desse momento histórico, surgem patriotismos de conteúdos variados. Para Tocqueville (2001), por exemplo, surge um patriotismo reflexivo que substitui a crença quase religiosa anterior e se desenvolve com a ajuda das leis e com o exercício dos direitos. O cidadão percebe a utilidade da nação, e do estado, e se dá conta de que defender os seus interesses significa defender os seus próprios interesses³.

Outro patriotismo de origem étnica e cultural se desenvolve em paralelo nas nações que se formam alguns anos mais tarde como a Alemanha e a Rússia. Essas nações, cujas elites se ressentiram pela perda de status, foram buscar outras fontes de identidade (GREENFELD, 1992). Com o conceito romântico de *Volk*, Herder (2002) identificou o patriotismo não como

³ "*de cet intérêt naît un patriotisme réfléchi où le citoyen se rend compte que défendre les intérêts de la Nation c'est aussi défendre les siens*".

uma virtude política, mas sim como uma ligação espiritual com a nação. Em sua perspectiva, em lugar de preservar a liberdade republicana, patriotas deveriam tratar de preservar uma cultura comum e a unidade espiritual do seu povo. Colocada nesses termos, a terra dos pais (*Vaterland*) é pré-política e baseada em traços étnicos e linguísticos.

Essas duas formas de entendimento do patriotismo permanecem até o século XX. Em especial, o patriotismo étnico-cultural foi uma das explicações dadas à forma negativa que o amor à pátria tomou com o nazismo e o fascismo. Esse modelo já havia sido denunciado por Renan, que considerava que “o chamado à fusão dos valores políticos da pátria em uma nação espiritual, dando prioridade à cultura, raça ou língua era uma degradação moral e intelectual” (citado em VIROLI, 1997). Mas foi quando foram mobilizadas milhares de pessoas usando uma forma de patriotismo que colocava a sua pátria e identidade étnica (no caso do nazismo) acima das demais, que este “mau” patriotismo, ou *pseudopatriotismo* nas palavras de Adorno et al. (1950), adquiriu conotações realmente negativas. Hobsbawm (1992) chama de “paradoxo trágico”, que a politização do patriotismo tenha criado as condições para o surgimento das paixões nacionalistas que se sucederam.

Podem ser entendidas nessa perspectiva as tentativas de resgatar as concepções positivas do patriotismo do final dos anos 1990. Um desses patriotismos “bons” se desenvolve na tradição republicana clássica e é baseado nos princípios de liberdade, cidadania ativa e autossacrifício pelo bem comum (VIROLI, 1997). Outra forma positiva de patriotismo propõe uma resposta ao fato de as sociedades se desenvolverem cada vez mais diversas culturalmente, em um contexto no qual a identidade nacional é cada vez menos consensual. Nesses casos, a identidade política só poderia ser direcionada às leis e às instituições democráticas na forma de um patriotismo constitucional pós-nacional (HABERMAS, 1992; MÜLLER, 2009). Uma terceira forma positiva de patriotismo parte da premissa de que, encontrando-se valores e ideais similares em um outro país, um indivíduo poderia se orgulhar de outras pátrias dando origem a um patriotismo cosmopolita baseado no orgulho de valores universais (APPIAH, 1998; NUSSBAUM; COHEN, 2002; NUSSBAUM, 1994).

Tal como se convencionou chamar, nacionalismo é o patriotismo com conotação negativa, em geral ligado a uma identificação étnico-cultural e quando se acredita que a sua nação é superior a outra. Tendo em vista o entendimento de que a identificação é condição necessária para a estabilidade das nações, tenta-se encontrar um patriotismo positivo e virtuoso,

seja na forma republicana (Viroli), de engajamento em uma disputa por direitos em troca de deveres (Mill, Tocqueville, Habermas) ou cosmopolita (Appiah, Nussbaum). Contudo, essa distinção parece difícil. Nas palavras de Canovan (2000), seguidamente o bom patriotismo é o nosso, enquanto que o mau nacionalismo é o dos outros.

2.1.1. O patriotismo no Brasil

São muitas as descrições da identidade nacional dos brasileiros e suas origens. Podemos citar os estudos produzidos por Skidmore (1999), Chilcote (1969), Burns (1968), Rowland (2003), ou ainda Jancsó (2003), Lessa (2008), Pamplona (2011), Oliveira (1990) e as várias obras de José Murilo de Carvalho (1982, 1998, 1974, 1999, 2001, 2008, 1992, 2017).

Não se trata aqui de propor uma nova narrativa, e mesmo fazer uma síntese desses estudos é tarefa difícil. Como colocado mais acima, o que parece útil é tentar extrair desses diversos estudos uma ou mais representações do patriotismo do brasileiro de forma a ajudar na interpretação dos dados (*i.e.* nas imagens que serão obtidas com o aprendizado de máquina).

Um primeiro aspecto dessa representação diz respeito à origem da nação propriamente dita. O patriotismo da América do Norte e da Europa foi transplantado de diversas formas. Na maior parte da América Latina, o movimento revolucionário republicano chegou inspirado pela Declaração da Independência dos Estados Unidos da América e pela Declaração dos Direitos do Homem e do Cidadão na França. Os dissensos quanto ao modelo de contrato social “vindo de cima” que foram sentidos nos países de origem foram sentidos com mais força nas colônias. Como Anderson (1991) descreve em sua obra, o nacionalismo nas colônias foi tipicamente uma resposta por parte dos grupos dinásticos e aristocráticos — classes superiores — aos tipos de nacionalismo mais vernaculares⁴. Apesar da resistência das elites locais, os valores das revoluções liberais foram incorporados de forma variada. A queda do Império Espanhol e as diferentes guerras e conflitos que a sucederam geraram a pulverização do território colonial em várias nações independentes (PAMPLONA, 2011).⁵

⁴...was typically a response on the part of threatened dynastic and aristocratic groups - upper classes - to popular vernacular nationalism.

⁵Não há espaço aqui para descrever de forma detalhada a formação da identidade nacional na América hispânica e suas diferenças com o Brasil. Mas, apesar das reservas, o próprio Anderson considera os movimentos na América espanhola como independências reais e, de certa forma, mais alinhados aos valores das revoluções na Europa (pp. 49-51). Ver também Pamplona, que descreve o processo de descentralização que se seguiu à Constituição de 1812, determinado pelo nível de consolidação da administração colonial.

No Brasil, a era das Revoluções teve efeitos distintos. Primeiro, tentou-se em vão retomar um sistema dinástico com a chegada da família Real em 1808. Durante o Império, uma forma de patriotismo liberal persistiu nos movimentos separatistas e independentistas e ameaçou desintegrar o país (SKIDMORE, 1999). Mas esses movimentos foram contidos pela elite governante da época, homogênea e ideologicamente e ainda muito ligada ao seu passado colonial (CARVALHO, 1998, 1982).⁶

Nesse período, de qualquer forma, não havia no Brasil um receptáculo para ideias republicanas e liberais. Brasileiros, brasilienses e brasileiros, escravos indígenas e negros com educação limitada, conviviam sem que uma identidade comum existisse e sem educação (PAMPLONA, 2011). Até mesmo as elites nacionais tinham poucas oportunidades educacionais e para os demais não havia nada que pudesse construir uma identidade. O alcance da imprensa era muito limitado e livros eram inacessíveis para a população, e os seus escritos serviam mais para reforçar as visões das elites do que convencer os demais (CÂNDIDO, 1968)⁷.

Além disso, as transformações da estrutura econômica que ocorreram na Europa não tiveram lugar no Brasil. A estrutura oligárquica e patriarcal nos engenhos e nas minas foi mantida mesmo depois da independência (PRADO JR, 1953). Também não havia a possibilidade de construir-se uma identidade devido ao acesso muito limitado entre as diferentes regiões do país (ALENCASTRO, 2000; CARVALHO, 2001). Pará, Maranhão, Piauí e Ceará, que tinham ligação mais forte com a Coroa, inclusive se mantiveram fiéis logo no início do processo de independência (SCHWARCZ et al., 2011).

Se a narrativa do período anterior à independência era elitista e celebrava os valores morais, religiosos e políticos vindos de fora, durante o Império teve que ser criada uma imagem de nação brasileira que pudesse unir sob uma mesma pátria cariocas, paulistas, mineiros, gaúchos, maranhenses, pernambucanos (OLIVEIRA, 1990). A primeira história oficial do Brasil, a *História geral do Brasil* (1854-1857), escrita por Varnhagen (1857), foi construída

⁶Entre 1831 e 1848, mais de 20 revoltas menores e 7 maiores aconteceram em diversas partes do país. Em três delas os líderes proclamaram governos republicanos independentes. Ver Ricci (2007); Bernardes (2006); Jancsó (2003). Segundo Skidmore (1999), esses movimentos podem ser retrçados até meados do século XVIII e revelam a influência das ideias iluministas no Brasil.

⁷Cândido cita: *O Peregrino da América* de Marques Pereira (1728), a *História da América Portuguesa*, de Rocha Pita (1730); os poemas *Uruguai*, de Basílio da Gama (1769), *Vila Rica*, de Claudio Manuel da Costa (escrito antes de 1776), e *Caramuru*, escrito por Durão (1781).

em cima da natureza exuberante e do mito da miscigenação virtuosa das três raças tentando legitimar o passado de segregação racial e conter a hostilidade aos portugueses, talvez por receio de que uma revolta de origem racial despertasse como no Haiti⁸. Foi recuperada a visão edênica do Brasil que sempre esteve presente desde o descobrimento (HOLANDA, S. B. de, 2000; CARVALHO, 1998). Já na República, a história se repetiu, e a narrativa foi reescrita de modo a não mais defender o sistema monárquico e, mas novamente dialogando com as narrativas ufanistas da nação (OLIVEIRA, 1990).

Além de excesso de ufanismo, a construção da identidade nacional tratou seu povo de forma peculiar. De acordo com Carvalho (1999, p. 233), a primeira visão de nação foi marcada pela ausência de um povo, e daí a necessidade de recorrer ao mito da miscigenação das três raças. Como visto acima, parece não ser uma característica original, já que a criação de nacionais também foi posterior na Europa. A segunda e a terceira visão, ainda segundo Murilo de Carvalho, parecem ser mais peculiares ao caso brasileiro. A segunda visão do povo é negativa, marcada pela necessidade de branqueamento. Na terceira, desenvolvida já no século XX, apesar de ter um papel importante na propaganda nacional, o povo era visto de maneira extremamente paternalista. A partir desse período, com maiores ou menores intensidades, o estado se colocou como o grande patriarca que “cria a nação, seja ela rural ou urbana, agrícola ou industrial, e que garante que os valores da grande família nacional permaneçam inalterados sob sua autoridade” (CARVALHO, 1991).

Não tendo raízes próprias, o patriotismo que acabou se criando no Brasil é decorrente de uma série de construções elaboradas pelas elites a cada etapa da evolução da nação e do estado brasileiro. A cada etapa, a identidade foi reestruturada e reconstruída de forma a substituir a anterior e ser usada na socialização das gerações seguintes (OLIVEIRA, 1990). Tendo em vista o caráter ufanista e impositivo das visões nacionais, e o papel de coadjuvante que o povo e sua cultura tinham na elaboração dos projetos de país, sempre vistos como inferiores com relação aos seus pares europeus ou norte-americanos, a pátria oficial que se tentou construir teve sempre um caráter muito artificial⁹.

⁸“Em geral busquei inspirações de patriotismo sem ser no ódio a portugueses, ou à estrangeira Europa, que nos beneficia com ilustração; tratei de por um dique a tanta declamação e servilismo à democracia; e procurei ir disciplinando idéias soltas de nacionalidade” (citado em GUIMARÃES, 1988).

⁹Sérgio de Holanda (1999) falou da “fabricação de uma realidade artificiosa e livresca, onde nossa vida verdadeira morria asfixiada”. Nas palavras de Ramos (1960), foi a construção de uma cultura “amorfa, alienada e inautêntica”.

Neste contexto, vale lembrar o período mais recente da história brasileira, em que o Estado tentou criar “brasileiros” e desenvolver o patriotismo de forma autoritária. Primeiramente durante o Estado Novo de Vargas, que teve como uma das características a propaganda oficial do Departamento de Imprensa e Propaganda (DIP)¹⁰; e depois resignificada durante o regime militar entre 1964 e 1985 sob aparência democrática, quando o Estado tentou, novamente, impor uma visão nacional por meio, por exemplo, da Educação Moral e Cívica (EMC) nas escolas e dos diversos dispositivos usados para controlar as políticas culturais. É essa herança, mais recente, que pode ser encontrada mais facilmente nos dados usados na pesquisa.

2.2. A PESQUISA EMPÍRICA SOBRE PATRIOTISMO

Patriotismo é estudados em diversas disciplinas das ciências sociais. História, sociologia, ciência política, cada uma dessas disciplinas tem sua própria abordagem metodológica, desde a elaboração de teorias até a coleta de dados, abarcando também técnicas de análises próprias.

As explicações na área da ciência política, ou as representações da sociedade que são elaboradas sob a ótica das relações de poder, são, em geral, organizadas por um quadro analítico que ora privilegia evidências relacionadas às instituições, ora relacionadas à cultura. No primeiro caso, buscam-se explicações sobre como emergiram e foram impostas instituições formais e informais, se são mantidas pela força ou não, sua qualidade etc. (BAQUERO; CASTRO; RANINCHESKI, 2016; BAQUERO, 2018; BAQUERO; GONZALEZ, 2016)

No segundo caso, as relações de poder e as próprias instituições emergem como uma consequência de ações descentralizadas e que podem ser explicadas por normas sociais, valores e crenças individuais etc. O patriotismo, visto por essas lentes, é uma característica da cultura política. De forma resumida, a explicação para os fenômenos nesse ramo da ciência política tem como base ideias, teorias e hipóteses formuladas em torno do que pensam e como agem as pessoas em relação às instituições políticas e na maneira como estas evoluem como uma causa e efeito da sua cultura política (CASTRO, 2008; MOISÉS, 2005; RENNÓ, 1998).

De forma um pouco mais precisa, a cultura política pode ser definida como um

¹⁰A revista “Cultura Política: Revista Mensal de Estudos Brasileiros” foi um dos produtos estratégicos do DIP.

subconjunto dos fenômenos culturais, moldado por crenças, normas e valores políticos. Crenças são definidas como o que as pessoas pensam dos *fatos* como certos ou errados, valores são o que as pessoas pensam que é *moralmente* bom ou mau, e normas são as orientações de comportamento que, ao contrário dos valores, são sancionadas pela sociedade, formal ou informalmente, estejam elas internalizadas ou não (WELZEL; INGLEHART, 2010b). Valores diferem das atitudes, pois são em menor número e mais resilientes¹¹. Uma vez internalizados, não são sancionados como as normas e operam relativamente independentes das instituições formais (INGLEHART; WELZEL, 2005). O repertório conceitual da cultura política inclui ainda o conceito de atitudes e opinião¹².

Apesar das muitas críticas, a cultura política permanece uma forma bastante aceita de reconstruir ou interpretar a realidade social e explicar o comportamento político. O faz com base na experiência histórica das sociedades e nos padrões de crenças, normas e valores das suas populações, considerando-os como componentes endógenos da tomada de decisão (CASTRO; CAPISTRANO, 2008; BORBA, 2005)¹³.

A escolha da cultura política como quadro analítico implica uma forma própria de observar a realidade e procurar evidências. As consequências das teorias e das hipóteses com base em valores, crenças, atitudes, ou opiniões individuais, podem ser observadas tanto no comportamento humano individualizado quanto nos padrões nos quais são encontradas em uma determinada sociedade. Na psicologia, a principal fonte de dados sobre comportamento são os experimentos em laboratório. Na psicologia social e na cultura política, uma das formas mais comuns de coletar evidências é a partir de *surveys* de opinião pública, que permitem a análise das *declarações* que os indivíduos fazem sobre as diversas questões que lhes são apresentadas. As declarações são consideradas reações aos estímulos (questões do *survey*); a parte observável do comportamento dos indivíduos. A partir delas, pode-se fazer análises de atitudes, crenças e valores. Pelo desenho das *survey* , que amostram uma parcela representativa

¹¹ "...values differ operationally from attitudes only in being fewer in number, more general, central and pervasive, less situation-bound, more resistant to modification and perhaps tied to developmentally more primitive or dramatic experiences" (ROBINSON; SHAVER; WRIGHTSMAN, 1991; citado em INGLEHART, R., 1977, p. 29).

¹²Para mais detalhes: ver a produção de Ronald Inglehart. Para opinião: ver Zaller (1992). Para uma análise da ligação entre estes conceitos, ver por exemplo Bergman (1998).

¹³Uma revisão completa da área de estudo está fora do escopo do presente estudo. Ver, por exemplo, Castro (2014) sobre a armadilha das tipologias fixadas à priori; Carole Pateman (1971) sobre o determinismo étnico da democracia na origem dos estudos em cultura política; Brian Barry (1988), James Bill e Robert Hardgrave (1973) e Bertrand Badie (1993) pela inabilidade de explicar mudanças de regime, entre outras críticas mais importantes.

de uma determinada população, essas declarações podem ser agregadas para uma análise do patriotismo em nível societal¹⁴.

Patriotismo, na definição apresentada acima, é uma forma de amor, portanto de sentimento, em direção à pátria. Mas a noção de amar alguém, ou algo como o seu país, antes de ser uma reação emocional ou sentimento episódico, pode ser considerada também como uma disposição afetiva de longo-prazo (SCHERER, 2005). Nesses termos, o patriotismo se aproxima de uma atitude tal como é definida na psicologia. Atitudes são estudadas em nível individual e em relação a um objeto específico. Autores da sociologia e da cultura política, com os quais esta tese se alinha, postulam que indivíduos raramente formam atitudes no vácuo. As atitudes estão sujeitas a forças sociais importantes e são adquiridas juntamente com valores e crenças durante as diferentes etapas da socialização. Nessa linha, o patriotismo surge atrelado a valores e crenças mais profundamente enraizados e compartilhados por aqueles que tiveram uma socialização comum.

A operacionalização mais comum do patriotismo em *surveys* são as declarações de orgulho da nacionalidade. Essa declaração de orgulho seria o comportamento visível do patriotismo, seja ele motivado por uma atitude, crença ou valor. Mesmo que alguns autores advoguem por uma medida multidimensional de orgulho que permita identificar de forma mais clara quais os fatores (estímulos) que fazem com que um indivíduo sinta orgulho, uma medida simples de orgulho parece ser suficientemente consistente para medir o patriotismo de forma geral (TILLEY; HEATH, 2007; DE FIGUEIREDO; ELKINS, 2003; MOADDEL; TESSLER; INGLEHART, 2008; MICHELAT; THOMAS, 1966).

Diversas pesquisas empíricas foram realizadas sobre patriotismo, nas disciplinas de sociologia, psicologia, ciência política e economia. A possível operacionalização do patriotismo em uma única pergunta tanto para identificar patriotismo como uma atitude quanto como crença ou valor amplia o leque de possíveis explicações do fenômeno. Assim, parece interessante ir além das teses que deram origem aos dados que foram usados no trabalho.

Assim, os trabalhos revisados a seguir têm dois quadros teóricos principais: o da cultura política e o da psicologia social. Foram revisados tanto os trabalhos empíricos voltados à

¹⁴Para o problema da falácia ecológica ver o debate entre Inglehart e Welzel (2003) e Seligson (2002).

explicação de atitudes, que têm como base questionários mais amplos em termos de estímulos ao orgulho (como é o caso do ISS), quanto os mais específicos da Cultura Política que se preocupam com a explicação de crenças e valores. Além disso, a abordagem da cultura política, apesar de propor um quadro que tenta unificar análises micros e macros, de forma geral foca na análise de dados agregados em nível societal, ou de culturas. A psicologia social, por sua vez, tenta explicar as influências culturais na formação das atitudes individuais, em que a ênfase da análise política em nível societal é secundária (Ver a tabela 10 no Apêndice C para mais detalhes sobre cada um dos estudos levantados).

O campo da psicologia social é, talvez, o mais amplo em termos de estudos empíricos sobre o patriotismo usando *survey*. O objetivo principal dessas pesquisas girou em torno da identificação de atitudes patrióticas e nacionalistas entre indivíduos, e as mudanças de padrão ao longo do tempo. Uma grande parte delas foi realizada com base no módulo de estudos de identidade nacional do *International Social Survey Programme (ISSP)*, conduzido em 1995, 2003 e 2013, que contém mais de uma pergunta sobre o orgulho nacional e permite identificar diferentes fontes do orgulho pela nação.

Boa parte desses estudos tem relação com a Teoria da Identidade Social (TAJFEL; TURNER, 2004; TAJFEL, 1979), que sugere que indivíduos gostam de ser positivos com relação ao seu grupo social. Consequentemente, quanto mais eles se identificam com pessoas do mesmo grupo, mais eles tendem a avaliar o grupo positivamente (MUMMENDEY; KLINK; BROWN, 2001). Em uma das interpretações dessa teoria, o patriotismo se restringiria a uma avaliação interna positiva e autônoma, enquanto que o nacionalismo seria acompanhado de uma dimensão relacional de comparação e competição com grupos externos (DE FIGUEIREDO; ELKINS, 2003).

Estudos nessa linha tentam distinguir o patriotismo verdadeiro de um nacionalismo que se opõe ao internacionalismo e à globalização¹⁵, e que, consequentemente, tem uma visão negativa de imigrantes e da imigração em geral (KOSTERMAN; FESHBACH, 1989; SMITH; KIM, 2006; DE FIGUEIREDO; ELKINS, 2003; MUMMENDEY; KLINK; BROWN, 2001). Na mesma perspectiva, países menos homogêneos quanto à religião, ou em conflito internos ou externos apresentam níveis mais baixos de orgulho devido à competição entre grupos sociais (SMITH; JARKKO, 1998; ARIELY, 2016, 2018). Na perspectiva da ciência política, isso

¹⁵Ver, por exemplo, a crítica contundente de Billig (1995) ao argumento de Kosterman e Feshbach (1989).

equivale a dizer que o orgulho é menor quando existe uma discrepância entre o território nacional e suas instituições e a comunidade de pessoas que a formam.

Este orgulho de grupo aparece ligado a fatores culturais que estão na base da identificação social e variam entre uma sociedade e outra (EVANS; KELLEY, 2002; ARIELY, 2011; KAVETSOS, 2012)¹⁶. Pelas bases culturais, esse patriotismo seria “cego” a eventuais desvios de conduta da nação e estaria mais associado ao desengajamento político e mais sensível à exposição mediática pró-nacional (SCHATZ; STAUB; LAVINE, 1999). O orgulho também estaria relacionado à representação política do seu grupo étnico-cultural e declinaria com preconceito, status político ou conflitos nacionais (WIMMER, 2017; GREEN, 2020).

Em paralelo às identidades étnico-culturais, outra fonte de orgulho nacional está relacionada às instituições do estado. Esse orgulho, mais visível em nações desenvolvidas (BLANK; SCHMIDT, 2003; GREEN et al., 2011; SATHERLEY et al., 2019; NINCIC; RAMOS, 2012), estaria mais relacionado a um patriotismo virtuoso, cívico, construtivo, tolerante e com uma certa distância crítica do governo e do regime político. Nesses casos, o orgulho seria motivado pela vitalidade do sistema democrático, pela existência de um estado de bem-estar e segurança social, e daria sustentação a políticas de igualdade no tratamento dos cidadãos (DAVIDOV, 2009/ed, 2011; ARIELY, 2011; SCHATZ; STAUB; LAVINE, 1999; HJERM, 1998). Em termos de cultura política, esse patriotismo se encaixa no chamado apoio difuso ao sistema político (EASTON, 1965).

Quando fundado em sólidas bases tradicionais (local de nascimento, religião), esse patriotismo não se sobrepõe totalmente aos fatores objetivos e pessoais na identificação nacional (JONES; SMITH, 2001). Mas alguns dados empíricos indicam que este patriotismo cultural poderia atingir os objetivos pretendidos pelos princípios normativos da democracia liberal, mesmo dentro de grande diversidade cultural nos termos propostos por Habermas para um patriotismo constitucional (ARIELY, 2011).

Ainda com base em dados empíricos, a construção deste tipo de patriotismo reflexivo depende de uma certa estabilidade e longevidade do sistema político e de um esforço continuado de construção da nacionalidade (SMITH; JARKKO, 1998; ARIELY, 2016, 2018).

¹⁶ e.g. Os Ingleses se orgulhavam mais da ciência do que das artes, orgulho do desempenho econômico em todas as sociedades, orgulho no esporte em nações menores. Kavetsos (2012) identificou crescimento do orgulho depois de eventos esportivos.

Ele é suscetível a mudanças de regime que alteram a sua natureza (MUÑOZ, 2009) e pode ser abalado após traumas gerados por regimes autoritários ou em casos específicos como a culpa de guerra nos países do Eixo (SMITH; JARKKO, 1998)¹⁷.

Na perspectiva de um patriotismo reflexivo, a identificação com a nação, *estar orgulhoso*, também depende do desempenho do país (JONES; SMITH, 2001; NINCIC; RAMOS, 2012) e varia de acordo com os acontecimentos políticos e econômicos mais salientes em um estado-nação, principalmente os negativos que afetam o sentimento de orgulho na mesma direção (NOH, 2018; BONIKOWSKI, 2010; BONIKOWSKI; DIMAGGIO, 2016).

Paradoxalmente, o orgulho tende a ser maior em países mais desiguais (ARIELY, 2016, 2018). Uma das explicações possíveis é a política de diversionismo praticada pelo estado e pelos seus governantes. Aproveitando a situação precária dos cidadãos, líderes fomentam a identificação nacional. Este esforço de diversionismo anularia a falta de coesão que poderia existir devido às disparidades econômicas ou gerada por apelos separatistas de empreendedores políticos que se beneficiam de fatores psicológicos do sentimento de pertencimento intranacionais (SOLT, 2011; SHAYO, 2009). Qian e Hung (2018) demonstraram que esses efeitos acabam por se anular entre si.

Em nível individual, o orgulho parece ser mais prevalente entre indivíduos mais pobres (SHAYO, 2009; HAN, 2013; SHAYO, 2009). Com exceção de alguns países mais desenvolvidos, o orgulho tende a vir acompanhado de uma redução no apoio a políticas redistributivas (QARI; KONRAD; GEYS, 2012). Esse seria também um efeito do diversionismo que os governantes exercem sobre os cidadãos de forma a criar ou reforçar mitos e contrapor demandas sociais (SHAYO, 2009; SOLT, 2011; ECKEL, 2014). A pobreza e a desigualdade parecem afetar indivíduos diferentemente segundo a classe social (SHAYO, 2009).

Os efeitos da globalização podem ser deletérios para a construção da identidade individual e do patriotismo. Diversos estudos indicam que o orgulho é menor entre os mais jovens pelo efeito da globalização, permanecendo orgulhosas coortes que foram socializadas em períodos menos globalizados e tendo atingido níveis de educação formal menores

¹⁷Usando dados de uma *survey* de expertos da Europa e países da antiga União Soviética, Dimitrova-Grajzl, Eastwood e Grajzl (2016) obtiveram resultados que indicam que o orgulho nacional depende da longevidade da noção de identidade nacional, e que é uma convenção social historicamente determinada e que leva tempo para emergir.

(TILLEY; HEATH, 2007; SMITH; KIM, 2006; EVANS; KELLEY, 2002; BONIKOWSKI, 2010; BONIKOWSKI; DIMAGGIO, 2016; NOH, 2018; EVANS; KELLEY, 2002; ARIELY, 2011; KAVETSOS, 2012; LAN; LI, 2015). A globalização dificilmente rompe a ligação do patriotismo quando este tem bases étnicas (ARIELY, 2019). Apesar disso, alguns estudos apontam para a possibilidade de existirem identidades híbridas leais tanto a grupos subnacionais quanto a formas de identidades cosmopolitas (DUCHESNE; FROGNIER, 2007; FABRYKANT, 2014).

Assim como para a psicologia social, a expressão do orgulho pela nacionalidade é usada como operacionalização do patriotismo. E desde o *The Civic Culture*, estudo seminal e fundador da Cultura Política de Almond e Verba (1963; 1984), o patriotismo está presente nos estudos de cultura política usando *survey*. Nesse estudo, ingleses, alemães, italianos e mexicanos foram perguntados sobre as coisas do seu país que os deixavam mais orgulhosos. As alternativas eram entre o sistema político, legislação, cultura, valores espirituais, características físicas do país etc. Com base nos altos níveis de orgulho encontrados entre norte-americanos e britânicos, em contraste com italianos e alemães, Almond e Verba enfatizaram a importância do orgulho como parte da “cultura cívica”, especialmente pelo fato de que estes últimos eram orgulhosos também pelas conquistas políticas de sua nação¹⁸.

Nos estudos que sucederam o *The Civic Culture*, em especial os conduzidos por Ronald Inglehart (1971, 1977, 1997; 2005; 2018), o patriotismo passou a ser incluído em uma síndrome de valores tradicionais. Segundo Doob (1976), o patriotismo se trata de uma crença “mais ou menos consciente de uma pessoa de que sua vida depende de alguma forma da preservação e manutenção do poder e cultura da sociedade onde ele vive” (citado em VERBA, 1965). Para Inglehart, ao se socializarem em momentos de mais insegurança econômica, essa crença se consolida de forma mais clara e em conjunto com valores tradicionais.

Dentro da perspectiva de modernização proposta por Inglehart, o desenvolvimento econômico contribui para mudanças sociais, culturais e políticas (Marx) e a herança cultural de

¹⁸Outras conclusões dos autores podem ser destacadas: mexicanos e italianos se mostraram orgulhosos das qualidades ambientais/paisagísticas do país, enquanto britânicos e sobretudo americanos se mostravam mais orgulhosos do seu sistema político. Norte-americanos e mexicanos foram os mais orgulhosos do sistema econômico e das oportunidades que o país apresenta. Os americanos se mostraram os menos orgulhosos com seu povo. Italianos e mexicanos mais orgulhosos da sua contribuição às artes. Os autores argumentaram que esses resultados apontam para a existência de uma cultura cívica entre norte-americanos e britânicos, e uma alienação e paroquialismo dos italianos e dos mexicanos. Segundo os autores, estes últimos tinham o orgulho do país no que diz respeito à política devido à Revolução Mexicana.

uma sociedade molda suas crenças e motivações predominantes (Weber) (INGLEHART, R. F., 2018). Tendo em vista o efeito da afluência do pós-guerra nos países industrializados na redução dos níveis de patriotismo em direção a valores pós-materiais, Inglehart argumenta, com embasamento em dados empíricos, que a prevalência de sentimentos patrióticos é uma característica de países menos desenvolvidos e desiguais, e, em particular, os Latino-Americanos. Estes países, mesmo tendo alcançado algum nível de modernização (em termos de instituições democráticas, e, em menor escala, de igualdade de gênero, tolerância da homossexualidade), não perderam a religiosidade e o orgulho nacional.

Outra vertente da cultura política trabalhou o orgulho da nação dentro da perspectiva do apoio à democracia, seguindo a perspectiva dos autores do “Civic Culture”. O patriotismo, nessa perspectiva, é necessário para a estabilidade de um estado-nação e constitui uma fonte de apoio “difuso” ao sistema político em momentos de crise (EASTON, 1965). Essa abordagem foi desenvolvida por Dalton e Welzel (2014) e Norris (1999), que sustentam que existe uma transição das atitudes com relação ao sistema democrático, no qual cidadãos passam de leais a assertivos, e não necessariamente uma queda de apoio à ideia e a princípios democráticos. As declarações de orgulho tendem a ser menores nesses grupos. Outra relação estabelece que a democracia e as próprias instituições democráticas influenciam a cultura política e, conseqüentemente, o patriotismo contido nela (DAHL, 1973; MULLER; SELIGSON, 1994; citados em CASTRO, 2008).

Entretanto, a relação do orgulho com a democracia é ambígua. Estudando o caso dos países da antiga União Soviética, Haerpfer e Kizilova (2014) sustentam que orgulho pode simplesmente refletir a legitimidade do estado-nação, independentemente da forma do regime, aumentando a estabilidade do estado e não do regime em si. Assim como em estudos na psicologia social, os achados do estudo também indicam que o orgulho é geralmente mais fraco em países divididos em termos étnicos e linguísticos ou religiosos. Castro e Capistrano (2008), por sua vez, identificaram padrões de participação democrática diferentes nos países da América Latina.

Ainda nesta perspectiva de apoio à democracia, uma das características de cidadãos leais é a disponibilidade em lutar pelo país se for necessário. O crescimento de valores emancipativos faz com que cidadãos se sintam agentes individuais da política e menos afeitos a participar de guerras em nome do país (WELZEL; INGLEHART, 2010a). Puranen (2014)

identificou que em sociedades com valores mais emancipativos, onde os cidadãos são mais assertivos, lutar pelo país requer doses mais altas de orgulho.

Para os autores da cultura política, a globalização parece ter tido efeitos limitados em termos de mudanças em direção a orientações mais cosmopolitas na população. Uma das causas podem ser os apelos ao orgulho nacional que se mantiveram constantes (NORRIS, 2006, 2000; NORRIS; INGLEHART, 2009).

Enfim, uma das premissas básicas da cultura política é que existe uma congruência entre valores pessoais e instituições políticas (ECKSTEIN, 1997, 1988, 2015). As declarações de orgulho parecem reagir em função da congruência entre valores políticos pessoais e os valores em evidência no ambiente político (WOLAK; DAWKINS, 2016; FABRYKANT, 2014). A avaliação dos estímulos mais salientes (informação/notícias) é feita à luz dos valores políticos individuais (ZALLER, 1992).

2.2.1. Estudos empíricos sobre o patriotismo dos brasileiros

Em alguns poucos estudos empíricos comparativos na psicologia social, foram usados dados do Brasil, mas sem, necessariamente, apresentarem-se conclusões particulares ao país ou aos brasileiros¹⁹. Da mesma forma, os diversos estudos em cultura política desenvolvem pouco o tema para o Brasil.

Apenas alguns estudos foram identificados tratando unicamente do patriotismo no Brasil usando *survey*. Entre eles estão o estudo de Carvalho (1998) que testou e encontrou evidências da resiliência dos motivos edênicos na formação do orgulho nacional brasileiro, que partem do descobrimento e permanecem vivos até o final do século XX. Uma pesquisa de opinião patrocinada pelo WWF (2018) mediu o orgulho do brasileiro pelo seu meio ambiente e notou uma queda significativa.

Outro projeto realizado pelo Instituto República (2010) testou diversas teses sobre a identidade do brasileiro. Entre elas, encontrou evidências de que os brasileiros consideram o “jeitinho” como definidor da brasilidade positiva. Outro achado é que os brasileiros mais ricos

¹⁹Ver, por exemplo, Solt (2011), Ariely (2012), Ariely (2016), Fabrykant (2014), além dos estudos econômicos de Han (2013), Noh (2018) e Qian e Hung (2018).

(situados no topo da pirâmide de necessidades de Maslow)²⁰ se consideram mais patriotas, enquanto os mais abaixo se consideram alegres e “batalhadores”.

Em outro estudo, usando dados do Barômetro das Américas 2016/17, Noh (2018) encontrou que o Brasil é o país com prevalência mais baixa de pessoas muito orgulhosas (61%), ficando atrás apenas dos EUA e do Canadá. Informações mais salientes no ambiente político afetaram as respostas. Entre os principais motivos pelo baixo grau de orgulho estavam a piora da situação econômica, a corrupção e a insegurança, enquanto temas como desemprego foram pouco mencionados.

Por fim, em um estudo de psicologia social realizado com estudantes universitários brasileiros, Leite et al. (2018) testaram o papel do patriotismo, o nacionalismo e o essencialismo na perspectiva da construção de uma identidade social. Para esses autores, crenças nacionalistas e patrióticas são elementos constitutivos da identidade nacional dos brasileiros, mediadas pelo que os autores definem como essencialismo. Indo além do pertencimento a um grupo social, o essencialismo se refere à crença no pertencimento a um grupo específico, o nacional, com o qual compartilha de determinadas características ou atributos (PÉREZ-AGOTE, 1993). O estudo encontrou evidências da influência do pensamento essencialista quando os brasileiros são provocados a se pronunciar em questões sobre o nacionalismo e o patriotismo.

Quais representações da vida social e explicações do patriotismo que foram oferecidas acima e como podem servir de heurística para interpretar os dados que serão extraídos com o aprendizado de máquina? Para lembrar, na perspectiva metodológica proposta por Ragin adotada aqui, representação se refere à narrativa que descreve e explica um determinado fenômeno com base no diálogo entre ideias e as provas empíricas. A parte principal da representação é a imagem construída a partir das provas empíricas. Quadros teóricos ajudam na interpretação das evidências dando coerência à narrativa. É bom lembrar que, no caso

²⁰A pirâmide de Maslow (1943), usada também por Inglehart na tese do pós-materialismo, é uma teoria sobre as motivações humanas com base em necessidades. Maslow define cinco categorias de necessidades humanas: fisiológicas, segurança, afecto, estima e as de autorrealização. Esta teoria é representada por uma pirâmide onde na base se encontram as necessidades mais básicas pois estas estão directamente relacionadas com a sobrevivência.

dos dados do WVS que são usados a seguir, os próprios dados foram construídos a partir de alguma teoria e de certa forma a reproduzem. Isso implica que, de certa forma, o quadro teórico pode ser ele mesmo aprendido no aprendizado de máquina.

Tendo como premissa uma explicação baseada em crenças e valores, duas representações distintas parecem despontar com mais clareza deste diálogo entre teoria e empiria no patriotismo moderno. Ambas as versões são tributárias à criação do Estado-nação e da noção de comunidade “imaginada” proposta por Anderson, mas com pelo menos duas nuances principais. Na primeira delas, o patriotismo é representado como um amor incondicional à pátria, um pertencimento desinteressado e ainda muito próximo da lealdade direcionada às bases culturais e étnicas da família e do clã que foram transpostos à comunidade nacional como um todo. O orgulho de um indivíduo pela pátria é motivado pela crença na pátria como uma grande família ou clã que se imagina terem valores étnicos e culturais compartilhados.

Na segunda, talvez iniciada com Constant et Tocqueville, o patriotismo é representado como um sentimento mais racional e interessado. Um amor que, além de sentido, é pensado. Esse patriotismo seria virtuoso para os sistemas políticos, porque implica uma relação condicional e construtiva baseada na troca de direitos e deveres. O sentimento é interessado, pois tem base na consciência de que a vida individual depende da sobrevivência e da prosperidade da pátria como um todo. Assim, pequenas pátrias podem se constituir em pátrias maiores integradas seja pela crença na necessidade de abdicação de uma parte da liberdade pessoal em favor da vontade geral (versão republicana), seja pela crença na diversidade, na individualidade e na complementaridade dos seus membros (versão constitucional e cosmopolita). Esse patriotismo é sensível aos arranjos institucionais que permitem a construção da pátria comum e a convivência em uma mesma comunidade política. Esse patriotismo surge com a separação entre o Estado e a Igreja e, por isso, é menos motivado por valores e crenças religiosas, e disso se orgulha.

Essa distinção se enquadra na perspectiva de modernização proposta por diversos autores, entre eles Inglehart e os demais autores da cultura política. A primeira narrativa é compatível com valores tradicionais e a segunda, com valores seculares-rationais. Assim como a maior parte das inferências com base no sistema de valores, o patriotismo pode se manifestar nesses dois tipos ou em variantes intermediárias em uma sociedade em função da

maneira como cada indivíduo se relaciona com o sistema político.

Como visto acima, foram testadas diversas explicações para o enfraquecimento ou fortalecimento de uma ou outra forma de patriotismo. Valores religiosos e de família, valores políticos, globalização, desigualdade, pobreza têm efeitos distintos em uma ou outra forma de amor à Pátria. A figura 3 ilustra as duas narrativas de patriotismo e algumas das explicações já testadas empiricamente.

Figura 3: Patriotismos e suas influências. A figura ilustra a relação entre o indivíduo e o estado, objeto que encarna a pátria nos tempos modernos. Os adjetivos à esquerda e à direita das linhas pontilhadas correspondem às representações mais comuns do patriotismo. Do lado esquerdo, o patriotismo emotivo e cultural. Do lado direito, o patriotismo racional e cívico. Nas extremidades esquerda e direita, as principais hipóteses testadas empiricamente. Elaborado pelo autor com base nas referências teóricas e empíricas levantadas no trabalho.



No que diz respeito ao Brasil, os teóricos brasileiros enfatizaram diversas vezes a criação tardia de uma identidade verdadeiramente brasileira. A maioria das narrativas apresenta o patriotismo nacional aparecendo de forma incipiente ao final do período imperial, a partir da exaltação dos símbolos e dos heróis nacionais na campanha da Guerra do Paraguai. Antes, existiam apenas “patriotismos provinciais” e um antipatriotismo português antes da independência (CARVALHO, 2001). No século XX, pode-se dizer que os valores ligados ao patriotismo, ressentido pela escravidão, pela desigualdade, pela miscigenação, passaram por uma “transmutação”, fazendo com que o sistema patriarcal e a miscigenação de raças pudessem ser vistos de forma positiva (principalmente com Freyre).

No que diz respeito ao patriotismo reflexivo, apesar das tentativas de criá-lo de cima para baixo, o ressentimento também se fez presente. Pela falta de povo, ou pela inadequação do povo, o liberalismo nos moldes europeus não pôde ser implantado e uma solução autoritária prevaleceu até muito recentemente (ver O. Viana).

Para o caso brasileiro, para além de narrativas históricas, dados de *survey* foram pouco explorados de maneira a testar essas teses. As poucas evidências empíricas disponíveis apontam para um patriotismo ainda vinculado a aspectos simbólicos do Brasil como a imagem do paraíso e do eldorado e aspectos culturais como o samba e o futebol.

O efeito da saliência de temas como a corrupção, a insegurança e a desigualdade no orgulho nacional parece ser explicado mais por fatores históricos ligados à desconfiança crônica nas instituições políticas do que pela existência de um patriotismo cívico fundado em uma noção de cidadania mais ativa e que tem vergonha da inoperância do estado em resolver estes problemas. Reforça esse ponto a identificação do brasileiro com o jeitinho, que vai de encontro ao ideal de um patriotismo negociado com as instituições que, no caso brasileiro, não respondem às demandas dos cidadãos.

Não se pode descartar que um efeito de natureza cívica tenha exacerbado a queda no orgulho nas gerações mais jovens (pós-constituição de 1988). Infelizmente, poucos trabalhos foram realizados de maneira a entender qual é o patriotismo que restou da ação do estado autoritário, manipulador de valores e desenvolvedor de uma cidadania às avessas (CARVALHO, 1996).

Visto que a lógica metodológica adotada no trabalho é indutiva, nenhuma dessas teses será “testada”. A contribuição do trabalho vai no sentido de consolidar algumas dessas teses ou gerar novas hipóteses, o que precisa ser feito à luz do conhecimento preexistente sobre esse tema.

3. METODOLOGIA

Como descrito na introdução, a proposta da tese é explorar a integração da inteligência artificial nas ciências sociais. O primeiro capítulo apresentou o conceito de inteligência artificial e o foco nas técnicas de aprendizado de máquina. Para explorar a integração dessas técnicas no processo de construção do conhecimento nas ciências sociais, o patriotismo foi escolhido como caso de estudo e foi definido e descrito no capítulo 2.

O processo de construção do conhecimento nas ciências sociais, tal como entendido neste trabalho, ocorre em um processo contínuo que alterna momentos indutivos e (hipotético-)dedutivos. Este capítulo descreve os métodos e as técnicas que foram utilizados nesses dois momentos e as diferentes decisões que foram tomadas ao longo deste processo. A primeira parte do capítulo se dedica à descrição dos dados que foram usados na tese. Em seguida são apresentadas, brevemente, as diferenças que envolvem a modelagem estatística em cada momento. Feita essa distinção, são apresentadas as técnicas aplicadas na modelagem dedutiva e os métodos e as técnicas do aprendizado de máquina na terceira e na quarta parte do capítulo. A parte dedicada ao aprendizado de máquina é mais extensa, e inclui o estudo comparativo que foi realizado para selecionar um algoritmo dentre os muitos existentes; o treinamento ou aprendizado propriamente dito; e as técnicas usadas para interpretação dos resultados. Este passo final, essencial, permite fechar o ciclo e introduzir o capítulo seguinte, que trata da inferência e da integração dos achados no processo de criação de conhecimento em um novo momento dedutivo.

3.1. A PESQUISA MUNDIAL DE VALORES

Foram usados neste trabalho os dados da Pesquisa Mundial de Valores, chamado aqui pelo acrônimo WVS, do inglês *World Values Survey* (INGLEHART; HAERPFER et al.,

3 Vs (Volume, Velocidade e Variedade) (FAVARETTO et al., 2020; OLLION; BOELAERT, 2015).

A base de dados gerada por estas diversas ondas é organizada em uma matriz retangular onde cada linha representa uma observação (ou caso) e cada coluna, uma variável. A base em sua versão mais atual contém 618 405 observações e mais de 1 500 variáveis. As sete ondas realizadas entre 1981 e 2019 foram conduzidas em mais de 110 sociedades em seis continentes habitados (ver Figura 4).²

Como a *survey* não foi coletada especificamente para a tese, os dados usados são considerados como secundários. Mesmo usando dados secundários, a tese apresenta dados inéditos para a onda mais recente de 2018-2019. Para a lista completa de sociedades e ondas, ver o Apêndice A.

3.1.1. Amostras de treino, validação e predição

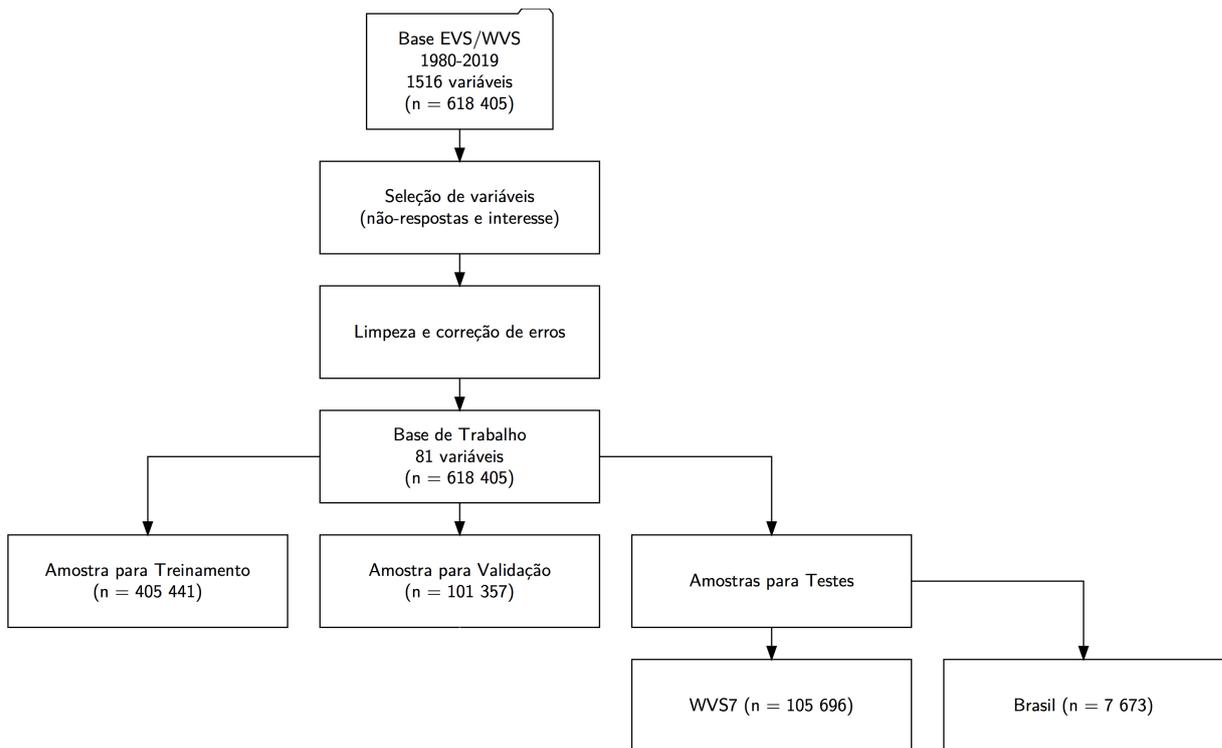
Para evitar o problema do sobreajuste (*overfitting*) no aprendizado de máquina, descrito mais detalhadamente nos próximos parágrafos, a base de trabalho foi separada, primeiramente, em duas amostras em uma abordagem chamada conjunto de validação, típica da modelagem preditiva (JAMES et al., 2013, p. 176). A amostra de treinamento contém 80% das observações ($n=405\ 441$) e a amostra de validação contém 20% dos dados ($n = 101\ 357$), proporção usada em outros trabalhos (KUHN; JOHNSON, 2013). Existem diversos critérios usados para fazer a divisão das amostras (MARTIN et al., 2012). A escolha aqui foi fazer uma separação aleatória das amostras respeitando-se apenas as proporções existentes da variável operacionalizada como patriotismo (*i.e.*, mantendo-se a mesma proporção em ambos os conjuntos da variável “Orgulho da nacionalizada” dicotomizada).

Além disso, foram separadas duas amostras para fazer testes da capacidade preditiva dos modelos. Um dos critérios de uma boa predição é olhar adiante em busca de dados (BEAUCHAMP, 2017). Modelos preditivos poderosos ou robustos devem predizer a ocorrência de casos em amostras ainda não vistas. Para simular essa situação usando os dados disponíveis,

²Os seis continentes são Europa, América do Norte, América do Sul, Ásia, África, Oceania. A escolha do termo sociedade em vez de país é devido ao fato de que em alguns casos foram conduzidas *surveys* em diferentes sociedades dentro de um mesmo país, como por exemplo na Catalunha.

foram separadas da base de trabalho a amostra da onda 7 do WVS ($n = 105\ 696$), realizada em 2019, e de todas as ondas no Brasil ($n = 7\ 673$). A Figura 5 resume a separação que foi efetuada da base de trabalho nas amostras de treinamento, validação e testes de predição:

Figura 5: Separação da base de dados do WVS. A figura indica as diferentes amostras usadas no trabalho e a quantidade de observações, ou casos, em cada uma delas.



3.1.2. Seleção e transformação das variáveis

Uma vez separadas as amostras, um aspecto comum da modelagem estatística é a transformação de algumas variáveis de interesse antes de fazer o ajuste ou lançar o aprendizado. A etapa de transformação dos dados envolve a aplicação de técnicas que permitem extrair o potencial máximo das variáveis explicativas de uma base de dados.

Essas transformações, ou engenharia (*feature engineering*), como é chamada no aprendizado de máquina, implica o uso de conhecimento do pesquisador para trabalhar a entrada dos dados (variáveis explicativas) para que os algoritmos funcionem melhor (DOMINGOS, 2012; KUHN; JOHNSON, 2019)³.

³Um exemplo apresentado em Kuhn e Johnson (2019) é o caso do endereço de uma propriedade à venda. A localização é uma informação crucial que pode ser representada de várias formas. Podemos tanto usar o CEP

A variável resposta escolhida para supervisionar o aprendizado de máquina foi a pergunta que mede o orgulho nacional autodeclarado (variável G006 na base Longitudinal). A variável, operacionalizada como uma aproximação do conceito de patriotismo, foi transformada em uma variável binária contendo duas classes. Foram considerados patriotas os que se declararam “Muito Orgulhosos” e “Orgulhosos” da sua nacionalidade. Os demais, que se declararam “Não muito orgulhoso” ou “Não sou orgulhoso”, foram considerados com baixo nível de patriotismo. A perda de informação gerada pela dicotomização é explicada por duas razões principais. A primeira delas é pela facilidade de tratamento estatístico e interpretação dos resultados. Outra razão, mais substantiva, é que se assumiu que existe uma dicotomia entre orgulhosos e não orgulhosos para a qual queremos um modelo preditivo.

Quase todas as variáveis explicativas também sofreram alguma transformação. Na base do WVS, as variáveis são, em sua maioria, categóricas, o que em geral não é compatível com os algoritmos de aprendizado de máquina. As variáveis categóricas que continham apenas duas respostas, em geral “Sim” ou “Não”, foram transformadas em variáveis binárias, como, por exemplo, a série de questões sobre a qualidade desejada dos filhos, a posição sobre os vizinhos, a confiança interpessoal, o trabalho feminino e migrante, os objetivos do país, mudanças futuras, ação política, religiosidade, sexo, estado civil e ter filhos ou não etc.

As variáveis categóricas ordenadas, em escalas tipo *Likert* de 4 ou 5 níveis, em sua maioria, foram transformadas em uma escala numérica, como, por exemplo, a série sobre o que é importante na vida, o estado de saúde, as posições sobre a família, o interesse na política, a confiança nas instituições, o sentido da vida, a prática religiosa.

As variáveis categóricas não ordenadas foram codificadas usando-se a técnica de codificação *one-hot*, que trata de transformar cada categoria da variável em uma variável binária distinta. Nesses casos o índice de referência foi retirado ($n-1$, onde n é igual ao número de categorias da variável.) para evitar a correlação entre as variáveis geradas. Algumas variáveis categóricas, que na sua forma original foram medidas em uma escala de 1-10, foram convertidas em escalas numéricas⁴.

para representá-la quanto a coordenada geográfica exata. Do ponto de vista informacional, a coordenada exata contém mais informação do que o CEP. Da mesma forma, os dados do EVS/WVS podem ser representados na sua forma original ou transformados para melhorar as previsões. Os autores citados acima sugerem um ordem para o pré-processamento começando pelas imputações, transformações individuais, discretização, numerização ou binarização, interações, normalização e transformações envolvendo mais de uma variável.

⁴Devido à distribuição fortemente assimétrica, as variáveis da série F “Justificável”, foram testadas

Uma vez convertidas, as variáveis numéricas foram normalizadas. Enfim, variáveis com variância quase inexistente, bem como as variáveis muito correlacionadas entre si, foram retiradas devido à influência dessas variáveis em muitos algoritmos de aprendizado de máquina (GREGORUTTI; MICHEL; SAINT-PIERRE, 2017; STROBL et al., 2007).

Para que os resultados do aprendizado possam ser usados para predição, a base de treino, validação e testes devem ser representadas da mesma forma. Assim, o mesmo pré-processamento realizado na base de treinamento foi replicado nas amostras de validação e de testes. O pré-processamento de forma individualizada nas bases previne o chamado *vazamento* de informações entre as bases e violaria o princípio de se medirem erros de treino e validação, o que tem consequências no sobreajuste do modelo.

3.1.3. Tratamento de não-respostas

A proporção de não-respostas na amostra longitudinal do WVS é importante. Esse padrão é encontrado em outras bases de dados usadas nas ciências sociais (KING, G. et al., 2001; DURRANT, 2005). De acordo com a classificação em pesquisas do tipo *survey*, as lacunas são em nível das unidades e de itens. O primeiro, refere-se a quando uma pessoa que deveria ter respondido a uma *survey* não estava disponível, ou não quis responder; o segundo, quando a pessoa deixa de responder a uma ou mais perguntas do questionário (SALANT; DILLMAN, 1994; DILLMAN; SMYTH; CHRISTIAN, 2014; DILLMAN; PHELPS et al., 2009). No caso do WVS, as lacunas em nível das unidades são devido ao fato que das 1.500 questões do questionário completo, muitas delas só foram aplicadas em uma onda ou país. Dentro de cada uma das *surveys* sempre existe uma proporção de não-respostas em nível dos itens.

Na lógica dedutiva, não-respostas limitam o pesquisador na medida em que podem causar vieses e comprometer a possibilidade de fazerem-se inferências a partir dos dados. Isso porque as regressões aplicadas na modelagem dedutiva precisam descartar as não-respostas para fazer o ajuste.

Já a lógica indutiva envolve o estudo aprofundado dos possíveis padrões e relações entre variáveis. Na modelagem estatística, essa abordagem, que vemos com mais detalhe

algumas transformações dessas variáveis na escala logarítmica, sem que tenha havido melhora significativa nos resultados.

abaixo, envolve o uso de todas as observações e variáveis possíveis de uma determinada base de dados, de forma a aumentar o poder preditivo. A preocupação com a parcimônia dos modelos é secundária, no sentido que mesmo variáveis com pouco peso podem aumentar o poder preditivo do modelo. Mesmo variáveis que não teriam poder explicativo nenhum segundo as teorias são bem-vindas. Nas palavras de Leo Breiman (2001b), a seleção *a priori* de variáveis não é uma boa ideia na modelagem preditiva, pois seria como “procurar as chaves do carro perdidas (apenas) na parte iluminada da rua”.

Na perspectiva do aprendizado de máquina, as não-respostas apresentam um problema duplo. O primeiro é na construção do modelo, porque diversos algoritmos não administram as não-respostas de maneira automática (como o Florestas Aleatórias, por exemplo)(MERCALDO; BLUME, 2017). O segundo problema reside no uso do modelo para fazer previsões em observações futuras, quando nelas também existem dados faltantes. Por razões de tempo e praticidade, o estudo comparativo para escolha do algoritmo foi realizado usando-se a base de casos completos da amostra de treinamento. O *Extreme Gradient Boosting* (XGBoost), que foi usado para o treinamento final e que vemos com mais detalhes mais abaixo, não descarta as não-respostas e o treinamento foi realizado sem a necessidade de imputação (CHEN; GUESTRIN, 2016).

De forma a minimizar esses problemas, antes de iniciar-se a modelagem estatística, foram selecionadas as variáveis mais recorrentes. Mesmo após essa seleção, ainda restaram muitas lacunas. Entre as variáveis da base de trabalho, como foi chamada a base depois dessa primeira triagem, a proporção de não-respostas variou entre 0 e 35%, com uma média geral de 8,9% (4,5 milhões de células de um total de 50 milhões). Dos 618.405 casos da base, apenas 70.042 casos são completos (11,3%)⁵. No caso das variáveis explicativas escolhidas para o modelo dedutivo, a porcentagem de não-respostas variou entre 5 e 35%⁶.

Além do ajuste usando-se as observações completas, foi testada uma técnica de imputação múltipla para ajustar o modelo com a mesma quantidade de dados usada no

⁵Foram também retiradas diversas variáveis relacionadas à estrutura dos dados (principalmente da série S) e os índices computados *a posteriori* a partir de outras variáveis (*i.e.* Pós-materialismo, Índice Emancipação, etc.), de modo a privilegiar as variáveis originais da pesquisa. Foram também excluídas variáveis para as quais não há dados para o Brasil ou na onda 7. Essas variáveis ficaram de fora por causa da natureza da representação dos dados na modelagem preditiva, na medida em que seria inútil incluir variáveis no modelo que não poderiam ser usadas no momento da previsão.

⁶E001: 28,1%; E003: 4,28%; E005: 35,4%; X025: 34,1% e X047: 20,3%.

aprendizado de máquina. A técnica de especificação condicional completa entre os dados (VAN BUUREN et al., 2006) permite substituir os valores faltantes por valores estimados de modelos de imputação específicos para cada variável a ser imputada e suas covariantes. Usando esta técnica, foram geradas cinco bases imputadas. Em cada uma delas, foram imputados ao total 756 156 respostas em todas as variáveis. Mais informações sobre as não-respostas e a técnica da imputação múltipla testada aqui pode ser encontrada no Apêndice B.

3.2. DEDUÇÃO E INDUÇÃO NA MODELAGEM ESTATÍSTICA

A base longitudinal contém entrevistas realizadas nas últimas três décadas em diversos países do mundo. Esta riqueza de informações permite tanto operacionalizar construtos teóricos em uma modelagem explicativa quanto aplicar técnicas de modelagem indutiva para fazer inferências a partir delas. A cada um desses momentos corresponde a um tipo de modelagem estatística particular. Breiman (2001b) e Shmueli (2010) consideram que a explicação e a predição são os dois principais caminhos, ou culturas da modelagem estatística. Cada um deles tem características e objetivos distintos, mas complementares, que correspondem a esses dois momentos⁷. As disparidades que existem entre os caminhos explicativo e preditivos podem ser ilustrados em algumas dimensões principais (Quadro 2).

Tabela 2: Disparidades entre os caminhos preditivos e explicativos

| | Caminho Explicativo | Caminho Preditivo |
|-------------------------------|--|---|
| Causalidade- Associação | A função ajustada representa uma teoria onde x explica y | A função ajustada captura a associação entre as variáveis |
| Teoria- Dados | O modelo é construído com base na teoria, que apoia a interpretação dos resultados | O modelo é construído à partir dos dados e não requer interpretação |
| Retrospectiva- Prospectiva | Voltado a explicar o passado | Voltado a prever futuras ocorrências |
| Vieses- Variância | Foco na minimização dos vieses para representar a teoria com acurácia | Compromisso entre vieses e variância com foco na melhoria das predições |

^a Fonte: Adaptado de Shmueli, 2010

⁷Breiman inclui um terceiro: a modelagem descritiva, próximo da explicativa, que envolve capturar e descrever a estrutura dos dados e associações entre variáveis. O principal objetivo é apresentar de forma parcimoniosa a distribuição condicional de uma variável resposta, y , dada uma série de variáveis explicativas, x .

A modelagem explicativa corresponde ao método hipotético-dedutivo. São os modelos que servem para testar hipóteses ou possíveis explicações para um dado fenômeno. Em modelos explicativos, espera-se que uma série de variáveis explicativas, x , expliquem ou causem um determinado efeito, medido pela variável resposta, y . Os modelos explicativos são entendidos assim devido à existência de um construto teórico por trás de cada modelo. É a teoria que estabelece a causalidade, que é testada usando-se ferramentas estatísticas que identificam o nível de associações entre as variáveis.

Modelos preditivos, por sua vez, são essencialmente usados para prever novas ou futuras observações e obedecem a uma lógica indutiva ou de aprendizado com os dados. O objetivo da modelagem preditiva é prever um valor de saída, y , em futuras observações, dada uma série de valores de entrada, x . Em outras palavras, trata-se de estimar a função $f(x)$, onde a sua acurácia é medida usando-se uma função custo que mede a discrepância entre as previsões e as ocorrências reais. Os modelos preditivos não dependem de uma formulação teórica *a priori*, e a base de sustentação da coerência do modelo reside na sua capacidade preditiva.

Abaixo, são descritos os métodos em cada um dos momentos.

3.3. MOMENTO DEDUTIVO

No capítulo anterior, sobre o patriotismo, foram revisadas algumas das principais teorias e achados empíricos que estão disponíveis sobre o tema. A partir dessas informações, o caminho natural seguido pelos pesquisadores nas ciências sociais orientados aos métodos quantitativos é a identificação de uma ou mais teorias que se deseja testar; a operacionalização dentro da perspectiva teórica escolhida; a coleta de dados empíricos; e, finalmente, a análise dos dados usando métodos estatísticos.

Para o problema posto aqui, em uma lógica hipotético-dedutiva, o cientista social se interessa pela explicação do porquê um determinado indivíduo, ou a sociedade como um todo, é mais ou menos patriótico. A dimensão explicativa decorre da existência de um construto teórico subjacente à modelagem estatística. Uma vez construída a teoria, os modelos estatísticos buscam testar esses construtos. As associações que são percebidas após o ajuste do modelo são consideradas provas de que as hipóteses testadas são plausíveis. Assim, o fenômeno é

primeiro explicado em nível teórico para em seguida ser testado com dados empíricos.

Considere o exemplo do fenômeno social do patriotismo como construto teórico (Y) representado pela função teórica F , explicado por outros construtos teóricos (X), resultando em $Y = F(X)$. Em uma modelagem estatística explicativa a função F deve ser operacionalizada na forma de um modelo estatístico f . Para que a inferência estatística se aplique às hipóteses teóricas, tenta-se aproximar ao máximo as funções f e F , de forma que $F(y) = f(x)$.

Nesse modelo estatístico, as variáveis empíricas x e y são as operacionalizações das variáveis resposta e explicativas e são usadas para construir a função f , que é usada para testar as hipóteses. No caso presente, a variável y , “Orgulho da Nacionalidade”, é o construto escolhido para representar Y , o Patriotismo. As variáveis explicativas x_n , por sua vez, são operacionalizadas de modo a representar da melhor forma as hipóteses teóricas X_n .

Tendo em vista que a variável resposta tem apenas dois valores possíveis, Alto ou Baixo Patriotismo, é necessário aplicar uma técnica que seja capaz de ajustar um modelo entre 0 e 1. A técnica selecionada para testar as hipóteses teóricas com essas características é naturalmente a regressão logística binária que serve para estudar associações entre uma variável resposta desse tipo e outras variáveis explicativas (AGRESTI, 2007).

A regressão logística é um tipo de regressão linear que modela uma razão de chances (*odds ratios*): a fração de eventos, ou casos, que se encontram em determinado estado (*i.e.* neste caso, “Alto orgulho”) com relação às que não se encontram nesse estado (*i.e.* neste caso, “Baixo orgulho”). O modelo binomial *logit* estima o logaritmo das chances e os efeitos multiplicativos nas chances de um evento ser observado, que podem ser interpretados como a razão entre o número esperado de “sucessos” e “fracassos” (BUIS, 2015). De acordo com este autor, a regressão logística, pensada dessa forma, mede o quão plausível, ou quais as chances de um evento ocorrer, onde probabilidades, chances e a razão de chances são as maneiras de quantificá-la.

Feito o ajuste⁸, a interpretação do modelo de regressão logística se deu, inicialmente,

⁸Os parâmetros da regressão foram ajustados tanto na base contendo não-respostas quanto nas bases imputadas. Neste caso, a regressão foi ajustada em cada base imputada e os coeficientes agregados (*pooled*) usando-se as médias de Rubin (2004). Esses coeficientes agregados foram usados para calcular as probabilidades de classificação de cada caso e as predições.

pela avaliação da significância estatística de cada variável pela razão de chances. Para facilitar a comparação, os coeficientes foram normalizados de modo a que uma unidade de mudança represente uma unidade de mudança no erro padrão (GELMAN, 2008). Além disso, o modelo de regressão foi avaliado usando-se estatísticas globais e de qualidade do ajuste (*goodness-of-fit*, ou *GOF*)⁹.

3.4. MOMENTO INDUTIVO

No caso do patriotismo, de fato existem uma série de relações potenciais entre variáveis que podem ser deduzidas a partir das teorias existentes. Essas relações podem ser testadas para confirmar ou reforçar uma ou mais explicações para o patriotismo. Mas nem sempre isso é possível. A falta de teorias, ou a necessidade de atualizá-las ou de melhorá-las, muitas vezes faz com que seja pertinente a pesquisa em uma lógica indutiva com a aplicação de técnicas capazes de aprender com os dados existentes.

Em um modelo preditivo, a função decorrente do aprendizado não depende de nenhum construto teórico. A função f é *induzida*, ou aprendida, usando-se x e y e, em seguida, usada como uma ferramenta para gerar previsões em valores y ainda não observados. De forma resumida, busca-se, com base nas informações disponíveis, prever no futuro se um determinado indivíduo é patriótico ou não.

Na cultura da modelagem preditiva e na terminologia do aprendizado de máquina, este é um problema de classificação. Em termos formais, tendo em vista que o valor de saída (*output*), ou variável resposta, tem somente dois valores possíveis $y = \{0, 1\}$, correspondendo às duas classes possíveis de Alto e Baixo Patriotismo, trata-se de um problema de classificação binária. As classes, estando já classificadas (*labelled*) em termos de orgulho nacional, permitem que as iterações, ou o aprendizado em si, seja supervisionado. Busca-se aqui, portanto, encontrar uma solução para um **problema de classificação binária, supervisionado**.

Definido o problema, é preciso tomar uma série de decisões. A escolha de um algoritmo é o primeiro desafio no aprendizado de máquina. Existem diversos algoritmos aptos a tratar

⁹A avaliação global do modelo se dá em função da melhora que o modelo traz com relação ao modelo com somente o Intercepto (modelo nulo), como o Akaike Information Criterion (AIC). Em termos de qualidade do ajuste, foram realizados testes como os Pseudo R^2 (e.g. Cox & Snell, Nagelkerke, Tjur), e o teste de Hosmer-Lemeshow, que testa as evidências para descartar a hipótese de um modelo bem ajustado.

do problema de classificação binária proposto na tese, e a escolha não tem como ser feita de forma teórica *a priori* devido a resultados muitas vezes similares usando-se algoritmos concorrentes (WOLPERT, 1996). Além disso, cada “tribo” do aprendizado de máquina tem suas preferências e algoritmos que estão em contínuo desenvolvimento (DOMINGOS, 2015)¹⁰. O pesquisador pode não ter condições de conhecer plenamente cada um deles, ou uma combinação deles, o que dificulta a escolha (FERNÁNDEZ-DELGADO et al., 2014).

A abordagem escolhida neste trabalho foi de iniciar por um teste de uma série de algoritmos, a fim de avaliar, comparativamente, seu desempenho e escolher um ou mais deles para um treinamento mais completo. Um estudo desse tipo pode ser chamado de *benchmarking* ou comparativo. A implementação do estudo seguiu o procedimento descrito em Hothorn, Leisch et al. (2005) e Manuel J. A. Eugster, Hothorn e Leisch (2016). Em um estudo dessa natureza, os algoritmos devem ser treinados com os mesmos parâmetros de validação cruzada, e o melhor ajuste deve ser escolhido usando-se a mesma métrica de avaliação. As comparações podem ser usadas também para estimar o tempo de processamento necessário para o treinamento (e o seu custo), o que ajuda a determinar a viabilidade da implementação de um algoritmo para uma determinada tarefa.

Pelas limitações de tempo e processamento, foram utilizados apenas os casos completos da base de treinamento para a comparação ($n = 54\ 086$). Em princípio, o desempenho dos algoritmos pode variar com o tamanho da amostra, mas, segundo alguns autores, o desempenho tem mais relação com a complexidade ou com a representação dos dados (MACIÀ et al., 2013). Além disso, a base de casos completos é substancialmente maior do que as bases usadas em experimentos similares, o que garante certa credibilidade a essa abordagem.

O XGBoost foi o algoritmo que obteve os melhores resultados e, com ele, foi realizado um treinamento mais extensivo. O *eXtreme Gradient Boosting*, proposto por Chen e Guestrin (2016), é uma evolução recente das técnicas desenvolvidas nos algoritmos *AdaBoost* e *GBM - Stochastic Gradient Boosting Machines* (FREUND; SCHAPIRE, 1997, 1999), e na Máquina de Boosting Gradual (FRIEDMAN, 2001; FRIEDMAN, J. H.; HASTIE; TIBSHIRANI, Robert, 1998; FRIEDMAN, 2002). O XGBoost melhora essas técnicas a partir de otimizações e melhorias em termos de paralelização, cortes de baixo para cima (*pruning*), uso do cache,

¹⁰Como vimos anteriormente, o autor caracteriza as cinco tribos como: Simbolistas, Conexionistas, Evolucionistas, Bayesianos, e Analogistas.

regularização (*shrinkage* e *column subsampling*), e, também, a gestão de não-respostas.

O XGBoost é um algoritmo do tipo *ensemble* que combina diversas árvores de decisão e usa a técnica do *boosting* gradiente para otimização. A ideia principal dos *ensembles* é combinar diversos algoritmos fracos e simplificados. No caso do XGBoost são combinadas diversas árvores de decisão tipo CARTs (*Classification and Regression Trees*), usando-se a técnica de *bagging* (*Bootstrap AGGregatING*). Essa técnica consiste em incorporar a contribuição de cada uma das árvores por meio de um processo de votação.¹¹

A ideia principal do *boosting* é adicionar novos modelos (*i.e.* novas árvores) ao *ensemble* de forma sequencial. O aprendizado envolve ajustar/otimizar consecutivamente cada modelo de forma a encontrar estimativas mais acuradas da variável resposta. Cada nova árvore $N + 1$ do XGBoost foca no que não foi aprendido na árvore N (no Florestas Aleatórias, as árvores N e $N + 1$ são independentes).

O termo gradiente, por sua vez, vem do fato de que cada nova árvore treinada é otimizada em função da precedente usando-se a técnica do gradiente descendente. Em termos de aprendizado, isso consiste em avaliar o erro do modelo N (testado em uma sub-amostra de validação) e modificá-lo seguindo a maior inclinação descendente da curva de erro de modo a minimizá-la no modelo $N + 1$ ¹². Essa técnica permite escapar dos mínimos locais que interromperiam o aprendizado (NATEKIN; KNOLL, 2013). O hiperparâmetro “encolhimento” (*shrinkage*) garante que a influência de cada árvore seja equilibrada em relação às futuras/próximas, para que todas possam continuar contribuindo para a otimização do modelo final.

Como outros algoritmos de *gradient boosting*, o XGBoost usa a função *log loss* como função custo para aprender a relação entre as variáveis explicativa x_n e resposta y . Essa função mede o desempenho do modelo em termos de classificação, sendo o resultado uma probabilidade entre 0 e 1. Um modelo perfeito teria um *log loss* igual a 0. O custo aumenta

¹¹O XGBoost difere do Floresta Aleatória, que também é uma combinação de árvores de decisão. Em uma Floresta Aleatória, apenas uma parte das variáveis explicativas é sorteada de maneira aleatória para cada árvore (*column subsampling*) e são computadas médias para formar o ensemble.

¹²Pesquisadores descrevem diversos fenômenos naturais de otimização: um espermatozóide segue a direção de maior inclinação em termos de concentração dos marcadores químicos do óvulo; um mosquito segue a direção de maior inclinação da curva de distribuição luminosa para chegar até a fonte de emissão da luz. No caso do AM se trata-se do campo de parâmetros do modelo sendo treinado. Exemplos Boullier e El Mhamdi (2020).

quando são atribuídas probabilidades baixas a uma observação que tem a classe verdadeira. O *log loss* penaliza tanto falsos positivos (FP) quanto falsos negativos (FN), mas em especial as previsões que têm alta probabilidade e são erradas.¹³

Tanto no estudo comparativo quanto no treinamento com o XGBoost outras decisões importantes foram tomadas. A primeira delas se refere à minimização do sobreajuste e, conseqüentemente, uma melhor capacidade de generalização do modelo gerado.

Existem diversas estratégias para tentar encontrar o melhor compromisso e aumentar o poder de generalização do modelo. A primeira delas envolve a separação dos dados em amostras de treino e validação mencionadas acima que permite calcular os erros em ambas. A importância do aprendizado usando-se diversas amostras tem relação com o problema que Shmueli (2010) chama de compromisso viés-variância. No contexto do aprendizado de máquina, viés corresponde à diferença entre as previsões e os casos reais, e variância é a sensibilidade do algoritmo a variações nos conjuntos de dados. Na modelagem preditiva, a capacidade de generalização está no centro das preocupações. Espera-se que o modelo classifique casos reais com o menor risco possível. Estimar a função f envolve, assim, a busca por um equilíbrio entre variância e viés. Quando um modelo é menos complexo e mais rígido, ele tende a não ajustar bem com os dados (alto viés). Contrariamente, quando o modelo é mais complexo e flexível, ele tende ao sobreajuste (alta variância). Como a minimização do erro é o objetivo do treinamento, o modelo tende sempre a ter sobreajuste na amostra na qual está sendo treinado. Por isso, usam-se duas amostras diferentes, de modo a poder estimar o erro em uma amostra não usada no treinamento.¹⁴

Outra estratégia é o uso de subconjuntos (*k-folds*) aleatórios da própria amostra de treinamento durante o aprendizado. Algumas das possibilidades são a amostragem via *bootstrapping*, *leave one out* ou validações cruzadas repetidas ou não. A técnica escolhida aqui foi a de validações cruzadas com repetições (*repeated k-fold cross-validations*) proposta por Weiss e Kulikowski (1991). No presente estudo, a validação cruzada foi realizada em 10

¹³A *log loss* também é conhecida como entropia cruzada. Cada divisão das árvores (cada nova decisão) implica um ganho informacional na medida em que reduz a entropia dos dados. Trata-se, assim, de encontrar a combinação de variáveis que retorna o maior ganho de informação, tornando os ramos da árvore mais homogêneos e com menor entropia.

¹⁴A separação dos dados é uma das principais precauções que se deve tomar para minimizar o sobreajuste. Isso é o contrário do que acontece na modelagem explicativa, que envolve, geralmente, usar todos os dados disponíveis de forma a minimizar o erro do modelo.

subamostras e 10 repetições, o que garante um resultado mais robusto (WITTEN; FRANK, 2011, p. 151). Aplicar a validação cruzada dessa forma equivale a treinar o algoritmo 100 vezes usando as subamostras que são 9/10 do tamanho original da amostra de treinamento (o 1/10 restante é usado para medir o erro a cada iteração).

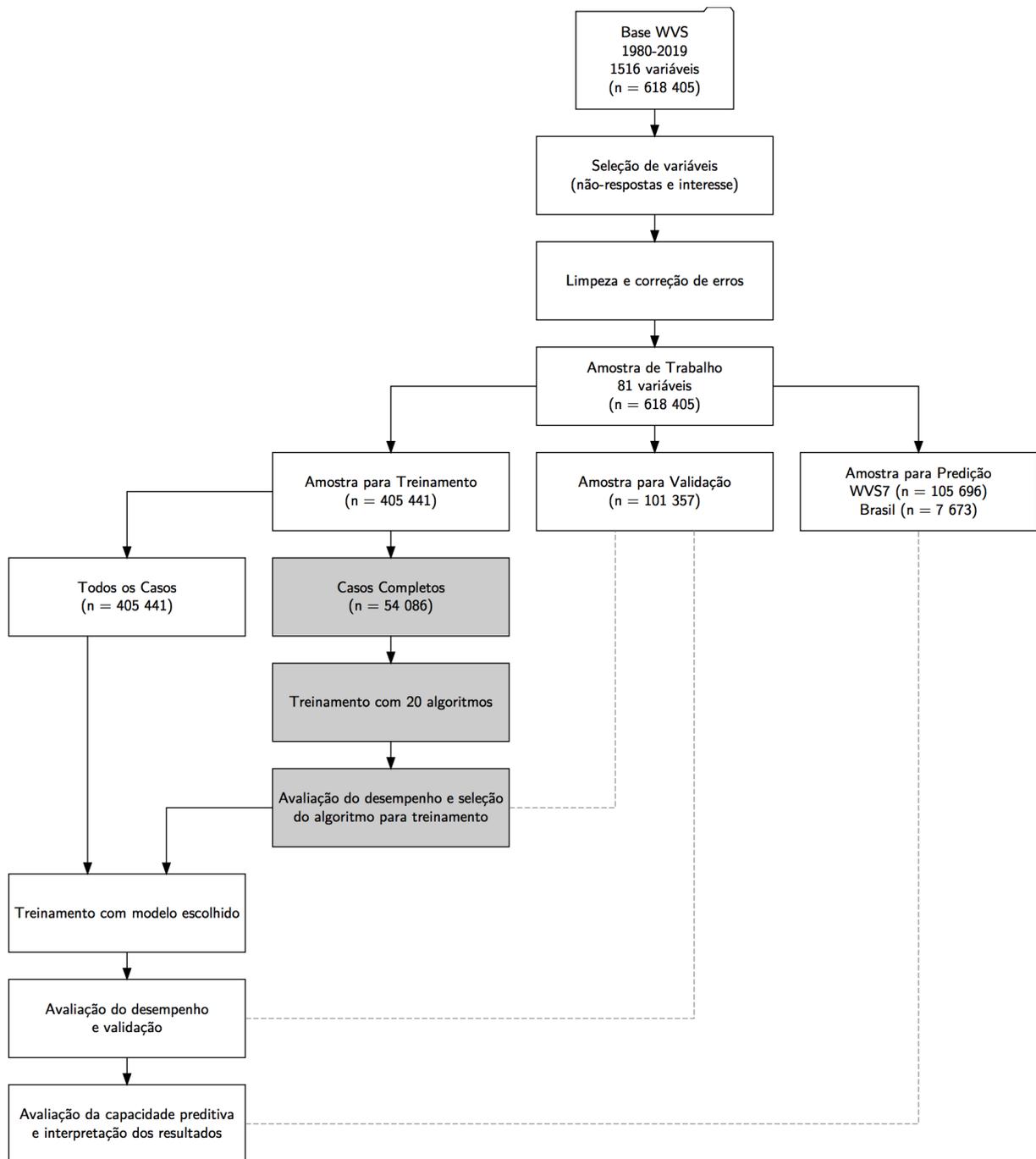
Os resultados finais de cada modelo treinado são as médias do número de repetições nas subamostras. A capacidade de generalização do modelo é avaliada (ou validada) comparando-se este desempenho médio (erro) com o desempenho nos dados de validação. Se o desempenho for significativamente melhor na amostra de treinamento, estima-se que há sobreajuste no modelo. A capacidade de predição, ou de generalização real, só pode ser testada e avaliada observando-se o desempenho do modelo com futuras observações. Para isso, foram separadas as bases do WVS7 e Brasil. A Figura 6 resume o procedimento de treinamento usando-se as diferentes amostras de treinamento, validação e testes.

Outras estratégias que permitem aumentar o poder de generalização estão no grupo das regularizações. Regularização pode ser definida como qualquer parte de um modelo que reconhece as limitações dos dados, ou a sua variância (ROSSET, 2003). Uma das técnicas de regularização consiste em aplicar restrições ou penalidades aos modelos de maneira a reduzir sua flexibilidade e a tendência a incorporar o ruído dos dados. São usados para isso os hiperparâmetros, que são configurações externas aos algoritmos e que não podem ser aprendidos durante o treinamento com os dados. As diversas combinações de hiperparâmetros podem ser testadas e a melhor configuração, selecionada de acordo com o seu desempenho nas amostras e nas subamostras.

Pode-se fazer a escolha dos melhores hiperparâmetros de várias formas. Pode-se fazer uma escolha manual, treinando diversos modelos sucessivamente até encontrar a melhor combinação. A principal desvantagem dessa abordagem é consumir muito tempo de processamento, sobretudo em bases relativamente grandes. Pode-se também definir uma série de valores para guiar o algoritmo na procura pela melhor combinação em um vetor de possíveis combinações (*grid search*). A escolha das combinações possíveis pode não ser eficaz na medida em que os hiperparâmetros não têm, necessariamente, a mesma importância em diferentes bases de dados¹⁵.

¹⁵Ver Swersky, Snoek e Adams (2013) e Horváth, Mantovani e de Carvalho (2017) para problemas com essa abordagem.

Figura 6: Procedimento adotado na modelagem preditiva. A figura mostra as etapas do processo desde a seleção da Pesquisa Mundial de Valores até a predição, além das subamostras que foram retiradas da base longitudinal. O primeiro passo foi a constituição da amostra de trabalho a partir da base longitudinal do WVS. Em seguida, a base foi separada em três partes: duas amostras aleatórias para as bases de treinamento e validação, e as amostras de todas as ondas do Brasil e do WVS 7 para os testes de predição. A partir da amostra de treinamento, foram extraídos os casos completos para o estudo de benchmarking (etapas na cor cinza). Com o melhor algoritmo foi realizado o treinamento usando-se a base contendo não-respostas. As linhas tracejadas indicam com qual amostra foram feitas as avaliações de desempenho e generalização.



Uma terceira maneira é a busca aleatória (*random search*) proposta por Bergstra e Bengio (2012). A busca consiste em estimar distribuições para cada hiperparâmetro das quais são sorteadas as combinações que serão usadas no treinamento. O algoritmo é treinado usando-se as combinações aleatórias dos parâmetros até encontrar-se a combinação “ótima” de hiperparâmetros.

Existem também as técnicas conhecidas como de busca automática como as *Sequential Model-Based Optimization (SMBO)*, como, por exemplo, a otimização Bayesiana, a *Sequential Model-based Algorithm Configuration (SMAC)*, a *Tree-structured Parzen Estimator (TPE)*, ou, ainda, a busca automática por hiperparâmetros usando-se algoritmos genéticos¹⁶. No estudo de benchmark, a maioria dos algoritmos foi treinada com 50 combinações diferentes, mas em alguns casos, devido ao tempo de processamento e à quantidade de hiperparâmetros, a procura foi limitada a 5 combinações. No treinamento final do XGBoost, foram testadas ambas as técnicas.

Para selecionar a melhor combinação de hiperparâmetros e, eventualmente, comparar o desempenho de cada um dos diferentes algoritmos, é necessário escolher a métrica para medir o quão boa é a função f que foi aprendida. É importante diferenciar a métrica de desempenho da função de custo, que é usada internamente pelo modelo para aprender a relação entre *input* e *output*. A métrica de avaliação é usada para avaliar o quão boa é a “teoria” que foi aprendida com o treinamento. A tabela 3 apresenta as principais métricas usadas no aprendizado automático.

A métrica deve ser escolhida em função do problema e do objetivo. Algumas das métricas são mais aptas a selecionar algoritmos que fazem mais previsões corretas (*i.e.* Acurácia), outras são mais direcionadas a evitar erros tipo I (Falsos Positivos) ou II (Falsos Negativos). Em todos os treinamentos realizados, foram computados Acurácia, Kappa (κ), Especificidade e Sensibilidade, F1, a curva ROC e a área sob a curva ROC (AUC). A métrica usada para escolher o modelo com a melhor combinação de hiperparâmetros foi a ROC/AUC, pois se trata de uma medida que equilibra os acertos em Verdadeiros Positivos (VP) e Verdadeiros Negativos (VN). Enquanto a Acurácia recompensa algoritmos pelo quanto são bons em prever positivos e negativos de forma global, e a sensibilidade e sensibilidade

¹⁶Para SMOB e SMAC ver por exemplo Bergstra, Bardenet et al. (2011), Hutter, Hoos e Leyton-Brown (2011) ou Pham (2016); Sundhararajan, Pahwa e Krishnaswami (1998) para algoritmos genéticos.

recompensam separadamente cada um deles, a ROC é uma medida que equilibra tanto sensibilidade quanto especificidade. Para esse tipo de problema, além da ROC, o Kappa (κ) e a F1 também poderiam ter sido usadas.

Tabela 3: A tabela abaixo lista as métricas mais utilizadas no aprendizado de máquina e a sua fórmula, que envolve principalmente as relações entre verdadeiros e falsos positivos e negativos.

| Métrica | Descrição | Fórmula |
|---|--|-----------------------------|
| Acurácia (accuracy) | Mede a taxa de predições corretas, como a soma de verdadeiros positivos (VP) e negativos (VN) sobre o total de casos positivos (P) e negativos (N). | $\frac{VP+VN}{P+N}$ |
| Kappa (κ) | Compara a acurácia (A) observada com a esperada (aleatória). A variação do κ é entre 0 e 1, sendo 1 o acordo perfeito entre o observado e o esperado. | $\frac{A_o - A_e}{1 - A_e}$ |
| Sensibilidade (sensitivity ou recall) | Mede a taxa de verdadeiros positivos, como o total de (VP) sobre a soma dos verdadeiros positivos (VP) somados aos falsos negativos (FN). | $\frac{VP}{VP+FN}$ |
| Especificidade (specificity) | Mede a taxa de verdadeiros negativos (VN), como o total de verdadeiros negativos (VN) sobre o total de verdadeiros negativos (VN) somados aos falsos positivos (FP). | $\frac{VN}{VN+FP}$ |
| Precisão (precision) | Mede a taxa de predições positivas | $\frac{VP}{VP+FP}$ |
| F1 | Média harmônica entre precisão e sensibilidade. A variação do F1 é entre 0 e 1, sendo 1 a classificação perfeita. | $\frac{2*VP}{2*VP+FP+FN}$ |
| Curva de operação de recebimento (ROC) | Relação entre VP e FP. Mede a probabilidade que um modelo de classificação dê mais valor à um caso aleatório 1-Sim do que um caso aleatório 0-Não. | $\frac{VP+VN}{2}$ |

^a Os termos usados na tabela provém de diversos campos de pesquisa. Para uma descrição mais detalhada, ver Witten, 2011, cap. 5.

Problemas nos quais as classes são desbalanceadas apresentam desafios particulares. Em geral, os modelos vão tentar minimizar os erros ou o custo de uma predição errada. Cada algoritmo tem sua própria função custo, onde o ponto de partida é atribuir o mesmo custo para uma predição positiva ou negativa. No caso de uma classe desbalanceada como a de que se está tratando, pode ser integrada uma matriz de custo para cada observação. Outra maneira de fazer os modelos serem sensíveis ao custo é variar a proporção de casos (WITTEN; FRANK, 2011, p. 166).

De forma a equilibrar os custos para predições negativas e positivas, na etapa de

escolha do algoritmo optou-se por uma redução dos casos positivos em cada subamostra. Para o treinamento final, foi atribuído um peso às declarações de baixo orgulho de forma a melhorar o equilíbrio entre baixo e alto orgulho¹⁷

3.5. ANÁLISE, INFERÊNCIA E INTERPRETAÇÃO DOS RESULTADOS

De acordo com Shmueli (2010), os modelos explicativos e preditivos são regularmente confundidos. Apesar de a distinção explicativa ter um impacto importante no processo de modelagem estatística, o autor argumenta que não existe uma distinção clara no que diz respeito às diferenças entre um modelo explicativo “coerente” e um modelo preditivo “poderoso”. Segundo ele, qualquer modelo dispõe de ambas as dimensões, que deveriam ser relatadas com a mesma importância pelo pesquisador. Apesar de não necessitar de uma explicação, o desempenho de um modelo preditivo pode ser relatado à luz de teorias explicativas com propósito de aprimorá-las. Da mesma forma, o poder preditivo de um modelo explicativo também pode ser avaliado. Essa também é a posição de Hofman, Sharma e Watts (2017), já citados anteriormente, que consideram inclusive que a entrada para resolução de um problema das ciências sociais possa ser diretamente pela modelagem preditiva.

Seguindo essa premissa, foram relatadas ambas as dimensões em ambas as modelagens estatísticas. A “coerência” do modelo explicativo é dada pela interpretação dos resultados da regressão logística em relação com a teoria que permitiu criar o modelo. Na modelagem indutiva, a explicação deriva da análise da importância das variáveis.

Em ambos os casos, a análise do poder preditivo foi feita de duas maneiras nas bases do WVS e do Brasil. A primeira delas é a análise das probabilidades usando-se as curvas ROC/AUC. A área sob a curva mede o índice de concordância c , ou a probabilidade de a predição e a observação serem concordantes. A segunda avaliação foi feita com matrizes de confusão, ou tabelas de classificação na terminologia estatística, que avaliam a capacidade efetiva de predição em função de um valor de corte. O ponto de corte foi escolhido de forma a otimizar as predições para casos positivos e negativos.

A construção do conhecimento científico precisa ir além do aprendizado com os dados,

¹⁷Os melhores resultados foram obtidos usando-se um valor arbitrário igual à raiz quadrada da proporção de casos positivos e negativos da base de treinamento $\sqrt{0,13} \approx 0,36$.

do ajuste dos modelos e da predição pura e simples. É preciso elaborar inferências sobre o que foi observado.

Para isso, a interpretabilidade (ou explicabilidade) dos modelos é a chave. Tim Miller (2019b,a) define interpretabilidade como “o quanto um humano consegue entender a causa de uma decisão em um modelo”¹⁸. Para diversos autores, a interpretabilidade dos modelos é o principal desafio dos algoritmos de aprendizado de máquina (MITCHELL, Tom M, 1986; RIBEIRO; SINGH; GUESTRIN, 2016; ROBNIK-SIKONJA; KONONENKO, 2008; BAEHRENS et al., 2010).

A interpretação é necessária por diversas razões. Segundo Doshi-Velez e Kim (2017), entre as razões estão a necessidade de avaliar a imparcialidade do algoritmo, a identificação do nível proteção de dados privados, a possibilidade de verificar a robustez dos resultados preditivos, o entendimento das relações de causalidade (ou as explicações no caso presente), e, não menos importante, para que se tenha confiança nos resultados. Esta capacidade de interpretação, característica da inteligência humana, é importante para tomar decisões e fazer julgamentos (GREENWALD; OERTEL, 2017).

Na perspectiva das ciências sociais, pode-se adicionar uma outra necessidade, que é poder transformar bons resultados preditivos em inferências que sirvam para aprimorar teorias existentes sobre fenômenos sociais ou criar novas hipóteses a serem testadas.

Um modelo pode ser mais ou menos interpretável. Alguns algoritmos são transparentes, ou seja, são intrinsecamente interpretáveis pela sua simplicidade. Outros são como uma “caixa-preta”, necessitam a aplicação de técnicas específicas para interpretá-los. O algoritmo que obteve a melhor performance nos testes e foi treinado com mais dados, o XGBoost, é um desses modelos. Para modelos como esse é preciso aplicar técnicas *post-hoc*, ou agnósticas, no sentido de que sirvam independentemente do algoritmo. Existem diversos métodos desse tipo, que extraem as explicações a partir da importância global das variáveis (FISHER; RUDIN; DOMINICI, 2018); dos seus efeitos locais acumulados (APLEY; ZHU, 2016); da dependência parcial entre elas (FRIEDMAN, 2001); das expectativas condicionais individuais (GOLDSTEIN et al., 2013); de modelos substitutos locais (RIBEIRO; SINGH; GUESTRIN, 2016); usando os valores Shapley (ŠTRUMBELJ; KONONENKO, 2014); ou ainda da interação entre as

¹⁸ “the degree to which a human can understand the cause of a decision in a model”.

variáveis explicativas (FRIEDMAN; POPESCU, 2008).¹⁹

Neste trabalho, uma vez feita a avaliação da capacidade preditiva do modelo treinado, foi realizada uma interpretação global e local dos resultados de forma a desvendar a teoria que foi aprendida.

A importância global foi computada de duas maneiras. A primeira usando-se os métodos específicos do XGBoost que calcula três índices. O índice mais importante é o Ganho, que se refere à melhoria em termos de acurácia trazida por cada variável em cada galho da árvore de decisão em que se encontra. A técnica utilizada é próxima do *Mean Decrease Impurity (MDI)* proposto por (LOUPPE et al., 2013). A importância global, nesse caso, é uma nota que indica o quão útil foi uma determinada variável na construção de cada árvore de decisão do modelo. Quanto mais uma variável é usada, maior sua importância²⁰.

Algoritmos com base em árvores de decisão e florestas aleatórias têm problemas específicos no que diz respeito à identificação da importância global. Um deles é que modelos desse tipo, depois de treinados, tendem a inflar a importância de algumas variáveis, especialmente as contínuas ou categóricas de alta cardinalidade (STROBL et al., 2007). Isso pode acontecer devido ao fato de essas variáveis terem mais valores possíveis de serem usados nas separações dos ramos (comparados com uma variável binária, por exemplo).

De modo a verificar esses vieses e verificar a robustez da importância do XGBoost, foi aplicada uma segunda técnica de importância global que consiste em fazer permutações de forma a encontrar quais variáveis importam mais para as previsões. As técnicas são

¹⁹Ver, por exemplo, Adadi e Berrada (2018), Carvalho, Pereira e Cardoso (2019) e Du, Liu e Hu (2019) para os diversos métodos e técnicas. Ver Payrovnaziri et al. (2020) para uma tentativa de classificação dos diferentes métodos e técnicas aplicados. Os autores identificam cinco categorias: 1) interação entre variáveis explicativas (*features*); 2) mecanismo de atenção; 3) redução da dimensão dos dados; 4) destilação do conhecimento e extração de regras; e 5) modelos intrinsecamente interpretáveis. Esses métodos podem ainda ser classificados em intrínsecos ou *post-hoc*, globais ou locais, e específicos ao modelo ou agnósticos. Ver também Molnar (2019).

²⁰O ganho não permite afirmar se uma variável deve ou não estar presente em uma previsão. Outros dois indicadores, Frequência e Cobertura, ajudam a interpretar o ganho. Frequência, ou peso, uma medida alternativa do Ganho, trata-se do número de vezes que uma variável explicativa é usada para separar os dados em todas as árvores. A Cobertura é também o número de vezes que uma variável é usada para separar os dados, mas normalizada pelo número de observações que são separadas. A Cobertura indica o número relativo de observações que tem relação com cada variável (em cada nó). Por exemplo, para 100 observações com 4 variáveis e 3 árvores, supondo que a variável x seja usada para decidir a separação nos nós para 10, 5 e 2 observações nas árvores 1, 2 e 3, respectivamente, a cobertura contabiliza para a variável x $10 + 5 + 2 = 17$ observações. Isso é calculado para todas as variáveis e a cobertura 17 é expressa como a porcentagem entre todas as variáveis.

baseadas em Breiman (2001b), e, mais recentemente, Fisher, Rudin e Dominici (2018), e são chamadas também de *Mean Decrease Accuracy (MDA)*. A importância das variáveis explicativas é calculada computando-se o aumento do erro da predição de um modelo depois de permutar uma variável explicativa. Se a variável é importante, o erro deve aumentar, pois o modelo se apoia nela para fazer a predição. Se o erro não muda, a variável pode ser ignorada para predição.

No método aplicado, proposto por Fisher, Rudin e Dominici (2018), a redução no ajuste é a diferença da AUC do modelo final e o resultado permutando-se cada uma das variáveis explicativas. Valores maiores significam maiores perdas e uma maior importância da variável. De modo a reduzir o erro introduzido pelas permutações, a importância contabilizada é uma média de 10 permutações aleatórias.

Em geral, os resultados em termos de importância não podem ser analisados pelo valor absoluto (Ganho ou *1-Log Loss* nos casos acima), mas apenas em sua ordem.

Foi também aplicada uma técnica para verificar a importância local das variáveis. A premissa dessa técnica é que a relação entre variáveis em observações próximas é similar e se pode ajustar localmente um modelo interpretável (e.g. regressão linear). Analisando-se essas estimativas, é possível ter uma visão das variáveis que têm influência nas predições individuais e podem ajudar a confirmar que a importância global das variáveis faz sentido e que se pode confiar no modelo (RIBEIRO; SINGH; GUESTRIN, 2016). Enquanto a importância global das variáveis é calculada com base nas validações cruzadas usando-se a amostra de treinamento, a importância local foi calculada nas bases do WVS7 e Brasil.

Na forma de uma classificação da mais influente à menos influente, as importâncias globais e locais facilitam a elaboração de possíveis explicações a partir de modelos que visam à predição em primeiro lugar, e, conseqüentemente, ampliam a possibilidade do uso dessas técnicas nas ciências sociais.

Elaborar novas teorias com base nessas análises tem limitações. Em primeiro lugar, a crença na importância global das variáveis é estabelecida somente pelo poder preditivo do modelo, e não por uma teoria estatística subjacente, como é o caso na modelagem explicativa (BREIMAN, 2001b). Além disso, a importância das variáveis pode ser diferente em cada modelo testado e varia com as subamostras testadas, mesmo quando os modelos

têm desempenho similar em termos preditivos — um exemplo do Efeito Rashomon²¹. A importância local também apresenta diversos problemas como a definição da vizinhança com a qual será ajustado o modelo local (ALVAREZ-MELIS; JAAKKOLA, 2018; LAUGEL et al., 2018). Apesar disso, essas formas de inferência parecem ser, por enquanto, a única maneira de abrir a caixa preta da aprendizagem de máquina. É com base nelas que se podem induzir teorias no sentido das ciências sociais.

3.6. AMBIENTE DE DESENVOLVIMENTO E PACOTES

O trabalho foi desenvolvido na linguagem de programação *R*, versão 3.5.1/3.5.2 (R CORE TEAM, 2018; IHAKA; GENTLEMAN, 1996; R CORE TEAM, 2016; R DEVELOPMENT CORE TEAM, 2004) no ambiente de desenvolvimento *RStudio* (RSTUDIO TEAM, 2019), versão 1.2.5033. Foram usados um laptop, com sistema operacional OSX (Processador 2,5 GHz Intel Core i5, 8 Go 1600 MHz DDR3), e, para o aprendizado de máquina, um desktop, com sistema operacional Windows 10 (Processador 3,41 GHz Intel Core i7, 64 GB 2133 MHz DDR3).

A biblioteca *Classification and Regression Training (CARET)* (KUHN, 2019) foi usada para implementar o aprendizado de máquina. Essa biblioteca permite fácil comparação entre algoritmos devido à facilidade em padronizar validações cruzadas, em definir o método de busca por hiperparâmetros, em definir as métricas de desempenho, etc. O pacote é, basicamente, uma função destinada a chamar uma ou mais funções (*wrapper*), e permite a programação em alto nível para operações simples de aprendizado automático aplicável a diversos problemas. Foram usados os seguintes pacotes contendo os os algoritmos de aprendizado de máquina: redes neurais (BERGMEIR et al., 2019; MOUSELIMIS; GOSSO, 2018; RIPLEY, 2020), classificação usando regras (HORNIK et al., 2020), métodos Bayesianos (PACKAGE), 2018), árvores de decisão (KUHN; QUINLAN, 2020; HOTHORN; HORNIK; STROBL et al., 2020; THERNEAU; ATKINSON, 2019), Florestas Aleatórias (BREIMAN; CUTLER et al., 2018), boosting gradiente (GREENWELL et al., 2019), regressão linear (FRIEDMAN; HASTIE et al., 2019), análise discriminante regularizada (ROEVER et al., 2020), e, para o treinamento mais

²¹O Efeito Rashomon tem origem no filme *Rashomon* de Akira Kurosawa, realizado em 1950. No filme, diversas testemunhas oculares têm diferentes interpretações do que viram, e, assim, não se pode saber o que de fato aconteceu.

extensivo, o XGBoost (2020).

A tese foi escrita em *Rmarkdown* (2020) e *LaTeX* usando os pacotes *bookdown* (2020a) e *knitr* (XIE, 2020b). A manipulação dos dados foi realizada usando as rotinas desenvolvidas na gramática e lógicas do *Tidyverse* (2020; 2020; 2018; 2019; 2019; 2019a; 2020; 2020; 2019b). O *ggplot* (2020), *DiagrammeR* (2020) e *kableExtra* (2019) foram usados para realizar as figuras e as tabelas. O programa *Zotero*²² foi usado para gerir as citações e inseridas no *RMarkdown* com auxílio da extensão *Better BibTex*²³.

²²https://www.zotero.org/support/credits_and_acknowledgments#about_zotero

²³<https://retorque.re/zotero-better-bibtex/>

4. RESULTADOS

Nos capítulos precedentes, foram apresentados os conceitos de inteligência artificial e aprendizado de máquina. Foi também descrita a natureza exploratória da tese que busca integrar essas ferramentas na produção de conhecimento das ciências sociais, usando o patriotismo como caso de estudo. No capítulo anterior, foi apresentada a metodologia usada neste trabalho.

Esse capítulo, descreve os resultados, na mesma ordem do capítulo anterior. As primeira e segunda partes apresentam os resultados da modelagem explicativa, e, em seguida, os resultados obtidos no momento indutivo com as técnicas de aprendizado de máquina. A terceira parte do capítulo descreve os resultados das técnicas de extração do conhecimento obtido durante o aprendizado.

4.1. RESULTADOS DA MODELAGEM EXPLICATIVA

O objetivo da modelagem estatística seguindo o caminho explicativo é “provar” uma ou mais teorias ou hipóteses. Cada uma das teorias sobre o patriotismo apresentadas no capítulo dois permite deduzir possíveis explicações para o fenômeno, que poderiam ser observadas na forma de associações entre a variável resposta e as variáveis explicativas selecionadas na base do WVS.

Para ilustrar a abordagem dedutiva, foi escolhida uma dessas teorias fundamentada na tradição da cultura política. Na perspectiva dominante da pesquisa nessa linha, o patriotismo é parte de uma síndrome de valores tradicionais em declínio e que pode ser medido pelo nível de orgulho que os indivíduos sentem pela sua nacionalidade. Segundo Ronald Inglehart (1971, 1977, 1997; 2005; 2018), que foi quem propôs originalmente a tese da modernização e do pós-materialismo, o declínio desses valores tradicionais é explicado pelo efeito da afluência

na geração pós-guerra, principalmente nos países desenvolvidos, que levou à modificação de valores em duas dimensões principais. Por um lado, valores emancipativos cresceram, juntamente com tendências antidiscriminatórias e igualitárias, predispondo os indivíduos ao individualismo e ao altruísmo. De outro, valores seculares cresceram em direção à uma maior abertura aos estrangeiros e à cooperação interpessoal.

A teoria de Inglehart tem como premissas básicas a importância das crenças e valores como objetivos pessoais e sociais e a influência da socialização dos indivíduos na formação desses valores. O autor sustenta a tese em duas hipóteses principais. A hipótese da escassez, inspirada na hierarquia das necessidades de Maslow, defende que as prioridades da ação humana são resultado das condições socioeconômicas. A hipótese da socialização, por sua vez, diz respeito à maneira como os valores são internalizados e distribuídos na sociedade. Segundo essa hipótese, os valores básicos de um indivíduo são resultados das condições socioeconômicas presentes no seu período de formação. De forma resumida, à medida que as sociedades se tornam mais afluentes, as necessidades dos indivíduos diminuem e os valores internalizados mudam de questões materiais à pós-materiais.

Inglehart (1989, 1971) e Abramson e Inglehart (1995) propuseram um índice do pós-materialismo usando três perguntas. Cada uma delas tem quatro alternativas que permitem identificar as principais prioridades pessoais e para a nação. O quadro 4 indica as três perguntas e as transformações que foram realizadas para o seu uso na regressão.

O orgulho pela pátria é em si um valor tradicional. Nessa perspectiva, pode-se esperar que valores e prioridades materiais guardem uma correlação positiva com o nível de orgulho que uma pessoa tem pela sua nacionalidade. Espera-se que indivíduos que tenham escolhido as alternativas materialistas destas perguntas se declarem também mais orgulhosos. Ainda segundo a teoria de Inglehart, com base na observação dos efeitos sentidos no pós-guerra nos países mais desenvolvidos, a afluência propicia a mudança de valores em direção a valores de auto expressão e racionais-seculares. Indivíduos mais afluentes, em sociedades que puderam investir mais em educação, teriam maior chances de se emancipar com relação ao estado-nação, e, em consequência, serem menos orgulhosos deles. Nessa perspectiva, espera-se, também, que a renda tenha relação inversa com orgulho nacional. Enfim, sociedades pós-materiais tendem a ser também mais globalizadas e as pessoas que se socializaram mais recentemente tenderiam a ser menos orgulhosas.

Tabela 4: Variáveis usadas na regressão logística. A primeira coluna se refere às perguntas do WVS e as alternativas em termos de prioridades. Cada uma delas corresponde ao construto teórico materialista ou pós-materialista. A última coluna se refere ao nome dado à variável na equação representando os resultados. Para cada variável estão indicadas as frequências.

| Perguntas | Alternativas (Prioridades) | Construto Teórico | Variável |
|---|--|-------------------|-------------|
| Objetivos para a pessoa: primeira escolha (E003 - Aims of respondent: first choice) | 1. Manter a ordem no país (41,2% das respostas) | Materialismo | Ordem |
| | 2. Dar mais voz às pessoas em decisões do governo (21,1%) | Pós-Materialismo | |
| | 3. Lutar contra a alta dos preços (23%) | Materialismo | Preços |
| | 4. Proteger a liberdade de expressão (10,4%) | Pós-Materialismo | |
| Objetivos para o país: primeira escolha (E001 - Aims of country: first choice) | Não-respostas (4,1%) | | |
| | 5. Alto nível de crescimento econômico (40,8%) | Materialismo | CEcon |
| | 6. Forte poder de defesa (7,6%) | Materialismo | Defesa |
| | 7. Pessoas devem participar mais da tomada de decisões (13,4%) | Pós-Materialismo | |
| Mais Importante: primeira escolha (E005 - Most important: first choice) | 8. Tentar fazer com que as cidades e campo sejam mais bonitos (4,9%) | Pós-Materialismo | |
| | Não-respostas (33,1%) | | |
| | 9. Uma economia estável (37,7%) | Materialismo | EconEst |
| | 10. O progresso em direção à uma sociedade mais impessoal e humana (11,8%) | Pós-Materialismo | |
| | 11. Idéias contam mais do que dinheiro (5,7%) | Pós-Materialismo | |
| | 12. A luta contra o crime (12,8%) | Materialismo | Crime |
| | Não-respostas (31,8%) | | |
| Escala de Renda | Numérica de 1-10 (média ~4,7) | Afluência | Renda |
| | Não-respostas (29,9%) | | |
| Idade | Numérica 16-103 (média ~42 anos) | Globalização | Idade |
| | Não-respostas (0%) | | |
| Orgulho de ser Brasileiro | Binária 0-1 Alto (88%) e Baixo Orgulho (11,7%) | Patriotismo | Patriotismo |
| | Não-respostas 5% | | |

^a As prioridades pós-materialistas de cada questão foram agrupadas em uma categoria que foi usada como valor de referência na regressão

Em termos mais específicos, a hipótese que foi escolhida para ser testada é que o patriotismo tem uma relação estatisticamente significativa e positiva com as prioridades materialistas de uma pessoa, e significativa e negativa com uma idade menos avançada e maior nível de renda.

Este modelo foi ajustado na amostra de treinamento, a mesma amostra usada para o aprendizado de máquina, excluindo-se as observações incompletas que foram retiradas automaticamente da regressão ($n = 234.580$). Os resultados de uma regressão logística, com os coeficientes apresentados no logaritmo natural das chances, são representados pela equação (4.1) abaixo.

$$\begin{aligned}
 f(\textit{Patriotismo}) = & 1,76 + 0,14 * CEcon + 0,51 * Defesa \\
 & + 0,25 * Ordem + 0,00 * Precos \\
 & + 0,14 * EconEst + 0,34 * Crime \\
 & + 0,14 * Renda + 0,78 * Idade
 \end{aligned}
 \tag{4.1}$$

Os resultados de uma regressão logística podem ser apresentados usando a razão de chances para que sejam visualizados de forma mais clara¹. A figura 7 apresenta a razão de chances para cada variável, com intervalo de confiança de 95%. Os coeficientes foram normalizados de forma a poder compará-los diretamente (GELMAN, 2008).

Os resultados do ajuste do modelo indicam que a hipótese do pós-materialismo encontra sustentação empírica. Todas as prioridades individuais e nacionais de natureza materialista, com exceção da luta contra preços altos, têm relação estatisticamente significativa ($p < 0$) e positiva (razão de chances > 1) com o um maior orgulho pela nação.

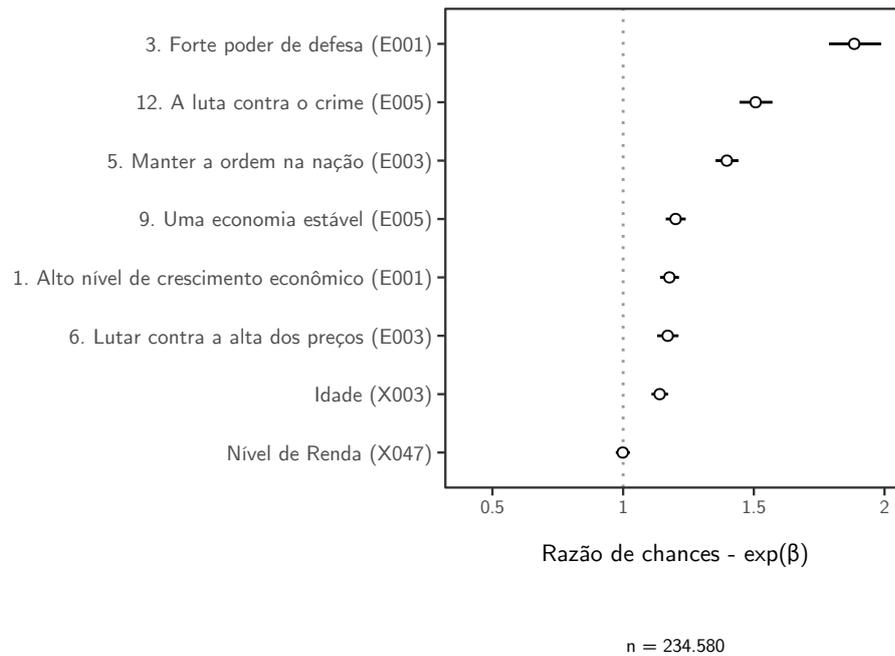
Claramente, as questões relacionadas à defesa, à criminalidade e à ordem contribuem mais para aumentar as chances de alguém também se declarar orgulhoso. Interpretando os resultados pela razão de chances, pode-se notar que, mantendo todas as outras variáveis constantes, as chances de que um indivíduo que priorize o poder de defesa também se declare orgulhoso da sua nacionalidade são 1,7 maiores (70%) em comparação com um indivíduo que priorize questões pós-materiais (da mesma série de perguntas). Além disso, pode-se notar uma certa distância entre as prioridades materialistas de cunho econômico e social. Essas são mais importantes do que aquelas.

Outro resultado esperado era uma relação significativa, mas negativa, com o nível de

¹O modelo pode também ser avaliado em termos da própria capacidade preditiva, tanto pelas probabilidades geradas quanto pelas predições. Essa análise é feita abaixo em comparação com os resultados do modelo indutivo.

renda. Surpreendentemente, o nível de renda quase não alterou as probabilidades de alguém se declarar mais orgulhoso. Confirmou-se o efeito previsto no que diz respeito à idade. Pessoas mais velhas têm mais chances de se declararem mais orgulhosas.

Figura 7: Resultados da regressão logística apresentado pela razão de chances e com intervalos de confiança. O valor da razão de chances foi calculado após normalização dos coeficiente



Além da avaliação da significância estatística de cada variável, o modelo de regressão produzido pode ser avaliado usando estatísticas globais e de qualidade do ajuste (*goodness-of-fit* GOF). A avaliação global do modelo se dá em função da melhora que o modelo traz com relação ao modelo com somente o Intercepto (modelo nulo). No modelo nulo, todas as predições são para a classe majoritária. Os testes *Likelihood Ratio*, Teste de Wald e de *Score* revelaram que o modelo é um avanço, mesmo que limitado, com relação ao modelo nulo. Entretanto, os testes de qualidade mostraram que o modelo tem pouca capacidade explicativa².

²O teste de Hosmer–Lemeshow revelou que não há evidência para descartar a hipótese de um modelo bem ajustado. Entretanto, os resultados obtidos com os testes como os Pseudo R^2 (e.g. Cox & Snell, Nagelkerke, Tjur) foram baixos indicando que o modelo é pouco ajustado aos dados. Além disso, foram obtidos resultados muito similares com as múltiplas bases imputadas. A imputação gerou ajustes e intervalos de confiança ligeiramente menores indicando, talvez, alguma eficiência adicional quando são usados todos os dados disponíveis. Entretanto, a qualidade do ajuste foi bastante menor. Os testes pseudo- R^2 globais tiveram resultados similares, mas o teste de Hosmer–Lemeshow, por exemplo, não descartou a hipótese nula de um modelo não bem ajustado. As análises pelos resíduos também mostraram um mau ajuste com os dados imputados.

Claramente, outras variáveis podem ser incorporadas ao modelo. As questões usadas na regressão foram somente algumas das consequências observáveis da teoria do pós-materialismo. O próprio Inglehart identifica várias outras. Como a ideia aqui não é melhorar esta regressão usando métodos tradicionais (i.e. voltar a teoria), na verdade a tese propõe uma alternativa usando aprendizado de máquina, para fins de comparação, assumimos o modelo com essas limitações.

4.2. RESULTADOS DO APRENDIZADO DE MÁQUINA

4.2.1. Estudo comparativo

No estudo comparativo, foram testados no total 20 algoritmos. Levando em consideração todas as combinações de hiperparâmetros, foram testados 492 modelos diferentes (entre 5 e 50 combinações de hiperparâmetros para cada algoritmo). Para cada um deles, foram calculadas as métricas principais do aprendizado de máquina como a Acurácia, Especificidade e Sensibilidade, curva ROC/AUC, F1 e κ . A comparação foi feita usando a melhor combinação de hiperparâmetros de cada algoritmo, definidos após uma procura aleatória, e selecionados pelo seu desempenho usando como métrica a curva ROC.

A lista dos algoritmos testados, seguindo a classificação proposta por Fernández-Delgado et al. (2014), pode ser encontrada no Apêndice C. Os resultados da comparação, bem como o tempo de processamento para cada algoritmo, estão expressos na Tabela 5 e na Figura 8. Os valores, tanto na tabela quanto na figura, estão ordenados pelo desempenho na métrica AUC-ROC.

Pela leitura da tabela, podemos observar que no caso do algoritmo xgb, por exemplo, a acurácia é de 0.72. Isso significa que em 72% dos casos são identificados corretamente as classes alto ou baixo orgulho. Pode-se notar, igualmente, que nenhum dos algoritmos chegou perto do que se pode chamar de acurácia de base. Essa acurácia seria atingida se o resultado mais comum (Alto Orgulho) fosse atribuído a todos os casos, sem exceção. No caso do orgulho nacional o caso mais comum, o Alto Patriotismo (as pessoas que se declararam “Muito Orgulhosos” ou “Orgulhoso”), representa aproximadamente 88% das respostas (ver Tabela 1).

Tabela 5: Resultados da comparação dos algoritmos na base de casos completos. A tabela apresenta os resultados em termos de Acurácia, κ , ROC/AUC, Sensibilidade e Especificidade. O valor final computado para cada métrica é a média das 100 iterações da validação cruzada, em 10 subamostras, repetidas 10 vezes.

| Algoritmo | Acurácia | κ | AUC-ROC | Sens. | Esp. | Hiper. | Processamento |
|-----------|----------|----------|---------|-------|------|--------|------------------------|
| xgb | 0.72 | 0.43 | 0.79 | 0.72 | 0.72 | 50 | 92836s (1.07 days) |
| rf | 0.71 | 0.41 | 0.78 | 0.71 | 0.71 | 5 | 62813s (17.45 hours) |
| xgbDART | 0.71 | 0.40 | 0.78 | 0.70 | 0.71 | 5 | 82664s (22.96 hours) |
| c50 | 0.71 | 0.40 | 0.78 | 0.71 | 0.70 | 50 | 12076s (3.35 hours) |
| gbm | 0.70 | 0.40 | 0.77 | 0.70 | 0.71 | 5 | 118924s (1.38 days) |
| avNNET | 0.69 | 0.37 | 0.76 | 0.68 | 0.71 | 5 | 197910s (2.29 days) |
| nnet | 0.69 | 0.37 | 0.76 | 0.68 | 0.70 | 50 | 252800s (2.93 days) |
| rda | 0.69 | 0.37 | 0.25 | 0.69 | 0.69 | 50 | 41121s (11.42 hours) |
| glmnet | 0.68 | 0.35 | 0.74 | 0.67 | 0.68 | 50 | 72165s (20.05 hours) |
| bglm | 0.68 | 0.35 | 0.74 | 0.68 | 0.68 | 1 | 1035s (17.25 minutes) |
| mlpML | 0.68 | 0.34 | 0.74 | 0.65 | 0.70 | 50 | 199919s (2.31 days) |
| Log | 0.67 | 0.33 | 0.73 | 0.66 | 0.68 | 5 | 9655s (2.68 hours) |
| nb | 0.67 | 0.33 | 0.73 | 0.64 | 0.69 | 5 | 15417s (4.28 hours) |
| ctree2 | 0.66 | 0.31 | 0.72 | 0.67 | 0.65 | 50 | 87486s (1.01 days) |
| rpart | 0.67 | 0.33 | 0.71 | 0.68 | 0.66 | 50 | 1455s (24.25 minutes) |
| JRip | 0.67 | 0.34 | 0.70 | 0.67 | 0.68 | 5 | 28050s (7.79 hours) |
| elmNN | 0.63 | 0.24 | 0.67 | 0.60 | 0.64 | 50 | 41003s (11.39 hours) |
| OneR | 0.61 | 0.23 | 0.62 | 0.67 | 0.57 | 1 | 223s (3.72 minutes) |
| bart | 0.69 | 0.37 | 0.24 | 0.69 | 0.69 | 5 | 117733s (1.36 days) |

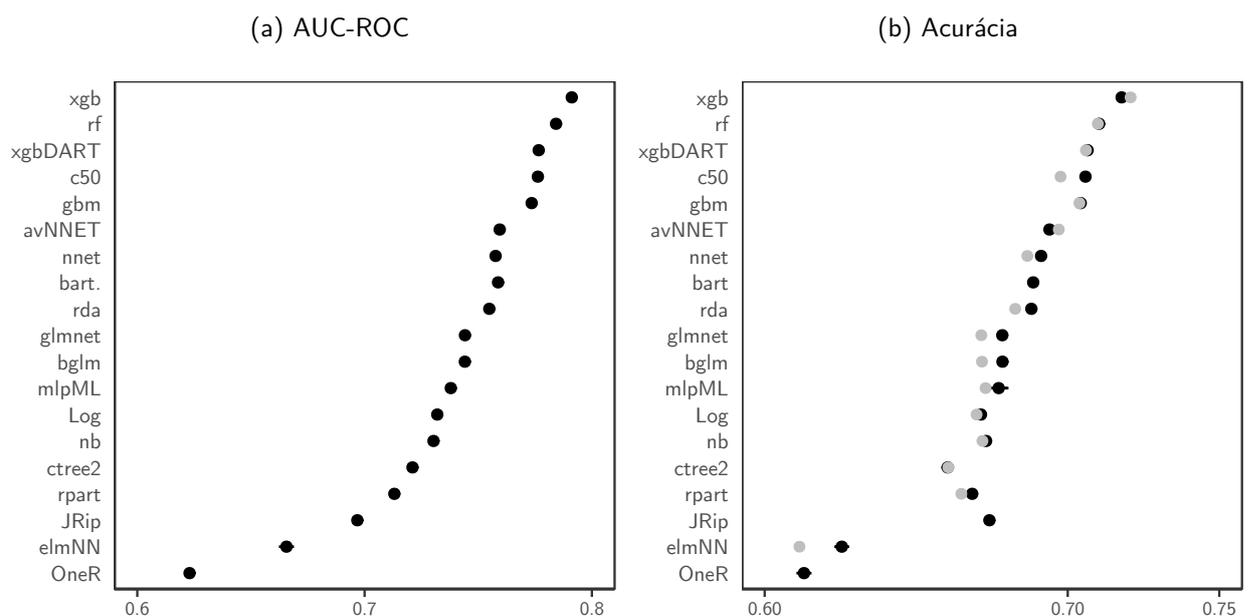
^a A coluna hiperparâmetros se refere ao número de combinações de hiperparâmetros testados

^b O tempo de processamento se refere ao tempo total usado para estimar todos as combinações com o algoritmo (Computado usando R3.5.2, Processador i7-6700 CPU @ 3.40 GHz, 64GB RAM)

Entre os algoritmos testados, os que tiveram melhor desempenho foram os classificadores *ensemble*. Não se trata de uma surpresa, pois estes algoritmos têm um registro de bons resultados na literatura (e.g. MARQUÉS; GARCÍA; SÁNCHEZ, 2012; TSAI; HSU; YEN, 2014; MACLIN; OPITZ, 1999). O classificador melhor colocado foi o XGBoost, ou *eXtreme Gradient Boosting* (“xgb” na tabela 5), com uma acurácia média de 68% (+-1,2%), e a máxima em 72%, entre todos as combinações de hiperparâmetros. Este também foi o algoritmo que atingiu o melhor desempenho usando a curva ROC como métrica. Em seguida, o algoritmo melhor colocado foi o Floresta Aleatória (rf), que atingiu uma acurácia média de 70,5% (+-0,05%) e máxima de 71%. O Floresta Aleatória foi também o segundo melhor colocado medindo-se a área abaixo da curva ROC.

Para verificar se o modelo tem sobreajuste ou não, ou seja, se ele tem capacidade de generalização em dados novos, foram comparados os resultados em termos de acurácia na amostra de casos completos da base de validação ($n = 13\ 497$). Na figura 8(b), pode-se constatar que o XGBoost teve acurácia um pouco maior na amostra de validação. Os algoritmos Floresta Aleatória (rf), Boosting Gradiente (gbm) e Boosting Gradiente Extremo DART (xgbDART) tiveram praticamente o mesmo desempenho em ambas as amostras. Os demais tiveram desempenho inferior na base de validação, com destaque para o C5.0 que teve boa performance na amostra de treinamento, mas um pouco menor na amostra de validação, indicando algum sobreajuste.

Figura 8: Resultados dos algoritmos testados em termos de Acurácia e ROC. Os modelos estão listados no eixo y pela ordem dos resultado da área abaixo da curva ROC. A Figura (a) mostra no eixo x a a área abaixo da curva ROC dos diferentes modelos; a figura (b) mostra no eixo x Acurácia, tanto na base de treinamento quanto de validação.



A partir desses resultados, foi possível identificar os algoritmos com melhor desempenho. Entretanto, a diferença entre os primeiros colocados foi ainda pequena. Foram aplicados testes estatísticos nos dois melhores colocados, o XGBoost (xgb) e o Floresta Aleatória (rf), que revelaram uma diferença estatisticamente significativa entre eles, em favor do XGBoost.

4.2.2. Aprendizado com o XGBoost

Como foi visto acima, o algoritmo XGBoost teve o melhor desempenho no estudo comparativo. Com este algoritmo, foi realizado um treinamento mais extenso.

Treinamento, em aprendizado de máquina, significa encontrar os melhores hiperparâmetros ou regularizações na forma de restrições e penalidades. Estes hiperparâmetros ajudam a estimar os parâmetros da função f que gerariam menos sobreajuste no momento da predição. A busca pelos hiperparâmetros do XGBoost foi feita de forma aleatória³. O tempo total usado para treinamento foi de 307,56 horas para o ajuste das 25 combinações de seis hiperparâmetros possíveis. A melhor combinação foi escolhida usando a área abaixo da curva ROC.

A performance do modelo treinado nas amostras de treinamento e validação pode ser verificada na tabela 6. Nesta tabela, pode-se notar a tendência ao sobreajuste do algoritmo. Quando se calcula a média obtida ao longo das validações cruzadas, os resultados se aproximam dos resultados que foram contabilizados na amostra de validação.

Tabela 6: Performance do Extreme Boosting usando a amostra de treinamento e validação

| Métrica | Base de Treinamento | | Base de Validação |
|----------------|------------------------------|--|------------------------------|
| | Desempenho com melhor ajuste | Desempenho médio nas validações cruzadas | Desempenho com melhor ajuste |
| Acurácia | 0.900 | 0.822 | 0.818 |
| Sensibilidade | 0.893 | 0.856 | 0.850 |
| Especificidade | 0.952 | 0.566 | 0.578 |
| AUC | 0.976 | 0.811 | 0.808 |
| Kappa | 0.637 | 0.329 | 0.328 |
| F1 | 0.940 | 0.894 | 0.892 |

^a Os resultados de desempenho são as médias das 10 sub-amostras de validação cruzada (10 x 10 = 100 valores ao total), usando os melhores hiperparâmetros na métrica ROC

³Os hiperparâmetros também foram testados dois a dois, sendo que os dois mais importantes, o *nrounds*, o número de árvores de decisão no modelo final, e, o *max_depth*, sua profundidade (o número de perguntas que o modelo respondeu por sim/não até chegar à decisão final ou folha), foram deixados para o final para uma busca com um *grid* mais restrito. Esta técnica gerou um modelo com maior sobreajuste do que a busca aleatória.

Podemos observar no caso do XGBoost, por exemplo, que se permitirmos que o modelo classifique 50% dos casos como falsos positivos, o modelo identifica praticamente 90% dos casos positivos (verdadeiro positivo). Para um valor análogo de falsos positivos, no caso da regressão, encontramos um pouco mais de 50% de verdadeiros positivos. Essa interpretação, aplicada para cada valor do eixo x , nos permite concluir que quanto maior a área sob a curva ROC maior é o poder de predição e discriminação do modelo.

4.3. PREDIÇÃO NAS AMOSTRAS WVS7 E BRASIL

Para verificar a capacidade de generalização de cada modelo, foram utilizadas as amostras do WVS7 e do Brasil, que, como não foram usadas no treinamento, são consideradas como “futuras”, ou ainda não vistas.

Existem duas maneiras de estimar a capacidade de generalização. A primeira delas é usando as probabilidades geradas pelos modelos, a segunda, pelas predições propriamente ditas, uma vez escolhido um valor de corte.

4.3.1. Probabilidades

Como vimos acima, o modelo treinado apresentou bom desempenho preditivo tanto na base de treinamento quanto na base de validação, mesmo com algum sobreajuste. Pelo valor baixo do κ , e em menor escala do F1, o modelo parece não ter necessariamente uma boa capacidade de discriminar, ou classificar, corretamente tanto ocorrências positivas quanto negativas. Em termos de probabilidades, isso significa que o modelo não consegue, simultaneamente, atribuir probabilidades altas de ocorrência às observações positivas, e probabilidades baixas às observações negativas.

Uma das maneiras de avaliar a capacidade geral de generalização e de discriminação do modelo é usando a curva ROC, e a área abaixo dessa curva AUC, também chamada de índice de concordância c . A curva ROC é o resultado da combinação das probabilidades para os indicadores de sensibilidade e especificidade. Para lembrar, sensibilidade é definida como a probabilidade de uma predição positiva para uma observação realmente positiva (Taxa de Verdadeiros Positivos); especificidade é a probabilidade de uma predição negativa para uma

4.3.2. Predições

A predição usando o modelo depende da definição de um ponto de corte. Este ponto depende do “custo” atribuído a uma boa ou má classificação. Se o classificador for usado, por exemplo, para definir sentenças de morte, ou diagnósticos que levam a tratamentos pesados, falsos positivos levam inocentes à morte, ou submetem pessoas a tratamentos desnecessários. Nesses casos, o valor de corte depende de uma escolha substantiva de forma a minimizar os Falsos Positivos, ao custo de uma menor taxa de identificação de Verdadeiros Positivos.

Não existe um risco desse tipo no caso de um erro de classificação de patriotas. Entretanto, de forma a comparar as predições de cada modelo usando as mesmas bases, optou-se por encontrar um mesmo valor de corte “ótimo” para ambos os modelos, definido como o ponto onde o custo de verdadeiros e falsos positivos é o mesmo. Assim a capacidade de discriminar entre verdadeiros positivos e verdadeiros negativos é a maior possível em ambos.

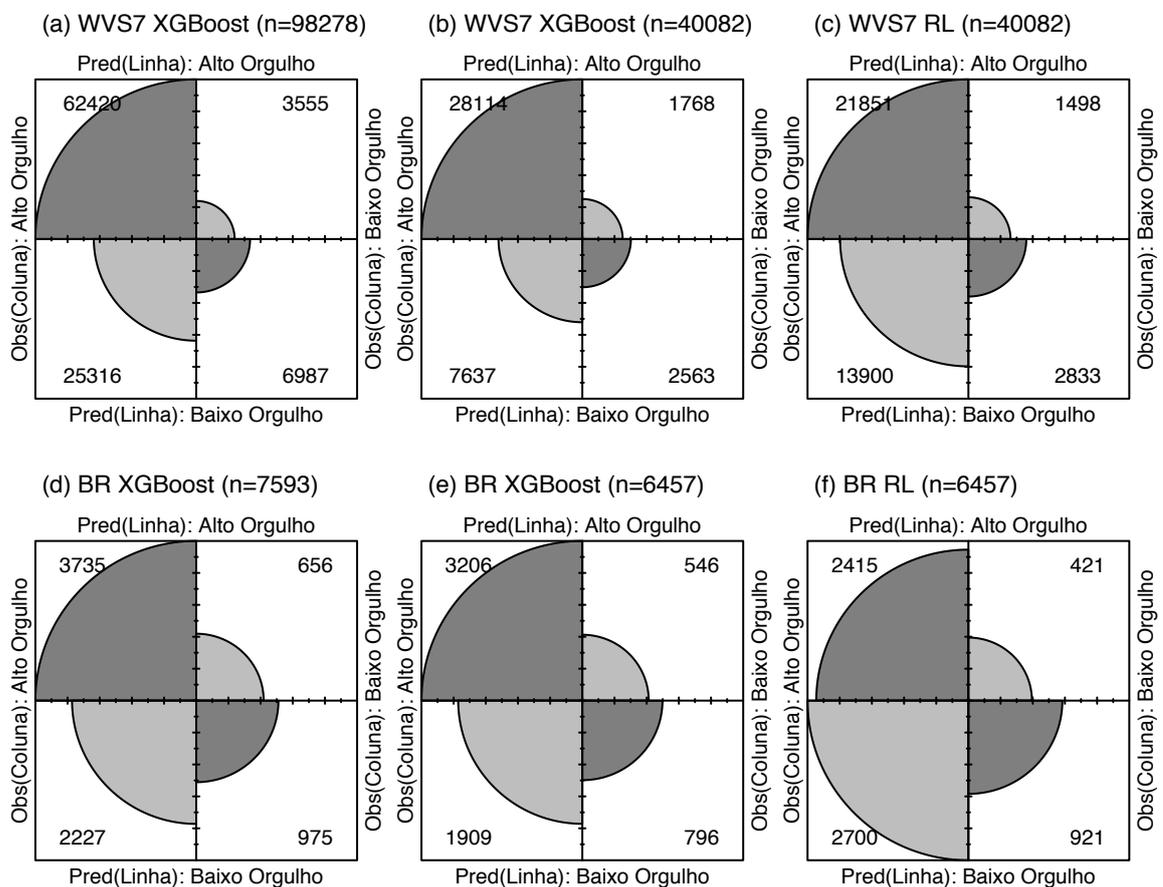
O ponto de corte é definido com base na curva ROC. A curva traça as taxas de verdadeiros positivos e falsos positivos em diferentes limiares de classificação. Existem diferentes limiares, ou pontos de corte. Cada um deles corresponde à escolha entre discriminar verdadeiros positivos e verdadeiros negativos.

Existem duas medidas que podem ser usadas para definir o limiar, ou ponto de corte “ótimo”. A primeira delas é obter o ponto menos distante do ponto superior esquerdo do gráfico. Esta medida pode gerar diferentes custos em diferentes situações (PERKINS; SCHISTERMAN, 2006). Assim, optou-se por outra medida, o J de Youden (1950). Graficamente, a estatística J representa a maior distância entre o limiar escolhido e a linha diagonal de 45° do gráfico.

Com os limiares estabelecidos para ambos os modelos, a avaliação da capacidade de predição é feita usando matrizes de confusão, ou “tabelas de contingência”, na terminologia estatística. Uma matriz de confusão é um resumo das predições de um modelo para um problema de classificação. A matriz agrega as predições para cada observação, cruzando as predições e os valores reais. As colunas da tabela representam as classes realmente observadas e as linhas representam a predição para cada classe. Cada célula contém a contagem do número de observações que respeitam as condições da linha e da coluna (*i.e.* os valores de VP, VN, FP e FN).

As previsões usando os limiares ótimos são representados na figura 11. O XGBoost é capaz de fazer previsões levando em consideração observações onde existam não-respostas em alguns itens. A regressão logística, por sua vez, descarta as observações quando existem não-respostas em alguma das variáveis explicativas. De modo a melhor comparar, foram elaboradas duas matrizes de confusão para o XGBoost, uma delas descartando-se as observações para as quais não foi possível fazer a previsão usando o modelo de regressão logística. Isso permite comparar os dois métodos com base na mesma amostra⁴.

Figura 11: Matrizes de Confusão para as amostras do WVS7 e Brasil usando ponto de corte definido pelo J de Youden. Para o XGBoost, existem duas matrizes, uma delas descartando as previsões para as quais não há previsão usando o modelo de regressão logística (RL). A área do semicírculo é proporcional à frequência de ocorrências em cada célula.



Na representação gráfica usada na figura 11, a área do semicírculo é proporcional à

⁴As previsões também foram realizadas usando os parâmetros de cada regressão agregados (*pooled*) com a base imputada (DRECHSLER, 2015; MILES, 2016). Foram obtidos resultados muito similares usando a outra técnica sugerida pelos autores: calcular as probabilidades para cada base imputada (5 probabilidades para cada caso) e, em seguida, agregá-las em uma probabilidade “média” seguindo as regras de Rubin. Em ambos os casos, as previsões foram bastante similares ao modelo ajustado sem as não-respostas, que acabou sendo escolhido para calcular as probabilidades e previsões expostas acima.

frequência de ocorrências em cada célula. As previsões corretas em termos de “alto orgulho” estão no primeiro quadrante e as previsões de “baixo orgulho” no terceiro quadrante.

A tabela 7, abaixo, resume as principais métricas calculadas depois da predição. A análise desses dados mostra que desempenho em termos de acurácia foi superior para a amostra do WVS7. Para a amostra brasileira, o XGBoost identificou menos orgulhosos do que a regressão (- ~100 indivíduos), mas um número maior de não-orgulhosos (+ ~400 indivíduos).

O desempenho do XGBoost em termos de especificidade, ou identificação dos verdadeiros negativos (ou baixo orgulho), foi bem melhor em ambas as amostras, provavelmente devido aos pesos usados para treinar o modelo. A diferença dos valores de κ mostra a superioridade geral do algoritmo de aprendizado de máquina. O F1 não contabiliza os verdadeiros positivos, por isso a proximidade entre os dois modelos usando esta métrica.

Tabela 7: Resultados das predições para o WVS7 e Brasil. As predições foram calculadas usando o valor de corte otimizado pelo método de Youden.

| Métrica | Base Brasil | | Base WVS 7 | |
|----------------|-------------|---------|------------|---------|
| | Regressão | XGBoost | Regressão | XGBoost |
| | Log. | | Log. | |
| Acurácia | 0.52 | 0.62 | 0.62 | 0.71 |
| Sensibilidade | 0.47 | 0.63 | 0.61 | 0.71 |
| Especificidade | 0.69 | 0.60 | 0.65 | 0.66 |
| Kappa | 0.10 | 0.17 | 0.12 | 0.20 |
| F1 | 0.61 | 0.72 | 0.74 | 0.81 |

^a A soma das predições usando GLM e XGBoost não são iguais. nA predição do GLM usa apenas os casos completos

^b O total de predições para os algoritmos XGBoost e GLM, para a Base Brasileira e do WVS7 são respectivamente 7593, 6457, 98278 e 40082 observações.

É importante notar que o XGBoost elaborou 50.000 predições a mais do que a regressão para a amostra do WVS7 (8% do total dos casos), e 567 a mais para a amostra do Brasil (45% do total dos casos). Isso porque consegue fazer predições em amostras contendo não-respostas, o que representa uma grande vantagem com relação à regressão⁵.

⁵Os resultados obtidos com as amostras imputadas foram praticamente os mesmos e foram omitidos da tabela.

4.4. INFERÊNCIA

O XGBoost teve resultados melhores do que a regressão logística. O modelo também foi relativamente melhor em prever novas e “futuras” declarações de orgulho, tanto nas amostras do WVS7 quanto do Brasil. Na medida em que estes dados foram representados da mesma forma com a qual foi realizado o aprendizado, pode-se dizer que o modelo foi capaz de generalizar o conhecimento ao qual foi exposto, e prever a ocorrência de patriotas em novas observações.

Entretanto, apesar de fazer um bom número de previsões verdadeiras, o modelo não deixa claro como nem por que as faz. Esse é o custo de usar um conjunto de árvores de decisão, florestas, gradientes descendentes e regressões lineares combinados em um único modelo. No jargão do aprendizado de máquina, esse é, tipicamente, um algoritmo que produz modelos “caixa-preta”.

A utilidade desses métodos para a ciência social depende não somente da capacidade de fazer previsões, que pode ser útil para construir *surveys* interativas, ou criar aplicações onde a previsão de valores e crenças tenha importância. A utilidade principal viria da possibilidade de fazer inferências explicativas a partir do modelo gerado.

Em outras palavras, o modelo foi capaz de aprender as relações existentes entre os casos e “induzir uma teoria” que possa reproduzir o padrão encontrado? É possível responder questões do tipo: quais as opiniões sobre a política têm maior relação com o orgulho? Pessoas religiosas têm maior propensão a crer no estado?

Essas perguntas se situam no paradigma da inferência estatística (JAMES et al., 2013). Ao contrário do modelo de regressão, para o qual a inferência é explícita e interpretável facilmente (erro, intervalo de confiança, resíduos), os modelos de aprendizado de máquina apresentam dificuldades a este respeito.

Abaixo, foram testadas duas técnicas desenvolvidas para extrair explicações do modelo gerado. A importância global e local das variáveis.

4.4.1. Importância global

O XGBoost é um modelo que não é explícito quanto aos seus resultados. Assim como todos os métodos que combinam árvores em florestas, estes modelos são de difícil interpretação (BIAU; SCORNET, 2016).

Afim de identificar como as variáveis foram usadas e contribuíram para o desempenho global do modelo, as informações sobre a importância global das variáveis foram extraídas de duas formas. A primeira delas é pelo “ganho”, ou *Mean Decrease Impurity (MDI)*, que indica o quanto uma determinada variável foi útil na construção de cada árvore de decisão do modelo.

A importância global das variáveis pelo ganho pode ser vista na figura 12. Um valor alto neste índice indica uma maior importância para a tarefa de classificação. As linhas horizontais pontilhadas indicam o ganho acumulado de todas as variáveis até chegar à 50%, 70% e 90%. Pode-se constatar que o modelo sustenta mais de 50% da sua capacidade de predição em apenas 22 variáveis, das 92 variáveis no total (após transformações). O modelo parece necessitar de mais variáveis à medida que tenta melhorar a acurácia.

Quatro variáveis aparecem com maior importância: a “Confiança nas Forças Armadas”, a “Importância de Deus”, a “Importância da Religião na vida” e o “Ano de nascimento”. Em seguida, com entre 2 e 2,5% do ganho total, aparecem as variáveis “Confiança na Polícia”, “Satisfação com a vida”, “Idade”, “Sentimento de Felicidade”, “Mudanças Futuras: Maior respeito pela Autoridade” (como uma coisa boa), “Confiança no Parlamento”, “O quanto uma pessoa acredita ter liberdade de escolha e controle da própria vida”, “Posicionamento político” e “Maior nível educacional atingido”.

Em seguida, em um grupo de variáveis com importância média, foram identificadas, entre outras, mais duas questões ligadas à religião, como a “Confiança na Igreja” e a “Frequência de participação em cerimônias religiosas”. Também apareceram em bloco as questões de normas e costumes ligadas ao aborto, divórcio, não pagar transporte público, sonegação fiscal, suicídio, homossexualidade e acesso a benefícios.

Também neste bloco, apareceram algumas variáveis de cunho econômico como a crença na propriedade privada e na competição, e a atuação do governo com relação à desigualdade. Nas questões institucionais, apareceram a confiança nas empresas, no sistema de justiça e

nos serviços públicos, para além da confiança na Organização das Nações Unidas (ONU). Apareceram também neste grupo a preferência por líderes fortes ou técnicos na gestão do sistema político, mas também com importância moderada. Enfim, três questões pessoais como a renda, a saúde e o número de filhos tiveram importância equivalente. As demais 56 variáveis contribuíram, cada uma, menos de 1% no modelo.

De modo a verificar eventuais vieses nesta classificação de importâncias, foi usada a técnica da permutação, ou *Mean Decrease Accuracy (MDA)*. Os resultados da importância pela permutação são apresentados na figura 13.

Algumas consistências foram encontradas comparando-se os dois métodos. Dentre as variáveis que figuram entre as primeiras 20 em ambos os rankings estão novamente a confiança nas Forças Armadas, na Polícia e na Igreja; a importância da religião na vida e a frequência de cerimônias religiosas; e a opinião de que um maior respeito pela autoridade seria uma mudança positiva no futuro. O ano de nascimento apareceu, novamente, como uma variável importante. Variáveis relacionadas ao sistema político, como a confiança no parlamento e no sistema de justiça, também foram identificadas usando ambas as técnicas. Por fim, em ambos os dois rankings, aparece o sentimento de felicidade como uma variável importante para o desempenho global do modelo.

Também vale notar a posição das variáveis que foram usadas na regressão logística relacionadas ao materialismo e pós-materialismo em ambos os rankings. A variável mais bem colocada em ambos foram a manutenção da ordem e o progresso em direção à uma sociedade mais humana (valor pós-material). Em ambos os casos figuram somente a partir da 30ª colocação no ranking das mais importantes.

Algumas variáveis se tornaram menos importantes e outras subiram sua posição no ranking, o que pode ser visualizado na figura 14. Uma explicação possível pode ser a representação dessas variáveis. Para as variáveis do pós-materialismo, as variáveis foram transformadas em variáveis binárias antes do treinamento e podem ter tido menos peso decorrente do viés conhecido do XGBoost para variáveis com maior cardinalidade. Entretanto, outras variáveis que foram transformadas em binárias *one-hot* foram usadas mais frequentemente pelo modelo, o que pode significar, simplesmente, que não são tão importantes para o desempenho global do modelo.

Figura 12: Importância das variáveis no modelo XgBoost. A figura apresenta a importância de cada variável calculada usando os índices Ganho. As medida é relativas, onde a soma total é igual a 1.

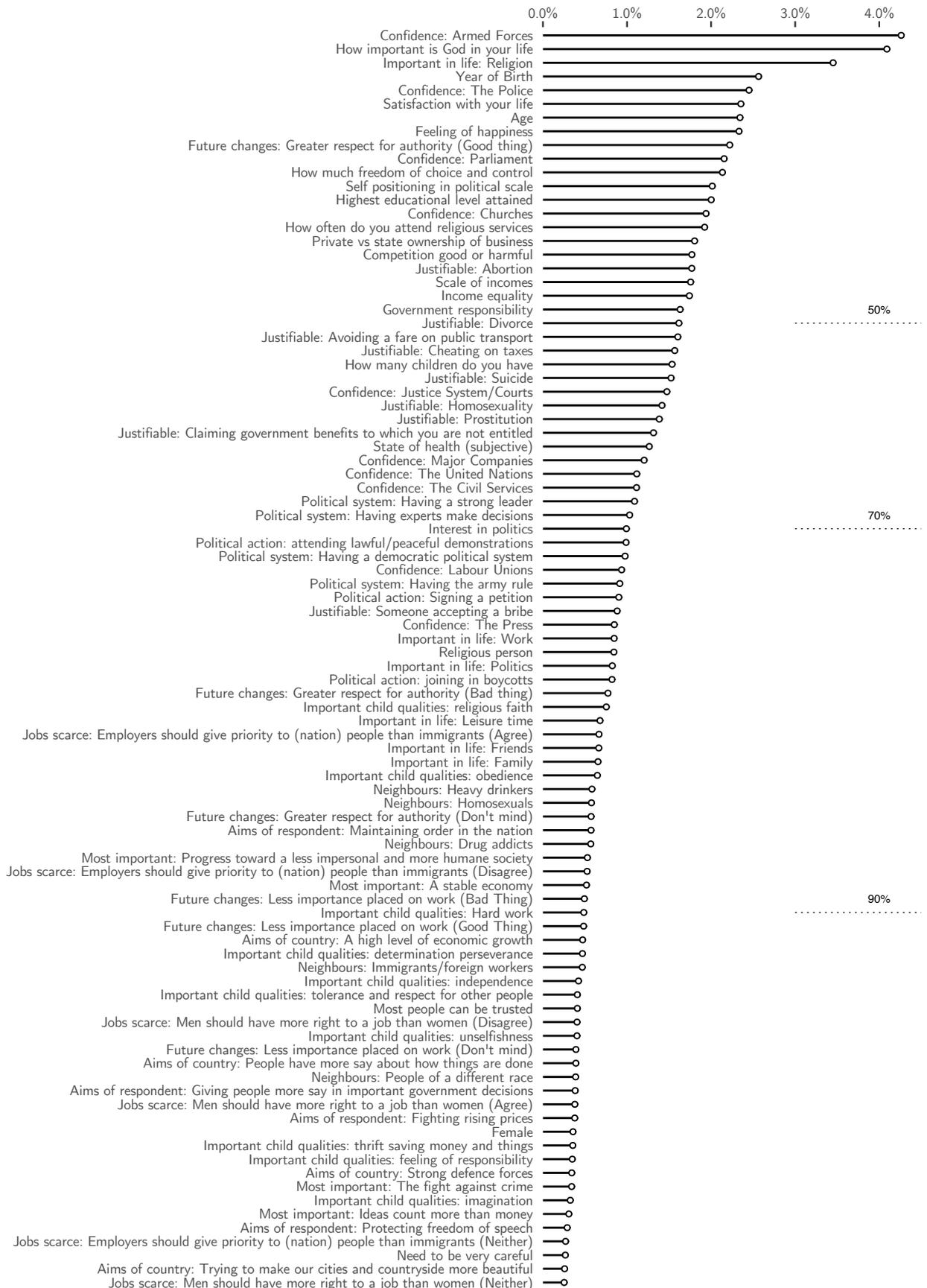


Figura 13: Importância das variáveis usando o método da permutação. A importância da variável corresponde ao aumento da penalidade quando a variável é retirada do modelo.

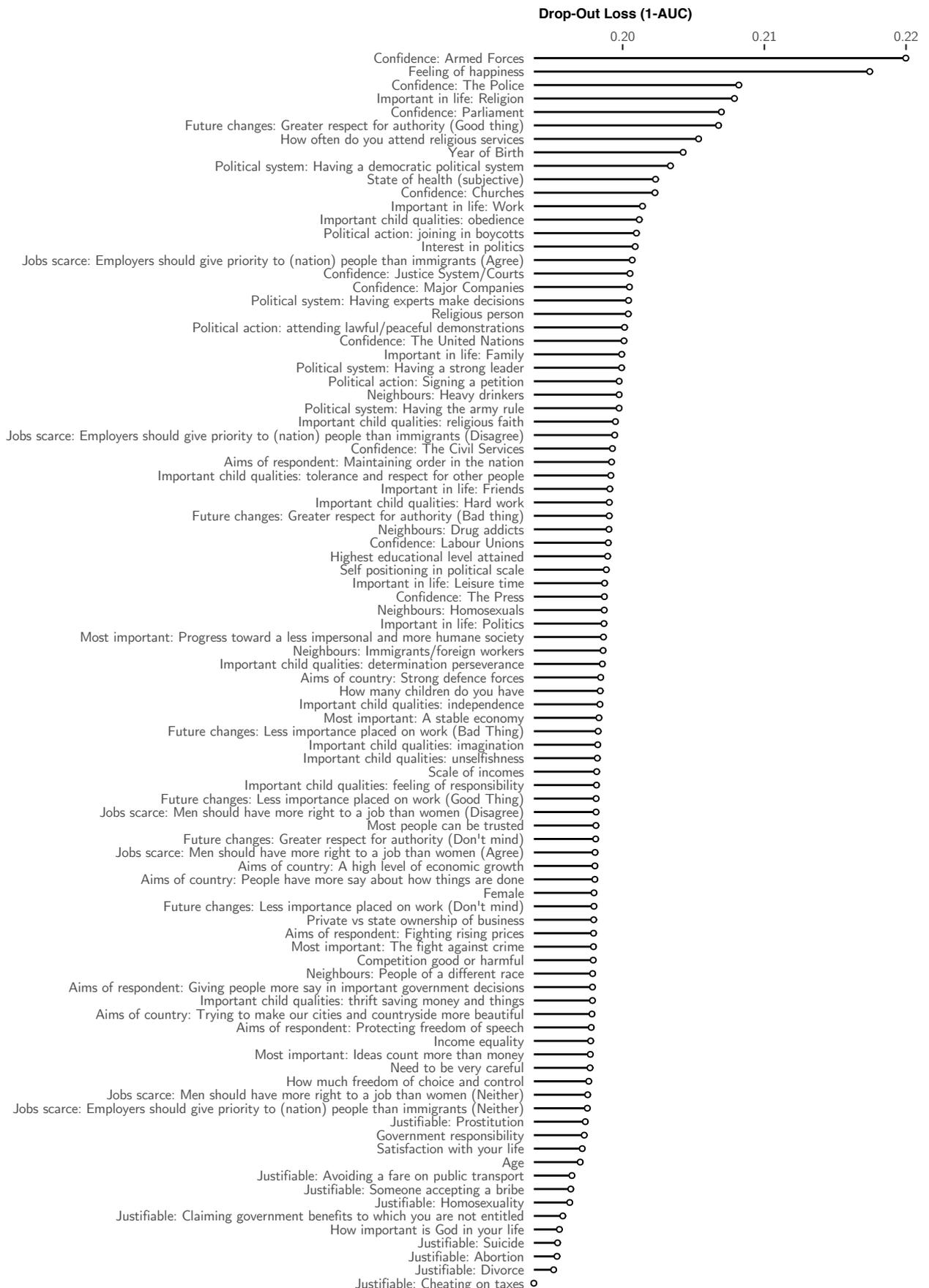
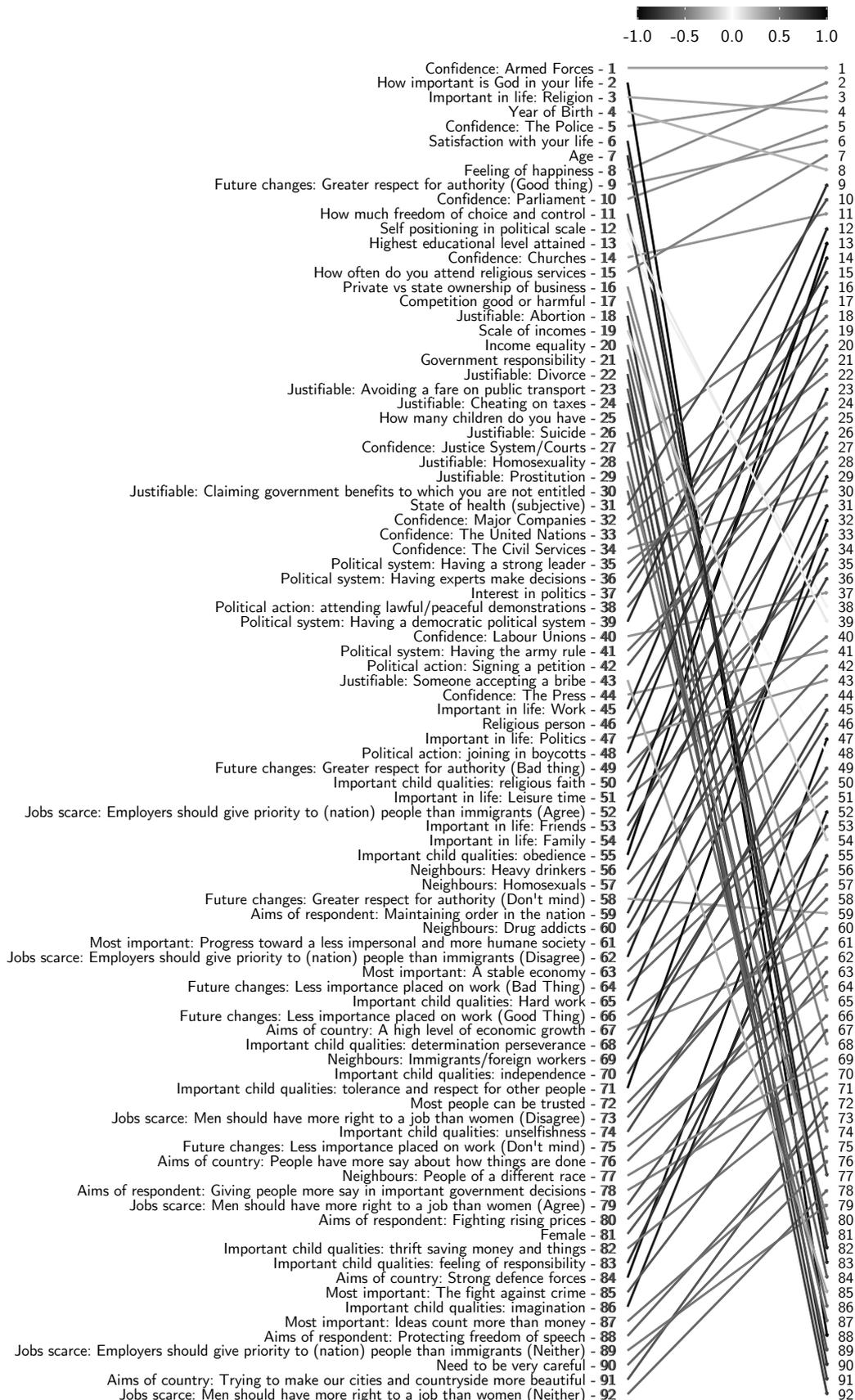


Figura 14: Diferença entre os métodos de importancia global para todas as variáveis. Na esquerda o ranking gerado pelo método interno do XgBoost, à direita o ranking gerado pelo método da permutação.



A importância pela permutação é uma medida global da importância das variáveis que pode ser agrupada (GREGORUTTI; MICHEL; SAINT-PIERRE, 2013). Assim, é possível ter uma ideia da redução da AUC quando são retirados grupos de variáveis importantes. Os resultados agregando algumas das variáveis mais importantes se encontram na figura 15.

Pela análise da figura, constata-se que ao retirar as variáveis relacionadas à autoridade e à religião existe um impacto importante no modelo, que perde em desempenho. As prioridades materiais ou pós-materiais se situam ainda próximas dessas variáveis, com maior importância global, por exemplo, do que algumas variáveis que medem a relação com o sistema político.

Figura 15: Impacto das variáveis agrupadas no desempenho do modelo usando o método da permutação. Foram agrupadas as variáveis mais importantes usando as técnicas MDI e MDA. A redução do ajuste é medida em 1-AUC.



4.4.2. Importância local

Uma técnica de explicação local foi aplicada nas predições realizadas na base do WVS7 e Brasil. Esta técnica, feita especificamente para modelos “caixa-preta”, foi desenvolvida para obter mais confiança no modelo. Segundo os autores, a técnica permite confirmar a importância das variáveis que foram identificadas no modelo treinado, verificando o quão importante elas são nas decisões individuais que o modelo faz no momento da predição. Esta técnica permite também fazer inferências mais precisas e elaborar hipóteses para um modelo mais parcimonioso selecionando apenas as variáveis mais importantes para fazer a predição.

A ideia é criar modelos de substituição aos modelos complexos para entender como são feitas as predições. No caso da técnica usada no trabalho, a premissa é que observações

próximas no espaço dos dados são similares e a predição pode ser explicada localmente com modelos mais simples.

Os modelos explicativos locais foram ajustados usando uma regressão tipo *ridge* para selecionar as 10 variáveis de maior peso local. Na base brasileira, o R^2 dos modelos locais variou entre 0,16 e 0,48, com uma média de 0,31. Estes valores da qualidade do ajuste indicam que os resultados não são muito confiáveis, pois em média apenas 30% da variância pode ser explicada com as 30 variáveis de maior peso no modelo. Na base com outros países do WVS, os resultados foram similares. Os modelos locais explicam entre 10% e 45% da variância.

Mesmo com resultados limitados, a possibilidade de inferência foi ampliada. Foram selecionados os casos onde as predições tinham alta probabilidade de serem corretas para Alto Orgulho. Essas predições foram agregadas por país (ou onda no caso do Brasil), de forma a obter indicações sobre como o modelo treinado usa as variáveis para membros de uma mesma nacionalidade ou entrevistados em um mesmo momento histórico no caso das ondas no Brasil. Nas figuras abaixo, as variáveis mais importantes localmente aparecem sustentando, ou contradizendo, a predição de “alto orgulho” feita pelo modelo gerado pelo XGBoost.

As explicações geradas pelos modelo de substituição locais confirmam a importância de algumas variáveis. Mas a inferência com essa técnica vai um pouco além delas. Na figura 16, são apresentados os resultados para os Estados Unidos, o Reino Unido, a Itália, a Alemanha, o México e a Suécia. Estes países foram escolhidos para compor a figura por serem exemplos recorrentes no estudo da cultura política. A Suécia é um dos países que atingiram os mais altos níveis de pós-materialismo.

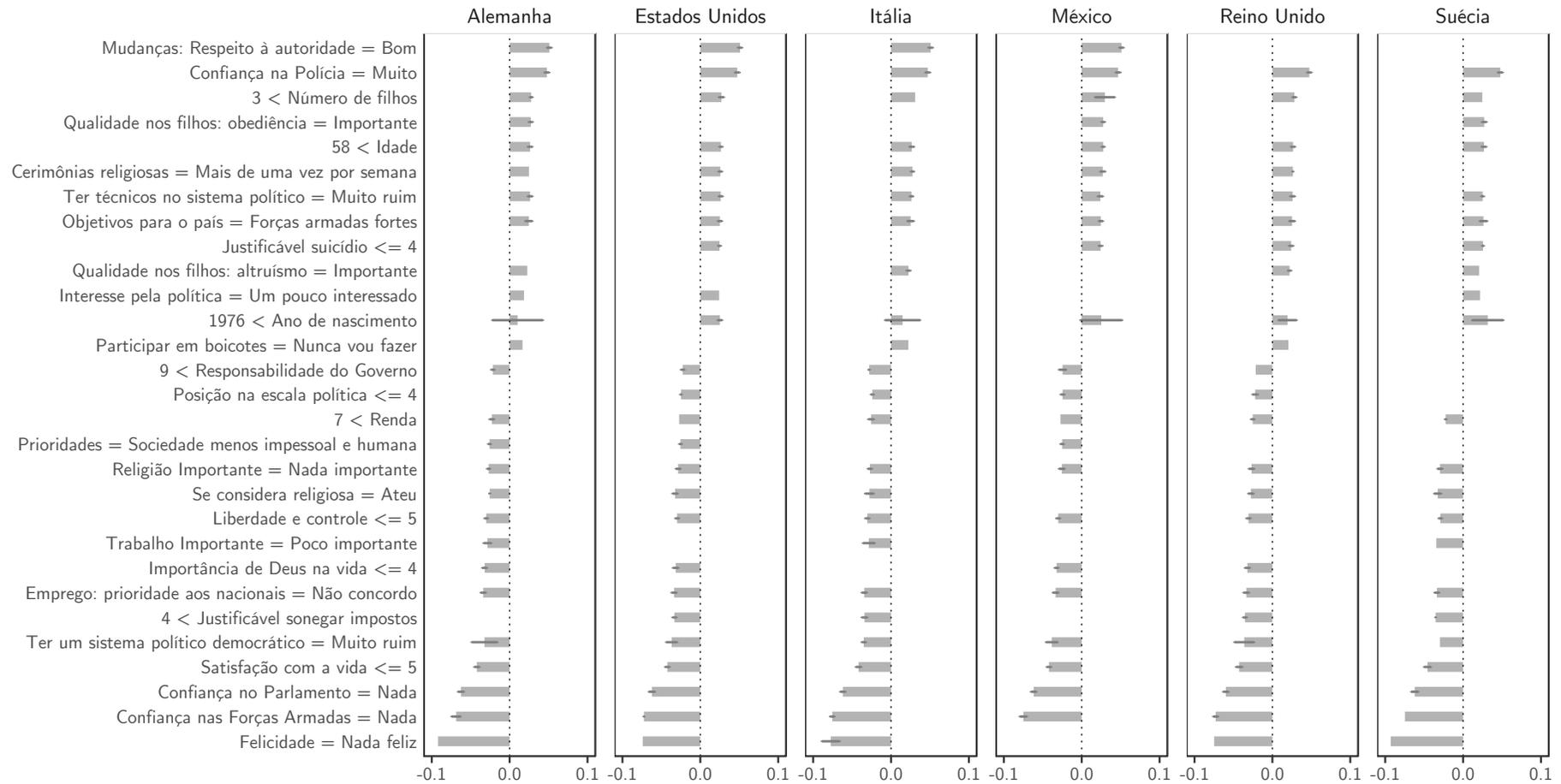
As explicações locais para estes países confirmam a importância de algumas variáveis já identificadas globalmente. A confiança na polícia, ter mais respeito a autoridade no futuro e forças armadas mais fortes são algumas delas. A maior frequência de cerimônias religiosas e a condenação de atos como o suicídio também aparecem corroborando as predições do modelo. Apareceram também variáveis relacionadas a valores tradicionais como as preocupações com a qualidade nos filhos de serem obedientes e altruístas. O próprio número de filhos, que pode ser indicador de maior orgulho, também apareceu. Também foram identificadas variáveis que indicam um desinteresse pela política, incluindo as novas formas de engajamento (boicotes, por exemplo). Enfim, a idade parece ser um indicador de maior orgulho nesses países.

Contradizendo as predições de alto orgulho aparecem o sentimento de infelicidade e satisfação com a vida, e a falta de confiança em instituições como as Forças Armadas e o Parlamento. As declarações de que a responsabilidade do governo deve ser maior reforçam a descrença no sistema político democrático e parecem estar ligadas a um posicionamento mais à esquerda no espectro político. Algumas variáveis relacionadas à secularidade também aparecem contradizendo as predições tais como a pouca importância dada à religião, o ateísmo e a pouca importância de deus na vida. As auto-declarações sobre um desejo de maior liberdade e controle da vida também aparecem contradizendo predições de maior orgulho. A pouca importância dada ao trabalho não aparece nos Estados Unidos e no México, mas aparecem na Suécia, na Itália e na Alemanha. Em todos os países, a tolerância a priorização de empregos a não-nacionais aparece contradizendo a predição de alto orgulho. Enfim, a renda aparece como um fator que contradiz as predições de alto orgulho.

Na figura 17 são apresentados os resultados para as cinco ondas da amostra brasileira do WVS. Assim como na figura precedente, as variáveis mais importantes localmente podem ser vistas com sua contribuição sustentando ou contradizendo a predição de “Alto Orgulho”. De forma geral, foram encontrados padrões semelhantes aos dos demais países. As variáveis relacionadas à confiança na política e o desejo de ter forças armadas fortes aparece relevante. A alta confiança nas Forças Armadas, entretanto, não aparece entre as variáveis com mais peso para as predições. Entretanto, a confiança nula (pouca confiança) aparece contradizendo as predições, confirmando a importância desta variável, mas de outra forma. A religião reaparece juntamente com a variável relacionada a maior frequência a cerimônias religiosas.

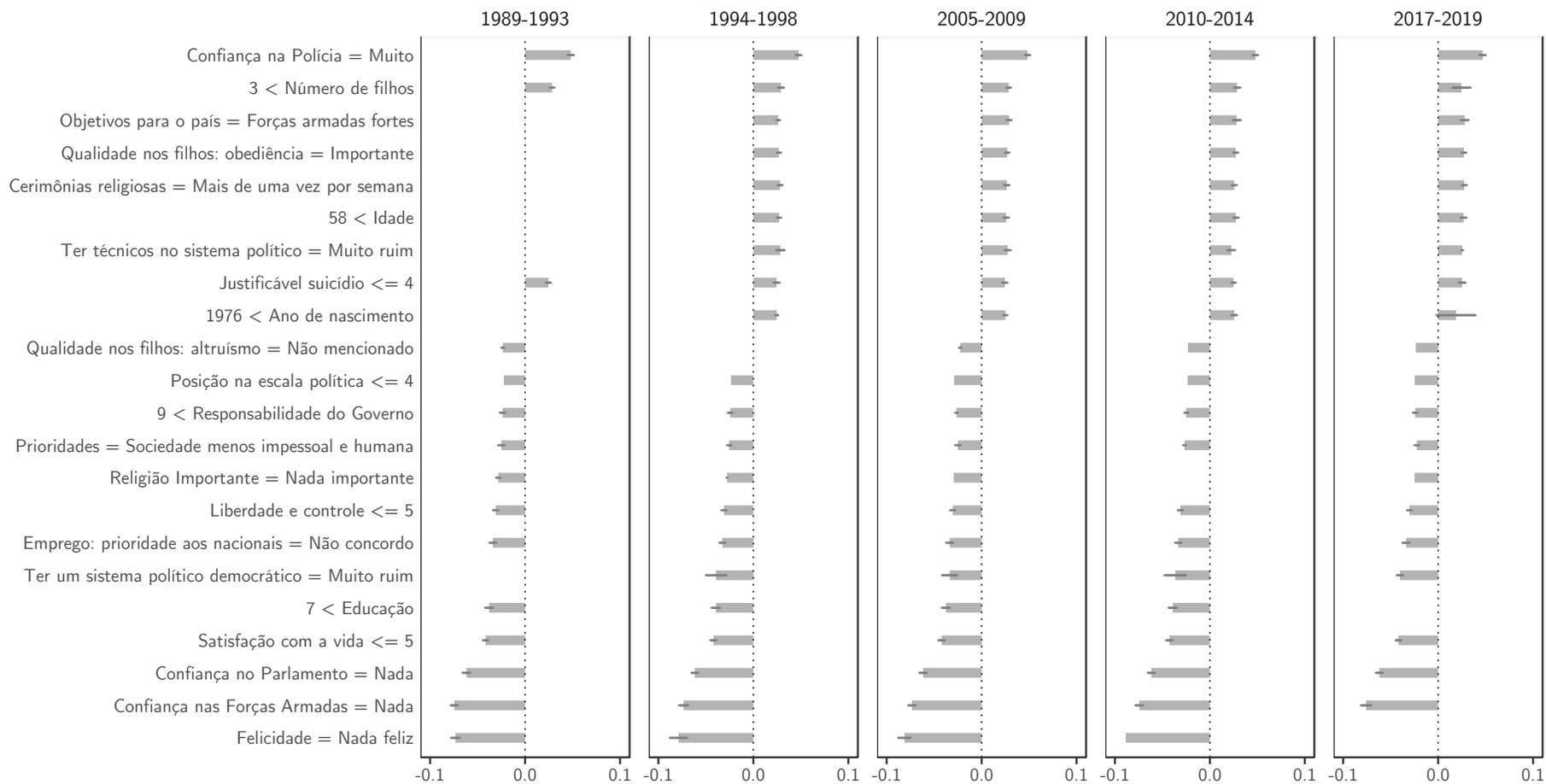
As variáveis usadas no modelo local que contradizem o orgulho se assemelham aos demais países: baixo sentimento de felicidade e satisfação com a vida, desconfiança no Parlamento e nas Forças Armadas. O desejo de maior responsabilidade governamental vem em conjunto com posições mais à esquerda. Contradizendo as predições de maior orgulho estão também a maior valorização da liberdade e controle, o desejo de uma sociedade mais humana e um maior nível educacional.

Figura 16: Importância local das variáveis para 6 países. A figura apresenta a importância local média de cada variável para a predição de Alto Orgulho usando as amostras do WVS7 dos Estados Unidos, Reino Unido, Itália, Alemanha e México e Suécia.



A importância local média foi computada para as predições com maior probabilidade em cada país e usando os modelos locais com melhor ajuste. O valor médio é a média dos coeficientes da regressão ridge para cada variável. Valores menores de zero indicam que a variável contradiz a predição, valores positivos indicam que a variável corrobora com a predição feita pelo modelo caixa-preta

Figura 17: Importância local das variáveis para o Brasil nas 5 ondas realizadas pelo WVS. A figura apresenta a importância local média de cada variável para a predição de Alto Orgulho usando as amostras do Brasil para as ondas 2, 3, 5, 6 e 7.



A importância local média foi computada para as predições com maior probabilidade e usando os modelos locais com melhor ajuste em cada onda do WVS para o Brasil. O valor médio é a média dos coeficientes da regressão ridge para cada variável. Valores menores de zero indicam que a variável contradiz a predição, valores positivos indicam que a variável corrobora com a predição feita pelo modelo caixa-preta

5. DISCUSSÃO

What is learned by a machine learning method is a kind of “theory” of the domain from which the examples are drawn, a theory that is predictive in that it is capable of generating new facts about the domain — in other words, the class of unseen instances. Witten, 2011

Descobrimos que toda inferência, conseqüentemente, toda prova, e toda descoberta de verdades não-evidentes em si mesmas, consistem em induções e na interpretação de induções; que todo o nosso conhecimento não-intuitivo provém, exclusivamente, dessa fonte. Mill, 2011

Uma vez apresentados os resultados, quais foram os achados que podem ser usados no aprimoramento das teorias do patriotismo? Além disso, é possível afirmar a utilidade da integração do aprendizado de máquina na produção de conhecimento nas ciências sociais?

Este capítulo traz a discussão sobre essas duas questões. A primeira parte do capítulo trata das implicações para as teorias do patriotismo e mais especificamente para a cultura política e a teoria do pós-materialismo. São discutidas as principais inferências que puderam ser feitas a partir do exercício indutivo e se estas podem complementar, e como, as hipóteses existentes, dando início a um novo ciclo dedutivo com base em um olhar renovado do fenômeno e novas expectativas em relação aos dados.

Na segunda parte do capítulo, discutem-se as implicações gerais dos resultados obtidos para a metodologia de pesquisa em ciência social. Em particular, são discutidos os pontos positivos e negativos da integração do aprendizado de máquina e os principais pontos de incerteza que foram identificados.

5.1. IMPLICAÇÃO PARA AS TEORIAS DO PATRIOTISMO

Diversos questionamentos surgem a partir das inferências do aprendizado de máquina. As inferências são coerentes com as teorias existentes do patriotismo? Os resultados trouxeram alguma novidade para o entendimento do fenômeno? Quais as implicações das inferências que foram extraídas do modelo de aprendizado de máquina para a cultura política e para teoria do pós-materialismo?

Foram feitas inferências em dois níveis, global e local. Uma primeira discussão pode ser feita com relação às variáveis que tiveram mais importância global no modelo preditivo. As variáveis que apareceram com maior importância indicam uma relação clara do orgulho da nacionalidade com a confiança que um indivíduo deposita em instituições como as forças armadas e valores e crenças religiosas. No primeiro caso, esses resultados não são surpreendentes, na medida em que as forças armadas fazem parte do conjunto de instituições ditas democráticas, fazem parte do estado-nação, e que são responsáveis pela proteção da nação contra inimigos externos. Em muitos casos, inclusive, trata-se de uma variável que é utilizada como um indicador de nacionalismo.

A presença da religião também não surpreende devido a sua importância na formação de valores básicos (DIAMOND, 1994; INGLEHART; WELZEL, 2005). O modelo identificou não somente a importância da confiança na instituição Igreja, mas também da importância que Deus tem na vida das pessoas e da própria frequência a cerimônias religiosas. Isso indica que os níveis de crença e prática religiosa têm relação com o orgulho nacional. Este parece vir de encontro à tese que apresenta o patriotismo moderno como uma crença no estado-nação que substituiu valores religiosos no tempo das revoluções (Hobsbawm, Tilly, Anderson, Smith). Em vez de substituição de valores religiosos por seculares, nota-se uma evolução em paralelo.

Nesse sentido, o patriotismo parece não ter deixado de ser um sentimento de paixão direcionado aos monarcas para se tornar um sentimento racional com relação ao estado-nação, como sugeriu Tocqueville. Crenças religiosas parecem conviver com uma crença secular no estado e na nação, sem que uma tivesse substituído a outra, como sugerem os estudos anteriores na corrente modernista da identidade nacional. A crença na providência que antes era direcionada aos monarcas, a Deus e à Igreja incorporados no estado, parece agora estar compartilhada com a crença em outra providência, não mais divina, mas real — o estado.

Pelas associações encontradas no modelo, as variáveis relacionadas ao *feedback* racional dos cidadãos ao estado são limitadas. Apesar de aparecerem entre as primeiras variáveis em termos de importância global, a confiança no Parlamento e o posicionamento político têm menos importância que as demais. Também são limitadas as indicações de que prioridades futuras influenciem orgulho nacional, que parece invariável, independentemente da direção política tomada pelos governantes. Pela análise local, abaixo, esses fatores parecem ser mais importantes para a declaração de baixo orgulho.

A relação do orgulho com o desejo de mudanças futuras em direção a um maior respeito à autoridade é a única variável que tem relação com projeções futuras e abstrações sobre a sociedade ideal do futuro. Aliada às demais variáveis, essa idealização reforça a relação entre orgulho nacional e um apreço a estruturas hierárquicas e de autoridade. Parece existir um entendimento de que uma relação de subordinação é uma condição necessária para a manutenção da ordem social. Tanto a importância de variáveis ligadas à Igreja quanto às Forças Armadas também indicam a relação do orgulho nacional com o apreço a estruturas hierárquicas. Esse é um fator já identificado em países como o Brasil (CASTRO, 2011; MOISÉS; MENEGUELLO, 2013), mas surpreende pela relação encontrada também em países onde valores dessa natureza estão em declínio tanto por razões de afluência, quanto pelos efeitos da globalização. A importância da confiança na polícia também reforça a relação do orgulho nacional com crenças de natureza hierárquica e de autoridade.

Talvez o efeito do declínio de valores dessa natureza, que podem ser considerados como tradicionais nas teorias da cultura política, seja captado pela importância global das variáveis idade e coorte. Ambas figuram como variáveis importantes na predição, podendo indicar o efeito previsto pelas teorias do pós-materialismo sobre determinadas coortes e também efeitos da globalização, no caso da idade. Na inferência realizada na importância global das variáveis, não é possível identificar a direção do efeito, mas a teoria parece indicar a hipótese mais plausível para o efeito da idade. Pessoas com idade maior de 58 anos parecem mais orgulhosas. A relação encontrada com os nascidos depois de 1976 é, entretanto, contraintuitiva. Os achados empíricos anteriores indicam que o efeito de coorte é mais forte nas gerações mais jovens.

Os níveis de satisfação e felicidade parecem estar envolvidos também na declaração de orgulho e patriotismo. A presença dessa variável entre as variáveis com importância

intermediária indica que essa relação é importante para predizer o orgulho de alguns indivíduos. Segundo as hipóteses da psicologia social e da cultura política vistas anteriormente, essas variáveis podem indicar tanto conformidade quanto autonomia. Por um lado, tanto uma quanto a outra declaração são características de sociedades pós-materiais, onde a satisfação é reforçada pela sensação de agência e controle sobre o curso da vida. A tese defendida por Welzel e Inglehart (2010a) é que o crescimento de valores emancipativos faz com que cidadãos se sintam agentes individuais da política e menos afeitos a participar de guerras em nome do país. Por outro lado, esses sentimentos podem estar ligados mais à anomia ou a uma certa resignação com a vida que teria como consequência uma identificação nacional mais exacerbada (FABRYKANT, 2014). Outra possibilidade, que tem sustentação nas inferências locais (ver abaixo), é de estes sentimentos estarem ligadas à insatisfação com o desempenho do estado. Nesse caso, corroborando a tese da existência de um patriotismo “racional” que se orgulha do desempenho do seu governo.

Outra discussão pode ser feita tendo como indicação as inferências extraídas dos modelos locais de substituição. Um primeiro ponto importante é que a análise local aumenta a capacidade indutiva do uso do aprendizado de máquina, pois permite identificar a direção dos efeitos, que no caso das variáveis globais não é indicado.

Outro ponto importante é que esta técnica de inferência no aprendizado automático tem relação mais direta com o campo da cultura política. Na cultura política, as análises são geralmente em nível agregado de culturas ou sociedades, enquanto o nível de análise do aprendizado de máquina foi, neste caso, individual. Como opiniões, atitudes, crenças e valores, coletados no paradigma das pesquisas de opinião, onde a legitimidade se dá pela representatividade da amostra sobre a população e pela padronização das questões, podem ser modeladas individualmente, mesmo sem “nacionalizar” as amostras, e, em seguida, retornar para uso em nível societal? Em outras palavras, como fazer com que uma teoria aprendida em nível individual suba ao nível da sociedade, ou das síndromes de valores sociais? As técnicas de agregação das inferências com modelos locais são uma pista para dar uma resposta a esse desafio, pois permitem a análise local de indivíduos de uma mesma nacionalidade ou onda de um *survey*.

Antes de discutir os resultados, é preciso dizer que se trata de uma abordagem experimental, visto que essa técnica foi desenvolvida para entender porque um modelo

caixa-preta consegue fazer previsões para casos individuais. A agregação não está prevista na técnica, pois os coeficientes gerados são para modelos preditivos individuais (no caso aqui foi calculada a média dos coeficientes para cada nacionalidade). Além disso, os ajustes encontrados com os modelos locais deixaram a desejar e a magnitude dos efeitos é pequena e de difícil interpretação. Enfim, as variáveis utilizadas no modelo local também dependem provavelmente da maneira como foram transformadas as variáveis explicativas (e.g. dummy, ordenada, numérica etc.). Por fim, as variáveis mais importantes no modelo local foram praticamente as mesmas para cada país, revelando uma dificuldade desses modelos de detectar diferenças entre os casos.

Dito isso, a agregação parece ter revelado relações coerentes tanto com a importância global revelada pelo modelo caixa-preta, quanto com as teorias existentes do patriotismo. De forma geral, as variáveis identificadas na explicação local são congruentes com as locais. As variáveis relacionadas à autoridade e à hierarquia foram praticamente as mesmas da análise global. Entretanto, foram encontrados os mesmos padrões para o orgulho nacional ligado a valores tradicionais, tanto em países mais pós-materiais, como a Alemanha e a Suécia, quanto para países como o México e o Brasil.

Nos dizeres de Anderson (2011), a condição nacional e o nacionalismo ora despertam um profundo sentimento de apego manifestado como orgulho, ora vergonha, implícita nas declarações de baixo orgulho. Para além dos fatores comuns do patriotismo visto acima, ligado a valores tradicionais, aparecem outros que permitem deduzir algumas nuances.

Uma das associações encontradas com o baixo orgulho parece indicar a existência de um patriotismo que se manifestaria pela insatisfação com o governo no seu papel de provedor. As variáveis confiança no parlamento (baixa) e a posição política (mais à esquerda), que contradizem as previsões de maior orgulho, permitem deduzir um “patriotismo patriarcal”. Das questões econômicas que apareceram tendo importância global, apenas a responsabilidade do governo reapareceu. A direção do efeito sugerida pelo modelo local é que a falta de responsabilização do governo leva a menos orgulho.

Outra associação encontrada pelos modelos locais são as duas variáveis sobre a liberdade individual. Essa entretanto, parece uma associação de outra natureza e um patriotismo mais liberal. O que parece ter sido identificado acima é um patriotismo que se

orgulha de um governo mais assertivo e provedor. Um outro “patriotismo liberal” parece esperar um governo menos intrusivo e com menos responsabilidade. A crença na necessidade de liberdade de escolha e controle também é relacionada a uma declaração de baixo orgulho, o que corrobora essa hipótese. Talvez algumas pessoas se sintam menos orgulhosas justamente porque o governo toma responsabilidades demais. Esses fatores indicam que podem existir um patriotismo ligado a valores liberais e antiestado que respondem positivamente em termos de orgulho no caso de as suas liberdades serem asseguradas e menor intervenção estatal. De certa forma, um orgulho de uma pátria pela falta de Estado.

Essas duas variantes podem talvez ser interpretadas como uma forma de patriotismo cívico, ou mais “racional” e utilitarista, que reage em função da relação estabelecida entre os governados e os que governam a pátria.

Questões ligadas às normas e aos costumes, ligadas à tolerância, como o aborto, divórcio, não pagar transporte público, sonegação fiscal, suicídio, homossexualidade e acesso a benefícios, tiveram influência limitada nos modelos locais. Em relação a essas questões, apenas a aceitação ou não do suicídio apareceu tanto na análise da importância global quanto local. Nesses casos, quanto menos justificável, mais orgulho, talvez como efeito da relação com valores religiosos e pró-vida.

Enfim, fatores como coorte, idade, renda e educação reapareceram, confirmando a importância dos efeitos de idade e coorte em ambos os modelos.

Os resultados apresentados pelo aprendizado de máquina levam, inevitavelmente, a uma discussão das hipóteses do pós-materialismo testada por meio da regressão. As variáveis com mais importância global parecem sustentar a tese do pós-materialismo na medida em que o orgulho nacional se relaciona com os demais valores tradicionais, como a família, a religião, etc. Parece haver uma disposição generalizada em aceitar o orgulho nacional “normativo” imposto pelo estado de forma menos crítica (FABRYKANT; MAGUN, 2015).

Entretanto, as cinco variáveis que foram testadas na regressão não foram encontradas usando-se as técnicas de inteligência artificial. A relação com o orgulho existe, mas com pouca importância na predição de novos casos, feita pelo modelo que obteve o melhor desempenho. A razão pode estar ligada à codificação dessas variáveis, mas talvez as prioridades individuais, apesar de medir corretamente valores pós-materiais ou emancipativos, não ajudam tanto

quanto outras variáveis mais explícitas para explicar o patriotismo. Em última análise, a modernização proposta por Inglehart e outros autores pode não ter tido o efeito esperado nesses aspectos do sistema de valores.

Quanto ao Brasil, as declarações relacionadas às forças de ordem, uma visão positiva do autoritarismo e da religiosidade foram as variáveis mais influentes na predição. No sentido contrário, o sentimento de pouca liberdade e controle da vida, a secularidade e uma menor satisfação e felicidade com a vida são indicativos de baixo orgulho pela pátria. Quanto aos militares e à polícia, pode-se dizer que a relação era esperada. Desde o fim do Império, os militares foram apresentados como o reduto do patriotismo brasileiro (OLIVEIRA, 1990). Eles protagonizaram a mais recente valorização patriótica durante a ditadura militar e também, no momento em que se escreve esta tese, na sua associação com o governo Jair Bolsonaro.

Além disso, desde os conflitos na Cisplatina e no Paraguai durante o Império, foi criada uma identidade nacional relacionada às forças armadas. Isso aconteceu mesmo que tendo um efeito diferente dos exércitos na Europa. Ao invés de criar apreço à nacionalidade pela disputa contra um inimigo externo, levou à apreciação pelo modelo de organização e pelo funcionamento (ordem, disciplina) da instituição (SKIDMORE, 1999; LESSA, 2008). Quanto à polícia, e mesmo à Igreja, os brasileiros tendem a ter apreço por instituições hierárquicas e autoritárias (CASTRO, 2011; MOISÉS; MENEGUELLO, 2013).

À luz das inferências, não parece clara a existência de um patriotismo “racional”. No caso de aceitarmos essa hipótese, assim como nos demais países analisados, parecem conviver no Brasil duas formas desse patriotismo: patriarcal e liberal. O primeiro deles ligado a uma expectativa paternalista com relação ao estado. O segundo faz lembrar as elites nacionais da primeira república. Estes têm outro objetivo, mais liberal do que aqueles, e prezam mais pela distância do estado dos seus *affairs* e pela preservação de liberdades individuais. Para estes, o patriotismo parece existir nos mesmos termos do início da república, em um misto de valores tradicionais e burgueses (CARVALHO, 2017).

Os dados do WVS, infelizmente, não permitem encontrar ligações com as identidades nacionais ligadas à natureza, à cultura nacional ou com os motivos edênicos incorporados nas visões de nação fundadas na qualidade do território brasileiro (CARVALHO, 1998; HOLANDA, S. B. de, 2000). Apesar disso, as ligações profundas desses motivos com a religião no Brasil

faz com que as relações encontradas com variáveis dessa natureza sejam a confirmação da existência, e permanência destes motivos no imaginário brasileiro, e influência no patriotismo.

Um último ponto merece atenção. Como já anunciado acima, a análise entre as ondas ou entre países não revela, necessariamente, diferenças. A agregação dos coeficientes dos ajustes indica, principalmente, similaridades entre as variáveis encontradas, o que pode estar relacionado ao tipo de modelo que foi ajustado. A ausência de algumas variáveis entre uma onda e outra, ou um país e outro, pode ser uma forma de interpretar mudanças ou diferenças. Como uma das motivações do trabalho, foi exposto o declínio do patriotismo brasileiro e uma vontade implícita de explicar o fenômeno como uma “mudança de valores”. Devido às limitações do método aplicado, explicações para essas mudanças não foram possíveis.

Contudo, esses resultados podem indicar uma pista de desenvolvimento de uma ferramenta que permita o escrutínio das inferências do aprendizado de máquina por cientistas sociais. Pode existir uma forma de agregar os coeficientes de uma forma que permita confirmar os resultados inferidos e também detectar mudanças de uma onda a outra ou entre culturas diferentes.

5.2. IMPLICAÇÕES METODOLÓGICAS

“se mais de uma teoria é consistente com as observações, guardemo-las todas”. *Epicúrio, 341—271/272 a.C.*

Quando Witten fala de *teoria*, na citação acima, ele trata do tema visto no paradigma do aprendizado estatístico, ou do “aprender com os dados” segundo Hastie, Tibshirani e Friedman (2017). Nesse caso, teoria é o termo que serve apenas para designar uma série de regras ou árvores de decisão de um modelo preditivo.

Esta *teoria* pode ou não ser revelada, dependendo do interesse do pesquisador em fazer somente previsões ou também inferências. Inferência, aqui também no paradigma do aprendizado estatístico, significa poder usar um segundo algoritmo que ateste a acurácia do primeiro (EFRON; HASTIE, 2016). Assim, quando o interesse é inferir as relações que foram aprendidas com os dados, a revelação da teoria depende da interpretabilidade ou explicabilidade do modelo, ou da possibilidade de aplicar outro algoritmo que permita atestar sua acurácia.

São os (bons ou maus) resultados da predição que motivam e dão sustentação a uma explicação posterior (BREIMAN, 2001b). Apesar de melhores e em maior número do que as predições da regressão, os resultados do modelo selecionado em termos de predição foram limitados. Por questões de tempo e limitações técnicas, a afinação do modelo não pode ser levada a cabo, o que limitou o poder preditivo do modelo. Além disso, o problema do orgulho, com classes bastante desbalanceadas, dificultou o trabalho de encontrar um algoritmo com boa capacidade de discriminar entre indivíduos muito ou pouco orgulhosos.

Mesmo assim, existiu o interesse, pelo menos para fins demonstrativos, de fazer inferências e explicar como o modelo consegue fazer essas predições. Para fins acadêmicos como neste caso, não se busca apenas encontrar uma maneira de prever quem será (estará) orgulhoso da pátria. Ao contrário, espera-se que se possa extrair destes modelos uma explicação dos resultados.

A partir desses resultados, pode-se notar, primeiramente que a falta de parcimônia dos modelos é um dos pontos que dificultam a integração do aprendizado de máquina nas ciências sociais. O modelo mais bem ajustado pelo XGBoost foi selecionado levando-se em consideração o problema apresentado em termos de viés-variância durante o processo de validação cruzada. Mas as *teorias* mais complexas acabam sendo mais eficientes e desejadas na predição. Os modelos selecionados acabam sempre tendo uma complexidade importante e que é aceita por causa dos melhores resultados em termos de predição. O princípio de Occam, quando se busca a predição, tem pouca sustentação empírica (DOMINGOS, 1998, 1999).

Essa característica do aprendizado de máquina faz com que o modelo gerado pelo XGBoost use praticamente todas as variáveis. Mesmo que algumas tenham importância marginal para a acurácia do modelo, todas elas podem ser usadas durante a predição (ver as figuras 12 e 13). Fica clara a razão de algumas críticas de que é difícil encontrar resultados negativos como a ausência total de relação com a variável resposta (OLLION, 2018).

Alguns modelos são menos obscuros, como as regressões, que permitem uma interpretação direta dos resultados a partir da multiplicação dos coeficientes e das análises dos desvios padrão, residuais etc. O modelo gerado pelo XGBoost para prever o orgulho é um modelo que não permite essa interpretação direta. Mesmo usando as técnicas de inferência pela identificação da importância global das variáveis e pela explicação local das predições, o

modelo continua opaco.

Como o cientista social pode fazer dessa complexa “teoria” que foi aprendida com os dados uma boa explicação com o mínimo? É talvez nesse ponto que resida a intersecção entre a lógica indutiva da ciência social e a inferência estatística do aprendizado de máquina. Assim como em outros campos, a simplicidade na explicação é um dos pilares da ciência social. Mas a complexidade gerada pelo modelo significa simplesmente que há mais informação a ser extraída e que necessita ser filtrada usando-se conhecimento social.

A abordagem indutiva à qual Mill se refere na segunda citação acima foi elaborada dentro do paradigma clássico da lógica científica, em que a indução faz parte de um ciclo contínuo que alterna momentos de indução e dedução (ver figura 2). Para Mill, no momento da indução não se trata apenas de descrever os fatos, mas de inferir sobre as observações, e que a principal dificuldade “*não é fazer induções, mas escolhê-las*”.

Assim, parece crucial o conhecimento do fenômeno analisado para discernir das inferências estatísticas o que é útil e inútil, o que é novo e o que pode ajudar a reforçar velhas teorias. No aprendizado de máquina, esse conhecimento se refere, normalmente, como *domain expertise*, mas com ênfase na seleção e na engenharia das variáveis antes de se treinar o modelo. No caso do uso por cientistas sociais, o conhecimento parece também ser necessário para interpretar as inferências. Para isso, são necessárias ferramentas específicas.

CONCLUSÃO

Neste trabalho foi explorada a aplicação de algumas técnicas de aprendizado de máquina nas ciências sociais. Foi escolhida uma pergunta sobre um fenômeno clássico da disciplina — o patriotismo — e, a partir dos dados, foram elaboradas algumas hipóteses sobre por que algumas pessoas são mais patriotas do que as outras.

As conclusões que podem ser tiradas da experiência estão divididas em dois níveis. O primeiro nível diz respeito aos resultados obtidos para o caso do patriotismo. O segundo nível de conclusões diz respeito ao uso das técnicas de aprendizado de máquina e a contribuição em termos metodológicos que puderam ser identificadas.

As conclusões que podem ser tiradas com base em uma pesquisa indutiva são diferentes de uma pesquisa que toma o caminho hipotético-dedutivo. Caso tivesse sido privilegiada uma abordagem hipotética-dedutiva, as conclusões seriam apresentadas em termos explicativos, identificando se determinados fatores testados podem ou não ser descartados como explicações para o patriotismo. No caso de um processo indutivo, as conclusões devem ser tratadas como hipóteses. As conclusões são fruto da observação dos dados e constituídas de uma série de generalizações empíricas. Quando agrupadas, transformam-se em hipóteses novas, ou que são capazes de completar, modificar ou refinar hipóteses existentes. Estas hipóteses, por sua vez, serão passíveis de teste empírico retomando o ciclo da produção do conhecimento.

Quais conclusões ou hipóteses puderam ser melhoradas ou reformuladas a partir do aprendizado de máquina? De que maneira os resultados ajudaram a ir além do que já se conhece sobre o tema, de maneira a formular novas hipóteses? Nos parágrafos que seguem, tentamos responder a essas duas perguntas.¹

¹Existem pelo menos duas maneiras de elaborar uma hipótese. O aprendizado de máquina, por mais

Como visto na revisão da teoria e empiria sobre patriotismo, parecem existir duas grandes hipóteses, ou grupo de hipóteses, que tentam explicar o patriotismo. Uma delas explica o patriotismo como uma continuidade do amor à pátria que mudou de objeto com o surgimento dos estados-nação, mas que se manteve quase inalterado em sua relação entre cidadãos e estado. Manteve-se uma relação quase que passional, incondicional e de lealdade que mimetiza o sentimento de um membro de um clã ou tribo ao seu líder. A outra, moderna, sustenta que o patriotismo não seria mais um sentimento incondicional e passional, mas reflexivo e interessado. Em termos de valores e crenças, essas duas hipóteses podem ser traduzidas em síndromes de valores tradicionais e modernos.

De certa forma, essa mudança desejada se inscreve na discussão paradigmática da época. Normativamente, esperava-se que as paixões, glória e honra que motivaram o amor pela pátria anteriormente fossem substituídos por um conjunto de interesses pessoais, mais “brandos” e cívicos. A busca era por uma forma de cimento que substituísse a força religiosa do passado, mas que ainda guardasse uma componente cívica. Uma força que fosse capaz de agir sobre governados e governantes, criando a previsibilidade e a constância necessárias para a estabilidade social.

Esta busca por um patriotismo moderno que ajudaria a criar ligações entre os membros de uma sociedade pode ser transposta aos tempos atuais. Primeiramente, no início do século XX, com a tentativa de diferenciação entre nacionalismo e patriotismo. Em tempos de globalização mais recentes, em que a diversidade étnica e cultural se tornou uma realidade em diversos países, existe, novamente, uma busca por um patriotismo que responda positivamente à necessidade de unir as pessoas em torno de uma mesma realidade imaginada. As formulações de Viroli e Habermas em busca de patriotismos republicanos ou constitucionais parecem se inscrever nessa mesma busca.

Em termos empíricos, a indução usando-se os dados do WVS apontou para um padrão que se repete tanto no Brasil quanto nos outros países. As inferências foram analisadas por modelos globais e locais que permitem extrair as explicações geradas por esses modelos.²

obscuro que possa ser, é implicitamente uma forma de gerar modelos formais na forma de funções matemáticas (abordagens bottom-up) ou de um conjunto de regras (em abordagens top-down como modelos simbólicos, por exemplo) que seriam capazes de reconstruir os dados originais. A tradição das ciências sociais, da ciência política e da cultura política não é, necessariamente, voltada à elaboração de modelos formais, mas sim de modelos menos restritos e informais. Será seguida esta linha aqui, mas já apresentando algumas pistas de operacionalização das hipóteses levantadas.

²As evidências empíricas que embasam as conclusões apresentadas neste capítulo estão nas figuras 12 e 13,

Primeiro, os dados revelam uma ligação forte do patriotismo com valores e a com a própria prática religiosa. Uma das hipóteses para essa associação seria a inexistência de uma verdadeira transição entre o patriotismo antes dos estados-nação e o atual, estando esse sentimento ainda fundamentalmente ligado a crenças tradicionais de pertencimento a uma mesma cultura, etnia, credo etc. A própria inclinação a crer na providência divina, que motiva a crença religiosa, se transfere ao Estado-nação de forma acrítica e mantém os cidadãos em um estado de anomia e de aceitação do seu destino.

A alta satisfação com a vida e a felicidade não aparece como determinante do patriotismo. Mas a sua presença nas análises pode ser declinada em duas hipóteses explicativas concorrentes ou antagônicas. Primeiro, os que se dizem insatisfeitos e infelizes tendem a se dizer também menos orgulhosos. Essa relação sugere que os mais satisfeitos e felizes, anômicos, são também os mais orgulhosos pela simples aceitação da ordem divina. Uma “religião civil”, nos dizeres de Rousseau, uma profissão de fé no interesse geral constituído no estado.

Outro padrão encontrado foi a relação do patriotismo com opiniões a respeito do valor de instituições responsáveis por manter a ordem interna e externa (*i.e.* Polícia, Forças Armadas). A ligação desta religião civil com valores e declarações acerca da autoridade e dos arranjos hierárquicos não surpreende, tanto pela própria constituição vertical e hierárquica da Igreja enquanto instituição, mas também pela própria importância desses órgãos em proteger o povo de dissidências internas e ataques externos. Parece haver sempre uma preocupação com o outro, seja ele estrangeiro, seja concidadão, e um reconhecimento de que é necessário e legítimo o uso desses poderes que detêm o monopólio da violência para conter eventuais dissidências.

Outro aspecto ligado a valores tradicionais é o aparecimento de valores ligados à família como o número de filhos e a obediência como uma de suas qualidades. Essa relação sugere que tanto as preocupações religiosas, quanto com a manutenção da ordem se declinam em nível da família com a perpetuação desses valores pelo processo de socialização. Nas famílias, permanece a expectativa de um futuro no qual o respeito à autoridade seja ainda maior.

Níveis de renda mais elevados parecem escapar desse padrão em todos os países

que indicam as variáveis mais importantes globalmente, na figura 15, onde essas importâncias foram agrupadas por sua natureza e, por fim, nas figuras 16 e 17, que são resultado das inferências em nível local.

analisados, mas não no Brasil. Pessoas com renda mais elevada tendem a se declarar menos orgulhosas. No Brasil, a renda parece não ser condição suficiente. Para o Brasil, o que aparece na categoria de variáveis socioeconômicas é a educação. Níveis de educação mais elevados parecem ser mais determinantes para a predição de maior orgulho indicando uma desconexão entre educação e renda no Brasil, algo menos pronunciado nos demais países onde a afluência média veio acompanhada de um maior nível de educação.

No Brasil, em especial, pode-se também sugerir a hipótese de que a ligação dos valores religiosos e autoritários com o patriotismo seja o resultado de dinâmicas mais recentes, como a luta anti-comunista, por exemplo. No Brasil como nos Estados Unidos, o comunismo foi percebido como uma ameaça à família e à sociedade como um todo. Foram períodos em que as forças armadas, que no Brasil também abrangem as polícias, assumiram a luta contra estes movimentos que atentavam contra os valores tradicionais da família e da religião.

Por fim, apesar de não terem sido estudadas variáveis explícitas sobre a questão edênica no Brasil, o fato de o édem brasileiro ter vindo regularmente decorado de símbolos religiosos pode ter relação com o motivo edênico proposto por Sérgio Buarque de Holanda e José Murilo de Carvalho.

Parece existir pouca evidência de um patriotismo racional. Apesar de que algumas variáveis políticas tenham aparecido nas análises de variáveis globais, a inferência local indica que essas variáveis são usadas para reduzir a probabilidade da predição de maior orgulho. A categoria da variável de confiança no parlamento, mais útil para as predições, por exemplo, é “Nada”, ou seja, baixíssima confiança no sistema político. A categoria mais importante para a predição no que diz respeito à importância de ter um sistema político democrático é “Muito Ruim” e contradizendo a predição de orgulho. A confiança maior e preferência pelas forças armadas reforça a sugestão de uma certa indiferença, ou até repulsa, pela democracia.

Com base nas variáveis que foram identificadas como tendo uma relação negativa com maior orgulho (contradizendo as predições de alto orgulho) e retirando-se as variáveis já mencionadas acima, pode-se identificar, mesmo que de forma moderada, algum indício do efeito de valores pós-materiais pela presença da renda e da educação. A relação com o desejo de uma sociedade menos impessoal e humana, a falta de liberdade e controle da própria vida, e o desacordo com questões protecionistas (emprego de não nacionais) parece também indicar

uma relação entre o baixo patriotismo e valores de autoexpressão.

Entretanto, a falta de patriotismo é marcada também pela infelicidade, pela insatisfação com a vida e pela descrença no sistema político. Essa descrença parece ter mais relação com fatores ligados a valores tradicionais do que seculares-rationais ou de auto-expressão, representados pela relação com o desejo de mais responsabilização pelo Estado. Pode se tratar da expectativa do Estado-providência da religião civil mencionada acima. No Brasil, onde a falta de responsabilidade governamental também contradiz as predições de alto orgulho, a insatisfação parece ligada ao paternalismo e à espera pelo Estado para reduzir desigualdades em detrimento de valores mais ligados à autonomia cívica.

As implicações para as teorias sobre o patriotismo são pelo menos duas. A primeira delas é que a tese de um patriotismo virtuoso que daria uma sustentação difusa ao sistema político parece um pouco frágil do ponto de vista empírico. Parece haver sim uma aceitação do estado-nação em si, na forma de uma quase religião. Entretanto, nenhuma indicação foi encontrada com relação aos regimes democráticos. Pelo contrário, o patriotismo parece ter um efeito deletério tendo em vista a sua associação com valores autoritários. Nesse sentido, as teorias normativas do patriotismo republicano, constitucional ou cosmopolita encontrariam barreiras importantes no que diz respeito aos valores a serem absorvidos e relacionados ao orgulho nacional e ao patriotismo.

Essas ideias e hipóteses, que puderam ser elaboradas a partir dos resultados do aprendizado, levam ao segundo nível de conclusões, que diz respeito aos resultados em termos metodológicos da experiência com o aprendizado de máquina aplicado às ciências sociais. De forma geral, pode-se concluir que essas ferramentas e técnicas são úteis, mas que ainda existem dificuldades importantes a serem superadas. Primeiramente, foi possível gerar novas ideias de pesquisa e refinar algumas dos achados já consensuais na literatura. O presente estudo por exemplo permitiu contestar, ao menos em parte, a existência de um patriotismo racional, quando olhamos para os dados do WVS como um todo. Dito de outra forma, mesmo que exista um patriotismo dessa natureza, parecem existir outros fatores mais importantes para prever o patriotismo, e, se acreditarmos nessas predições, explicá-lo.

A experiência de uso das técnicas de aprendizado de máquina revelou também as diferenças de abordagem causadas pela cultura preditiva da modelagem estatística. Esse

caminho implica outras formas de usar os dados e avaliar os erros dos modelos. Mais importante, a experiência revelou a dificuldade de fazerem-se inferências sobre modelos “caixa-preta”. Apesar de já existirem técnicas para isso, persiste a dificuldade em se extraírem explicações parcimoniosas, já que os algoritmos raramente descartam variáveis. Além disso, foi impossível identificar diferenças culturais e mudanças em termos de valores com as técnicas de inferência utilizadas. A presença de praticamente as mesmas variáveis com importância em todas as ondas do Brasil (ver Figura 17) ilustra a dificuldade de se identificarem mudanças com essas técnicas.

Dito isso, o uso das ferramentas tem utilidade, já que consegue, além de identificar relações conhecidas, adicionar complexidade e outras informações que podem não ser intuitivas ou deixadas de lado pelos cientistas sociais em um primeiro momento.

Os exemplos de modelagem indutiva e dedutiva apresentam os dois extremos que devem ser aproximados por uma experiência própria do cientista social, na medida em que é necessário tanto conseguir enriquecer os modelos dedutivos de forma a evitar modelos fracos demais quanto poder selecionar (reduzir ou aumentar) as dimensões de uma amostra de dados usada no aprendizado de máquina.

O cientista social aparece como um auxiliar tanto na engenharia dos dados, como sugerido na literatura das ciências da computação, quanto na interpretação dos resultados, no que diz respeito à robustez das inferências. O aprendizado de máquina, em comparação com a modelagem estatística explicativa tem a vantagem de ser desnecessário conhecer *a priori* as teorias e as variáveis de modo a elaborar um modelo dedutivo. Uma vez os dados disponíveis, as técnicas de aprendizado de máquina permitem uma abordagem “*kitchen sink*” rapidamente. Por outro lado, é necessário conhecimento no momento de fazer inferências, e esse conhecimento parece essencial também na engenharia das variáveis explicativas. Tendo em vista o contexto atual de ampla disponibilidade de dados, tanto estruturados quanto não estruturados, na forma de *big-data* ou não, a limitação das ferramentas existentes e a complexidade dos problemas a serem resolvidos, os resultados apontam para a necessidade da interdisciplinaridade quando o assunto envolve inteligência artificial e ciências sociais.

Técnicas de IA são difíceis de serem implementadas, mas não exigem conhecimento prévio da teoria social. Já o conhecimento de teoria social nem sempre vem acompanhado

do ferramental necessário para fazer análises empíricas. O caminho parece ser, portanto, que a pesquisa seja realizada em grupos que saibam usar as ferramentas e que conheçam a teoria social. A entrada dos dados pode ser melhorada com a teoria, e os resultados das inferências devem fazer sentido para um cientista político ou social. O cientista político pode em seguida aprimorar e testar teorias existentes ou ser capaz de elaborar novas teorias e buscar sua comprovação na realidade.

LIMITAÇÕES

A limitação com relação ao entendimento do patriotismo e suas diversas nuances já foi explicitada anteriormente. Apesar de os resultados terem contribuído para trazer alguma clareza, é necessário reforçar essa limitação. Os resultados empíricos revelaram relações de diversas naturezas com o orgulho nacional e não se pode analisar de forma aprofundada todas elas.

Aqui considera-se que as declarações fornecidas em um *survey* são suficientemente estáveis e podem ser comparadas entre indivíduos e culturas. Mas o uso de *surveys* como a fonte de dados principal tem importantes limitações. Pode-se destacar aqui o fato de que *surveys*, dentro da tradição hipotético-dedutiva de produzir conhecimento, são elaboradas de forma a poder coletar dados empíricos para uma determinada teoria. Em outras palavras, os dados recolhidos tendem a servir para provar uma teoria em particular. Esse pode ser o caso do pós-materialismo no teste apresentado nos capítulos anteriores. Quanto ao patriotismo enquanto orgulho da nacionalidade, os dados parecem servir não somente para provar a tese de Iglehart, mas diversas outras, pela forma como essa pergunta é usada nos questionários da psicologia social.³

Com relação ao aprendizado de máquina, o estudo se limitou a experimentar o aprendizado de máquina supervisionado. Além disso, tratou-se de um problema de classificação binária. Aprender a partir dos dados de forma não supervisionada também é possível, desejável tendo em vista a riqueza de dados do WVS. Além disso, é possível usar

³Outras delas são as instabilidades ou variabilidades em determinadas respostas (ZALLER, 1992), os efeitos de autoafirmação (BAUMEISTER, 1982), vieses de resposta para o socialmente desejável (*socially desirable*) ou para demonstrar um comportamento satisfatório para o entrevistador (*satisficing*) (CHANG; KROSNICK, 2009), a tentativa de causar boa impressão (TEDESCHI; SCHLENKER; BONOMA, 1971; GECAS, 1982), ou, ainda, problemas semânticos nas questões, etc. Wimmer (2017) discute a influência das normas sociais nas declarações de orgulho. Para mais limitações, ver também Dillman, Phelps et al. (2009), Lohuizen e Samohyl (2011), Brehm (1993).

o aprendizado de máquina para resolver problemas de regressão (onde a variável resposta é contínua) ou problemas de classificação mais complexos.

Outras limitações dizem respeito às técnicas de aprendizado de máquina. Uma delas é a representação dos dados do WVS. A representação dos dados se refere à forma como as informações são usadas e o quanto dela consegue ser acessado e preservado durante o aprendizado (BENGIO; COURVILLE; VINCENT, 2012). Em geral, dados mal representados, produzem mau resultados. Neste trabalho, foram testadas poucas transformações nos dados e no enriquecimento das informações antes do treinamento, o que pode ter limitado as conclusões.

Outra, já mencionada acima, é a dificuldade em se fazerem inferências nos resultados de algoritmos “caixa-preta”. A explicação extraída dos algoritmos usados no trabalho, mesmo que eles tenham sido relativamente bons em termos de predição, é limitada. A inferência usando-se importância global indica um ranking de como cada variável afeta a resposta, mas não indica a direção. Além disso, não descarta nenhuma das variáveis dificultando a geração de explicações parcimoniosas. As técnicas locais, usando-se modelos de substituição, indicaram a direção dos efeitos, mas as variáveis são usadas de outra maneira o que reduz a confiança nos resultados.

Por fim, as inferências que puderam ser elaboradas são em nível individual, o que difere da abordagem geral da cultura política que agrega as análises em nível das sociedades ou culturas. A tentativa de agregar explicações locais de forma a dar uma resposta a esse problema precisa de uma fundamentação matemática e estatística mais sólida para que os resultados possam ser considerados como confiáveis e robustos.

Uma última limitação é externa e tem relação com o chamado *big-data*. Desde o surgimento das ciências sociais, a maioria das pesquisas sobre a sociedade e a interação entre indivíduos foi produzida com base em informações autodeclaradas uma única vez em *surveys* ou pesquisas de opinião, entrevistas em profundidade, observação participante etc. Novas tecnologias como a vídeovigilância, o e-mail e as mensagens instantâneas, crachás inteligentes, objetos conectados, e, obviamente, os telefones inteligentes com seus diversos sensores, produzem uma série de informações e dados em tempo real. Coletados ao longo do tempo, esses dados podem fornecer informações sobre a estrutura e o conteúdo das relações

sociais. Nesse contexto, a pesquisa foi limitada aos dados disponíveis e estruturados de *survey*.

PERSPECTIVAS

No que diz respeito ao fenômeno do patriotismo, as perspectivas à luz dos resultados estão ligadas à exploração dos resultados usando-se técnicas dedutivas de modo a confirmar ou descartar algumas das hipóteses que foram identificadas durante a indução com os dados.

Também podem ser identificadas algumas pistas ou perspectivas do uso do aprendizado de máquina para as ciências sociais. Uma das perspectivas que parece promissora é a pesquisa sobre a representação dos dados de *survey* e a compatibilidade com os algoritmos disponíveis. Por anos a participação de *domain experts* foi essencial *antes* do aprendizado de máquina para a engenharia dos dados. Com as redes neurais e o aprendizado profundo, passou a ser desnecessário, em um primeiro momento, devido ao uso de classificadores capazes de aprender e otimizar as próprias representações (LECUN; BENGIO; HINTON, 2015). Parece terem ficado claras no trabalho as diferenças que podem aparecer entre duas técnicas diferentes, que foram atribuídas a representação das variáveis. Um exemplo disso são as variáveis ordenadas e numéricas nas diferentes medidas de importância.

Não sem relação com a representação, existe a perspectiva de explorar o tamanho da base necessária para o treinamento. Neste trabalho, foram usados todos os dados possíveis. Em nenhum momento antes do treinamento foi avaliada a necessidade de se usarem amostras de cada *survey* ou país ou, ainda, outro tipo de amostragem mais informada. Por questões de tempo e capacidade de processamento, um estudo da taxa de aprendizado e os níveis onde se atinge a saturação do aprendizado com diversos algoritmos e com subamostras de tamanhos diferentes facilitaria o uso do WVS em pesquisas de cultura política.

Outra perspectiva relacionada à representação dos dados está relacionada com o *big-data*. A adaptação de domínios (BLITZER; DREDZE; PEREIRA, 2007), entre o WVS e outras bases de redes sociais ou sensores poderia ser explorada de forma que os aprendizados nessa base possam ser usados para prever (e explicar) bases de dados de natureza e origem diferentes.

A imputação dos dados é outra perspectiva interessante e importante. O algoritmo que foi selecionado tem embutido nele uma solução para as não-respostas. Esse não é o caso de

diversos outros algoritmos, que podem ter resultados melhores ou que possibilitem inferências mais claras. A imputação de dados em um contexto onde são usados conjuntos de dados para treinamento, validação e teste, é mais complexa e incerta. É preciso assegurar que as técnicas sejam adaptadas e que o “vazamento” seja controlado. A imputação múltipla aumenta a complexidade na medida em que requer que os modelos sejam treinados em diversas bases de dados imputadas em paralelo. Nesses casos, é preciso automatizar o treinamento com diversas bases controlando o erro entre as imputações. O aprendizado de máquina pode ajudar também na imputação de dados para análises de regressão.

Por fim, uma perspectiva interessante é derivada da experiência de agregação de explicações individuais, sejam elas feitas a partir das inferências locais, como experimentado aqui, sejam a partir das inferências globais. Estudos com base nesse tipo de agregação são úteis aos cientistas sociais pois permitem uma análise das variações entre culturas e mudanças ao longo do tempo. A agregação testada neste trabalho trouxe resultados interessantes, mas carece de fundamentação matemática e de mais exploração técnica. A criação de uma ferramenta que permita interpretar as inferências dos algoritmos para as ciências sociais pode ser uma pista de trabalho para uma equipe multidisciplinar em IA e ciências sociais.⁴

⁴Existem iniciativas em outras áreas; ver, por exemplo, Krause et al. (2017).

A. APÊNDICE A - AMOSTRA DE TRABALHO

Esse apêndice contém a descrição dos dados da amostra extraída do WVS para o trabalho de aprendizado, validação e testes.

Tabela 8: Amostra de Trabalho extraída do WVS por sociedade e onda. A amostra de trabalho é constituída por todas as observações que foram usadas para treinamento, validação e teste, uma vez selecionadas as variáveis mais recorrentes e de interesse para a pesquisa.

| Sociedade | Onda 1 1981-1984 | Onda 2 1989-1993 | Onda 3 1994-1998 | Onda 4 1999-2004 | Onda 5 2005-2009 | Onda 6 2010-2014 | Onda 7 2017-2019 | Total |
|--------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|---------------------|--------------|
| Albania | 0 | 0 | 999 | 1000 | 1534 | 0 | 1435 | 4968 |
| Algeria | 0 | 0 | 0 | 1282 | 0 | 1200 | 0 | 2482 |
| Andorra | 0 | 0 | 0 | 0 | 1003 | 0 | 1004 | 2007 |
| Azerbaijan | 0 | 0 | 2002 | 0 | 1505 | 1002 | 1800 | 6309 |
| Argentina | 1005 | 1002 | 1079 | 1280 | 1002 | 1030 | 1003 | 7401 |
| Australia | 1228 | 0 | 2048 | 0 | 1421 | 1477 | 1813 | 7987 |
| Austria | 0 | 1460 | 0 | 1522 | 1510 | 0 | 1644 | 6136 |
| Bangladesh | 0 | 0 | 1525 | 1500 | 0 | 0 | 1200 | 4225 |
| Armenia | 0 | 0 | 2000 | 0 | 1500 | 1100 | 1500 | 6100 |
| Belgium | 1145 | 2792 | 0 | 1912 | 1509 | 0 | 0 | 7358 |
| Bolivia | 0 | 0 | 0 | 0 | 0 | 0 | 2067 | 2067 |
| Bosnia Herzegovina | 0 | 0 | 0 | 1200 | 1512 | 0 | 0 | 2712 |
| Brazil | 0 | 1782 | 1143 | 0 | 1500 | 1486 | 1762 | 7673 |
| Bulgaria | 0 | 1034 | 1072 | 1000 | 2501 | 0 | 1560 | 7167 |
| SrpSka Republic | 0 | 0 | 400 | 0 | 0 | 0 | 0 | 400 |
| Belarus | 0 | 1015 | 2092 | 1000 | 1500 | 1535 | 1548 | 8690 |
| Canada | 1254 | 1730 | 0 | 1931 | 2164 | 0 | 0 | 7079 |
| Chile | 0 | 1500 | 1000 | 1200 | 1000 | 1000 | 1000 | 6700 |
| China | 0 | 1000 | 1500 | 1000 | 1991 | 2300 | 3036 | 10827 |
| Taiwan | 0 | 0 | 780 | 0 | 1227 | 1238 | 0 | 3245 |
| Colombia | 0 | 0 | 6025 | 0 | 3025 | 1512 | 1520 | 12082 |
| Croatia | 0 | 0 | 1196 | 1003 | 1525 | 0 | 1488 | 5212 |
| Cyprus | 0 | 0 | 0 | 0 | 2050 | 1000 | 0 | 3050 |

Tabela 8: Amostra de Trabalho extraída do WVS por sociedade e onda. A amostra de trabalho é constituída por todas as observações que foram usadas para treinamento, validação e teste, uma vez selecionadas as variáveis mais recorrentes e de interesse para a pesquisa. (continuação)

| Sociedade | Onda 1 | Onda 2 | Onda 3 | Onda 4 | Onda 5 | Onda 6 | Onda 7 | Total |
|----------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|--------------|
| | 1981-1984 | 1989-1993 | 1994-1998 | 1999-2004 | 2005-2009 | 2010-2014 | 2017-2019 | |
| Cyprus (T) | 0 | 0 | 0 | 0 | 500 | 0 | 0 | 500 |
| Czech Rep. | 0 | 3033 | 1147 | 1908 | 1821 | 0 | 1812 | 9721 |
| Denmark | 1182 | 1030 | 0 | 1023 | 1507 | 0 | 3362 | 8104 |
| Dominican Rep. | 0 | 0 | 417 | 0 | 0 | 0 | 0 | 417 |
| Ecuador | 0 | 0 | 0 | 0 | 0 | 1202 | 1200 | 2402 |
| El Salvador | 0 | 0 | 1254 | 0 | 0 | 0 | 0 | 1254 |
| Ethiopia | 0 | 0 | 0 | 0 | 1500 | 0 | 0 | 1500 |
| Estonia | 0 | 1008 | 1021 | 1005 | 1518 | 1533 | 1304 | 7389 |
| Finland | 1003 | 588 | 987 | 1038 | 2148 | 0 | 1199 | 6963 |
| France | 1200 | 1002 | 0 | 1615 | 2502 | 0 | 1870 | 8189 |
| Georgia | 0 | 0 | 2008 | 0 | 3000 | 1202 | 2194 | 8404 |
| Palestine | 0 | 0 | 0 | 0 | 0 | 1000 | 0 | 1000 |
| Germany | 1305 | 3437 | 2026 | 2036 | 4139 | 2046 | 6935 | 21924 |
| Ghana | 0 | 0 | 0 | 0 | 1534 | 1552 | 0 | 3086 |
| Greece | 0 | 0 | 0 | 1142 | 1500 | 0 | 1200 | 3842 |
| Guatemala | 0 | 0 | 0 | 0 | 1000 | 0 | 0 | 1000 |
| Haiti | 0 | 0 | 0 | 0 | 0 | 1996 | 0 | 1996 |
| Hong Kong | 0 | 0 | 0 | 0 | 1252 | 1000 | 2075 | 4327 |
| Hungary | 1464 | 999 | 650 | 1000 | 2520 | 0 | 1516 | 8149 |
| Iceland | 927 | 702 | 0 | 968 | 808 | 0 | 2506 | 5911 |
| India | 0 | 2500 | 2040 | 2002 | 2001 | 4078 | 0 | 12621 |
| Indonesia | 0 | 0 | 0 | 1000 | 2015 | 0 | 3200 | 6215 |
| Iran | 0 | 0 | 0 | 2532 | 2667 | 0 | 0 | 5199 |
| Iraq | 0 | 0 | 0 | 2325 | 2701 | 1200 | 1200 | 7426 |
| Ireland | 1217 | 1000 | 0 | 1012 | 1013 | 0 | 0 | 4242 |
| Israel | 0 | 0 | 0 | 1199 | 0 | 0 | 0 | 1199 |
| Italy | 1348 | 2018 | 0 | 2000 | 2531 | 0 | 2277 | 10174 |
| Japan | 1204 | 1011 | 1054 | 1362 | 1096 | 2443 | 0 | 8170 |
| Kazakhstan | 0 | 0 | 0 | 0 | 0 | 1500 | 1277 | 2777 |
| Jordan | 0 | 0 | 0 | 1223 | 1200 | 1200 | 1203 | 4826 |
| South Korea | 970 | 1251 | 1249 | 1200 | 1200 | 1200 | 1245 | 8315 |
| Kuwait | 0 | 0 | 0 | 0 | 0 | 1303 | 0 | 1303 |
| Kyrgyzstan | 0 | 0 | 0 | 1043 | 0 | 1500 | 0 | 2543 |
| Lebanon | 0 | 0 | 0 | 0 | 0 | 1200 | 1200 | 2400 |
| Latvia | 0 | 903 | 1200 | 1013 | 1506 | 0 | 0 | 4622 |

Tabela 8: Amostra de Trabalho extraída do WVS por sociedade e onda. A amostra de trabalho é constituída por todas as observações que foram usadas para treinamento, validação e teste, uma vez selecionadas as variáveis mais recorrentes e de interesse para a pesquisa. (continuação)

| Sociedade | Onda 1 | Onda 2 | Onda 3 | Onda 4 | Onda 5 | Onda 6 | Onda 7 | Total |
|--------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|--------------|
| | 1981-1984 | 1989-1993 | 1994-1998 | 1999-2004 | 2005-2009 | 2010-2014 | 2017-2019 | |
| Libya | 0 | 0 | 0 | 0 | 0 | 2131 | 0 | 2131 |
| Lithuania | 0 | 1000 | 1009 | 1018 | 1500 | 0 | 1448 | 5975 |
| Luxembourg | 0 | 0 | 0 | 1211 | 1610 | 0 | 0 | 2821 |
| Malaysia | 0 | 0 | 0 | 0 | 1201 | 1300 | 1313 | 3814 |
| Mali | 0 | 0 | 0 | 0 | 1534 | 0 | 0 | 1534 |
| Malta | 467 | 393 | 0 | 1002 | 1500 | 0 | 0 | 3362 |
| Mexico | 1837 | 1531 | 1510 | 1535 | 1560 | 2000 | 1741 | 11714 |
| Moldova | 0 | 0 | 984 | 1008 | 2597 | 0 | 0 | 4589 |
| Montenegro | 0 | 0 | 240 | 1060 | 1516 | 0 | 0 | 2816 |
| Morocco | 0 | 0 | 0 | 1251 | 1200 | 1200 | 0 | 3651 |
| Netherlands | 1221 | 1017 | 0 | 1003 | 2604 | 1902 | 2721 | 10468 |
| New Zealand | 0 | 0 | 1201 | 0 | 954 | 841 | 0 | 2996 |
| Nigeria | 0 | 1001 | 1996 | 2022 | 0 | 1759 | 1237 | 8015 |
| Norway | 1051 | 1239 | 1127 | 0 | 2115 | 0 | 1123 | 6655 |
| Pakistan | 0 | 0 | 733 | 2000 | 0 | 1200 | 2000 | 5933 |
| Peru | 0 | 0 | 1211 | 1501 | 1500 | 1210 | 1400 | 6822 |
| Philippines | 0 | 0 | 1200 | 1200 | 0 | 1200 | 0 | 3600 |
| Poland | 0 | 1920 | 1153 | 1095 | 2510 | 966 | 1352 | 8996 |
| Portugal | 0 | 1185 | 0 | 1000 | 1553 | 0 | 0 | 3738 |
| Puerto Rico | 0 | 0 | 1164 | 720 | 0 | 0 | 1127 | 3011 |
| Qatar | 0 | 0 | 0 | 0 | 0 | 1060 | 0 | 1060 |
| Romania | 0 | 1103 | 1239 | 1146 | 3265 | 1503 | 2870 | 11126 |
| Russia | 0 | 1961 | 2040 | 2500 | 3537 | 2500 | 3635 | 16173 |
| Rwanda | 0 | 0 | 0 | 0 | 1507 | 1527 | 0 | 3034 |
| Saudi Arabia | 0 | 0 | 0 | 1502 | 0 | 0 | 0 | 1502 |
| Serbia | 0 | 0 | 1280 | 1200 | 2732 | 0 | 2706 | 7918 |
| Singapore | 0 | 0 | 0 | 1512 | 0 | 1972 | 0 | 3484 |
| Slovakia | 0 | 1602 | 1095 | 1331 | 1509 | 0 | 1435 | 6972 |
| Vietnam | 0 | 0 | 0 | 1000 | 1495 | 0 | 0 | 2495 |
| Slovenia | 0 | 1035 | 1007 | 1006 | 2403 | 1069 | 1076 | 7596 |
| South Africa | 1596 | 2736 | 2935 | 3000 | 2988 | 3531 | 0 | 16786 |
| Zimbabwe | 0 | 0 | 0 | 1002 | 0 | 1500 | 0 | 2502 |
| Spain | 2303 | 4147 | 1211 | 2409 | 2700 | 1189 | 1211 | 15170 |
| Sweden | 1908 | 1047 | 1009 | 1015 | 2190 | 1206 | 1194 | 9569 |
| Switzerland | 0 | 1400 | 1212 | 0 | 2513 | 0 | 3660 | 8785 |

Tabela 8: Amostra de Trabalho extraída do WVS por sociedade e onda. A amostra de trabalho é constituída por todas as observações que foram usadas para treinamento, validação e teste, uma vez selecionadas as variáveis mais recorrentes e de interesse para a pesquisa. (continuação)

| Sociedade | Onda 1 | Onda 2 | Onda 3 | Onda 4 | Onda 5 | Onda 6 | Onda 7 | Total |
|---------------------|--------------|--------------|--------------|---------------|---------------|--------------|---------------|---------------|
| | 1981-1984 | 1989-1993 | 1994-1998 | 1999-2004 | 2005-2009 | 2010-2014 | 2017-2019 | |
| Thailand | 0 | 0 | 0 | 0 | 1534 | 1200 | 1500 | 4234 |
| Trinidad and Tobago | 0 | 0 | 0 | 0 | 1002 | 999 | 0 | 2001 |
| Tunisia | 0 | 0 | 0 | 0 | 0 | 1205 | 1208 | 2413 |
| Turkey | 0 | 1030 | 1907 | 4607 | 3730 | 1605 | 0 | 12879 |
| Uganda | 0 | 0 | 0 | 1002 | 0 | 0 | 0 | 1002 |
| Ukraine | 0 | 0 | 2811 | 1195 | 2507 | 1500 | 0 | 8013 |
| Macedonia | 0 | 0 | 995 | 1055 | 1500 | 0 | 0 | 3550 |
| Egypt | 0 | 0 | 0 | 3000 | 3051 | 1523 | 1200 | 8774 |
| United Kingdom | 1167 | 1484 | 1093 | 1000 | 2602 | 0 | 1788 | 9134 |
| Tanzania | 0 | 0 | 0 | 1171 | 0 | 0 | 0 | 1171 |
| United States | 2325 | 1839 | 1542 | 1200 | 1249 | 2232 | 2596 | 12983 |
| Burkina Faso | 0 | 0 | 0 | 0 | 1534 | 0 | 0 | 1534 |
| Uruguay | 0 | 0 | 1000 | 0 | 1000 | 1000 | 0 | 3000 |
| Uzbekistan | 0 | 0 | 0 | 0 | 0 | 1500 | 0 | 1500 |
| Venezuela | 0 | 0 | 1200 | 1200 | 0 | 0 | 0 | 2400 |
| Yemen | 0 | 0 | 0 | 0 | 0 | 1000 | 0 | 1000 |
| Zambia | 0 | 0 | 0 | 0 | 1500 | 0 | 0 | 1500 |
| North Ireland | 312 | 304 | 0 | 1000 | 500 | 0 | 0 | 2116 |
| Bosnia | 0 | 0 | 800 | 0 | 0 | 0 | 0 | 800 |
| Kosovo | 0 | 0 | 0 | 0 | 1601 | 0 | 0 | 1601 |
| Total | 30639 | 62771 | 77818 | 100155 | 151761 | 89565 | 105696 | 618405 |

^a Fonte: Base Longitudinal do WVS, Janeiro 2020

^b Total de observações para cada onda e sociedade

B. APÊNDICE B - NÃO-RESPOSTAS E IMPUTAÇÃO

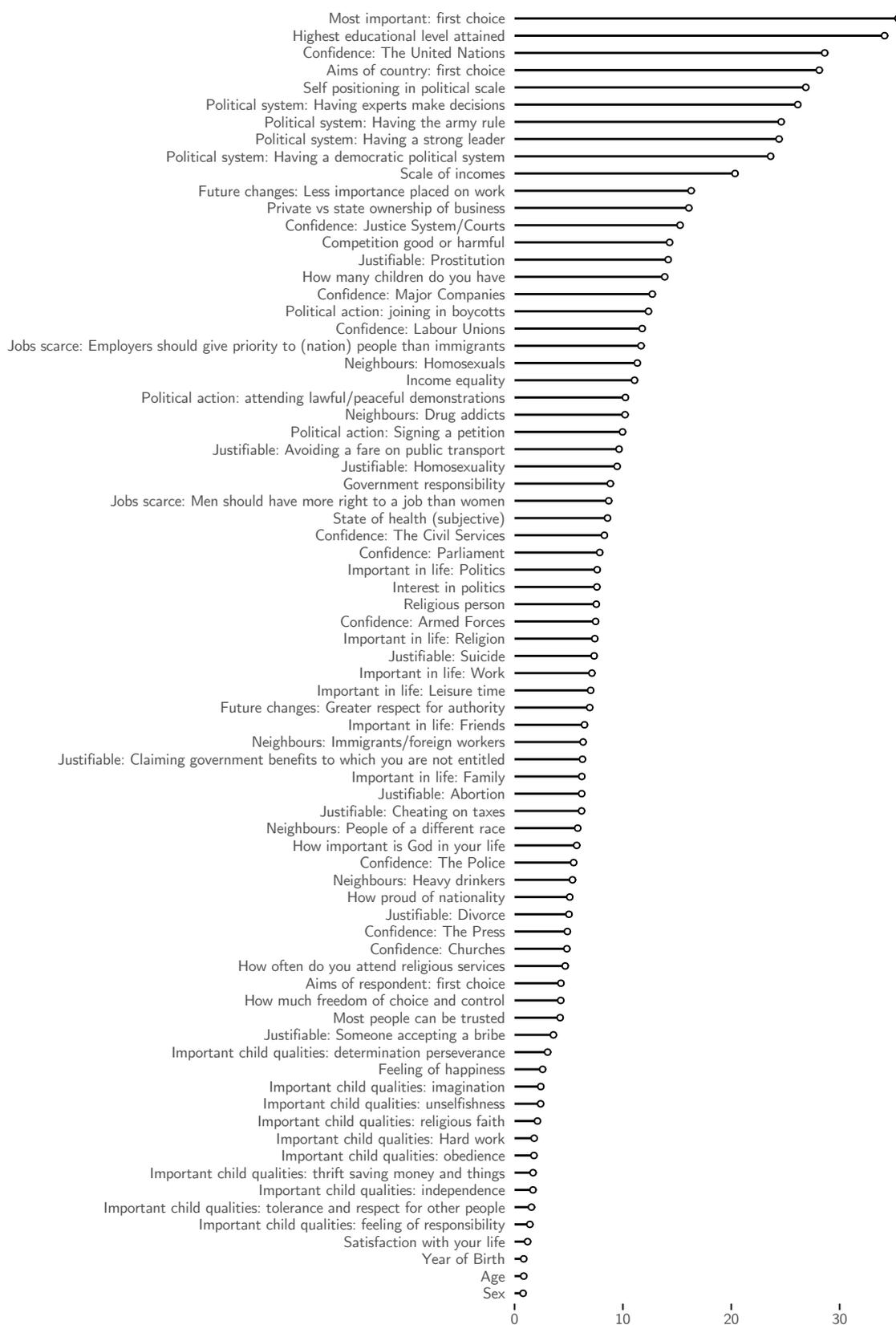
As não-respostas e a Imputação dos dados

Em pesquisas tipo survey, não-respostas são valores faltantes ou ausentes (*missing values*) que, no campo da estatística, fazem parte do conjunto de dados brutos, que sofreu agrupamentos, agregações, arredondamentos, exclusões, resultando na perda parcial das informações (SCHAFER; GRAHAM, 2002). A perda de informações indica que os dados estão incompletos, pois as medidas desejadas não foram conseguidas.

As não-respostas são classificadas de acordo com o nível em que se encontram em surveys, seja em nível das unidades ou de itens (*unit* ou *item nonresponse*). O primeiro se refere a quando uma pessoa que deveria ter respondido uma survey não estava disponível, ou não quis reponder; o segundo quando a pessoa deixa de responder uma ou mais perguntas do questionário (SALANT; DILLMAN, 1994; DILLMAN; SMYTH; CHRISTIAN, 2014; DILLMAN; PHELPS et al., 2009). Assim, o erro relacionado as não-respostas é um indicador de qualidade de uma pesquisa tipo survey (DILLMAN; SMYTH; CHRISTIAN, 2014). O erro relacionado as não-resposta é um indicador de qualidade de uma pesquisa tipo survey (DILLMAN; SMYTH; CHRISTIAN, 2014).

A maneira mais evidente de tratar o problema das não-respostas é descartar ou ignorar uma parte dos dados, seja usando apenas os casos completos ou então ignorando algumas das observações ou variáveis onde as taxas de não-respostas são mais altas (São chamados na literatura de *listwise*, *case deletion*, *pairwise deletion*, *available case analysis* ou ainda *complete case analysis*).

Figura 18: Taxa de não-respostas (NA) na base de trabalho. O gráfico mostra as 81 variáveis que foram mantidas na base (eixo y) e a taxa de não-respostas em percentagem no eixo x



Outra é imputar os dados. A imputação é o processo de substituir dados faltantes por valores estimados. O principal objetivo de uma imputação é reduzir o viés da amostra, que ocorre porque os dados estão incompletos, ou seja, porque a distribuição das não-respostas difere da distribuição observada (DURRANT, 2005).

Existem diversas formas de imputar os dados. Uma série de métodos simplificados são conhecidos como métodos de imputação única. Nesta categoria estão as técnicas de imputação pela média, moda ou pela mediana, para dados contínuos e categóricos respectivamente, ou ainda as técnicas *hot-deck* e *cold-deck*, que tratam de substituir os valores faltantes por um valor aleatório preservando a distribuição existente (ANDRIDGE; LITTLE, 2010). Estas técnicas são usadas de forma recorrente devido à sua facilidade de implementação, mas possuem limitações importantes (SCHAFER, 1997).

Também pode ser categorizados como métodos de imputação única e simples, os métodos que fazem imputações baseados em modelos dedutivos, isto é, que usam relações lógicas entre variáveis; métodos que usam regressões lineares; ou ainda métodos baseados na regra do vizinho mais próximo (*knn*)¹.

Em imputações únicas, uma vez imputados os dados, considera-se os valores imputados como observados. Ignora-se assim a incerteza do processo de imputação, o que por sua vez adiciona incerteza nas predições ou estimativas produzidas (SCHAFER; GRAHAM, 2002).

Para contornar esse problema, duas formas de tratar dados faltantes baseadas em modelos estatísticos foram desenvolvidas nos anos 1970. O primeiro deles é fundado na busca por parâmetros que maximizam a função de verossimilhança (*maximum likelihood*). Esses modelos estimam, por exemplo, a média e a variância, supondo uma distribuição normal por exemplo. Uma vez estimados, esses parâmetros podem ser usados para estimar valores ausentes e preencher as lacunas da base original. (LITTLE, R. J. A.; RUBIN, D. B., 2014; LITTLE, R.; RUBIN, D., 1987; DEMPSTER; LAIRD; RUBIN, 1977).

A imputação múltipla(IM) é uma técnica de modelagem estatística desenvolvida especificamente para imputação. A IM foi introduzida por Rubin (1976, 1978; 1987) e se trata de uma abordagem usando técnicas de Monte Carlo para gerar múltiplas versões dos

¹Esse método é usado também no aprendizado automático para problemas de agrupamento de dados e redução de variáveis (DURRANT, 2005; HARRELL, 2001)

dados imputados. Uma vez imputados, podem ser analisados de forma agregada para gerar inferências estatísticas (LITTLE, R.; RUBIN, D., 1987; HASTIE; TIBSHIRANI; FRIEDMAN, 2017).

A IM é vista atualmente como a técnica estado da arte para imputação dos dados, pois melhora a acurácia e o poder estatístico com relação a outras técnicas. Honaker e King (2010) afirma que os cientistas sociais tem evitado métodos que envolvem descartar casos incompletos ou imputações baseadas em *best-guesses*, preferindo o uso de IM.

Dentro da imputação múltipla, existem duas maneiras de estimar o modelo estatístico. O primeiro deles é a abordagem de modelagem conjunta (*Joint Modeling*) onde um modelo probabilístico único é estimado para imputar todas as variáveis. Um exemplo é o software Amelia II (HONAKER; KING; BLACKWELL, 2011) que assume um modelo gaussiano multivariado para imputação. A segunda abordagem é de modelagem condicional, que gera um modelo para cada variável a ser imputada².

Segundo Harrell Jr. (2001, pp. 49-50), é possível escolher os métodos de imputação de acordo com a proporção de dados faltantes em alguma das variáveis:

1. Proporção 0,05: Nesse caso pode ser usada a imputação única ou analisar somente os dados completos;
2. Proporção entre 0,05 e 0,15: Imputação única pode ser usada provavelmente sem problemas, entretanto o uso da imputação múltipla é indicado;
3. Proporção 0,15: A imputação múltipla é indicada na maior parte dos modelos.

Na perspectiva de Rubin, a causa da ausência é mais importante que a quantidade deles na tomada de decisão sobre como imputá-los (LITTLE, R.; RUBIN, D., 1987). É a correta definição da causa que permite fazer inferências estatísticas válidas, uma vez imputados os dados.

Em uma abordagem estatística moderna, proposta por Rubin, os dados faltantes são um fenômeno probabilístico: a probabilidade de observar a ocorrência ou não de uma resposta é descritas por uma função, assumindo uma distribuição particular entre os dados observados

²Além destas duas formas, existem outras técnicas de imputação múltipla como a imputação fracionada e a imputação múltipla (IM) baseada em regressões múltiplas. Ver Durrant (2005) e van_Buuren (2012)

e faltantes (SCHAFER; GRAHAM, 2002). De acordo com Rubin (1976), as causas ou as relações existentes entre respostas e não-respostas tem “mecanismos” diferentes: perdas completamente ao acaso (*missing completely at random* – MCAR), onde a probabilidade de valores faltantes é a mesma para todos os casos, ou seja, as causas são aleatórias. Perdas ao acaso (*missing at random* – MAR), menos restritivo, se refere aos casos onde as não-respostas estão também ausentes de forma aleatória, mas a probabilidade está relacionada com os valores de outras variáveis observadas. Enfim, existem os casos onde as perdas são não-aleatórias (*not missing at random* – NMAR, *missing not at random* - MNAR, ou *nonignorable NI*) que correspondem aos casos onde a probabilidade pode estar relacionada ao valor faltante em si, que não está ao alcance do pesquisador.

Para Rubin, além da correta definição das causas, os métodos de IM estavam previstos quando os dados faltantes não ultrapassavam 60% (ROBINS; WANG, 2000).

Em geral, a presunção para a maioria das imputações é MCAR (HASTIE; TIBSHIRANI; FRIEDMAN, 2017). Pode-se considerar, na melhor das hipóteses, que as perdas seriam ao acaso (MAR) (DURRANT, 2005; SCHAFER; OLSEN, 1998). Assumindo-se esse padrão os métodos de IM não geram viéses (LITTLE, R.; RUBIN, D., 1987; LITTLE, R. J. A.; RUBIN, D. B., 2014).

De acordo com a classificação em pesquisas do tipo survey não-respostas podem estar em nível das unidades e de itens. O primeiro em geral se refere a quando uma pessoa que deveria ter respondido uma survey não estava disponível, ou não quis responder; o segundo quando a pessoa deixa de responder uma ou mais perguntas do questionário (SALANT; DILLMAN, 1994; DILLMAN; SMYTH; CHRISTIAN, 2014; DILLMAN; PHELPS et al., 2009).

No caso do WVS, as não-respostas em nível da unidade (*unit nonresponse*), ou os membros da população que não foram entrevistados, que podem ser consideradas como perdas completamente ao acaso (MCAR), são corrigidos internamente usando pesos em cada uma das *surveys* realizadas ou técnicas de resampling dentro da população. Dessa forma, as inferências estatísticas para a população não tem viés para a sub-amostra, e os não-entrevistados não tem impacto nas estimativas estatísticas (KING, G. et al., 2001).

No caso da não-respostas em nível dos itens, significa de forma intuitiva que os indivíduos que não responderem alguma questão de uma survey são aleatoriamente diferentes

dos demais, e que para uma observação qualquer onde o valor de determinada variável esta ausente pode-se estimar o seu valor a partir dos dados existentes. O impacto de imputar os dados faltantes acaba sendo positivo nesses casos (FITZMAURICE et al., 2014).

Isso não exclui a possibilidade de os dados faltantes em determinadas perguntas se tratarem de perdas não-aleatórias (MNAR). As pessoas deixam de responder surveys por diversas razões, que nos são desconhecidas. A possibilidade de uma não-resposta sobre a renda por exemplo pode ser afetada pela própria renda do indivíduo: pessoas mais pobres podem ter vergonha de responder, e os mais ricos serem relutantes em passar a informação (KING, G. et al., 2001). Além disso, em pesquisas de opinião, pessoas com opiniões mais nuanciadas tendem a responder menos ou se abster mais vezes, indivíduos mais velhos respondem mais do que os jovens, etc. (DILLMAN; SMYTH; CHRISTIAN, 2014).

Considerar os dados MNAR tornaria o mecanismo de não-respostas *não-ignorável* e imputar os dados considerando-os dessa forma tornaria o procedimento mais complicado devido às inúmeras premissas (e novas informações) que teriam que ser agregadas para cada variável a ser imputada.

Outra dificuldade relativa à tomada de decisão pela imputação ou não dos dados está relacionada com a distribuição das não-respostas na base longitudinal do WVS. Como mencionado acima, a base é composta por diversas survey realizadas em vários países ao longo dos anos ($n = 424$). Devido ao caráter descentralizado no qual é realizada, o desenho das pesquisas não foi sempre idêntico entre uma onda e outra. Nem todas as perguntas foram aplicadas em todos os países, e algumas perguntas sofreram modificações durante os anos, ou foram simplesmente abandonadas. Existem assim dados ausentes tanto dentro de cada survey onda-país, quanto entre as diferentes surveys do conjunto. A melhor solução nesses casos seria imputar as surveys em um modelo hierárquico e depois combiná-los em uma matriz unica, como sugerido por Gelman, King e Liu (1998). Nesse trabalho, optou-se por uma imputação múltipla desconsiderando a hierarquia dos dados.

A quantidade de bases a imputar é outra dificuldade. Originalmente o número de bases imputadas proposto por Rubin era de 5. De forma a obter melhores estimativas de erro padrão geralmente é necessário gerar mais bases. Uma indicação é usar a porcentagem das não-respostas como base para o número de imputações. O que torna a IM rapidamente

ineficiente do ponto de vista das simulações dependendo da variável a ser utilizadas.

Outro critério importante na escolha da técnica de imputação está relacionado às premissas sobre a distribuição das variáveis, que devem ser estudadas de modo a poder parametrizar os modelos de imputação. Por exemplo, Honaker e King (2010) assumem uma distribuição normal multivariada, que implica que os dados são imputados de forma linear. Em casos de não-linearidade importante, o autor sugere uma abordagem única, ou então que os dados sejam transformados antes da imputação.

A imputação múltipla também dificulta a modelagem preditiva pois o número de vezes que a base deveria ser treinada aumenta em um efeito cascata. Mesmo contraintuitivo, é recomendado incluir a variável de supervisão na imputação porque ela pode ajudar na previsão de outras variáveis (ALLISON, 2002; JOSSE et al., 2019).

Honaker e King (2010) propuseram um algoritmo para IM específico para uso nas ciências sociais (Amelia II). Haggard, Kaufman e Long (2013), Adamczyk e Pitt (2009) e Ciftci (2010) usaram imputação múltipla na análise de dados de survey. Chen e Shao (2000) demonstrou a utilidade de técnicas de Nearest neighbourhood para imputar surveys em geral. Batista e Monard (2003) comparou uma abordagem de AM comparada com algoritmos C4.5 e CN2, com imputações múltiplas e simples pela moda e média.

No caso da modelagem dedutiva, a imputação tem consequência nas conclusões com relação à teoria explicativa. Não-respostas podem criar vieses e inflar a taxa de erros tipo I (*i.e.* rejeitar incorretamente uma hipótese nula verdadeira); e II (*i.e.* não aceitar uma hipótese alternativa verdadeira), devido à interferência nos intervalos de confiança. Na modelagem preditiva, as consequências da imputação são medidas no momento da predição.

De qualquer forma, imputar significa correr alguns riscos (HASTIE; TIBSHIRANI; FRIEDMAN, 2017). Como dizem Little e Rubin “*The idea of imputation is both seductive and dangerous*” (LITTLE, R.; RUBIN, D., 1987). Na pior das hipóteses, seria uma alternativa melhor do que descartar os casos (KING, G. et al., 2001). Isso porque a suposição de dados MNAR é raramente possível, e descartar dados geram vieses quando a quantidade de dados faltantes é grande. O exemplo usado pelos autores é o caso de as não-respostas não serem completamente aleatórios, mas relacionados ao fato de eleitores serem melhor educados, menos politizados e ou neutros em termos partidários.

C. APÊNDICE C - ESTUDO COMPARATIVO

Ao total, foram testados 19 classificadores, incluindo alguns algoritmos já conhecidos pelos seus bons resultados. A lista dos algoritmos testados se encontra na tabela 9.

A gama de algoritmos inclui modelos estatísticos e outros advindos da mineração de dados. Para além dos algoritmos individuais, foram testados meta-algoritmos que combinam diversos classificadores ou regressores. Esses estimadores são chamados *ensemble* porque combinam vários métodos de aprendizado de máquina em um modelo preditivo (SHMUELI, 2010). As técnicas principais de *ensemble* desenvolvidas são o *bagging*, usada para reduzir a variância dos modelos (BREIMAN, 1996), e o *boosting* que é usada para reduzir os viéses (FREUND; SCHAPIRE, 1997).

Pelos resultados do estudo, pode-se constatar em primeiro lugar que nenhum dos algoritmos chegou perto ao que pode se chamar de acurácia de base, ou a acurácia que seria atingida se o resultado mais comum (Alto Orgulho) fosse atribuído a todos os casos. Em se tratando de um estudo comparativo, onde o objetivo é encontrar o algoritmo mais competente, esse baixo desempenho foi ignorado.

Entre os algoritmos testados, os que tiveram melhor desempenho foram os classificadores *ensemble*. Não se trata de uma surpresa, pois esses algoritmos tem um registro de bons resultados na literatura (e.g. MARQUÉS; GARCÍA; SÁNCHEZ, 2012; TSAI; HSU; YEN, 2014; MACLIN; OPITZ, 1999).

O classificador melhor colocado foi o XGBoost, ou *eXtreme Gradient Boosting* (“xgb” na tabela 5), com uma acurácia média de 68% (+-1,2%) entre todos as combinações de hiperparâmetros, e a máxima em 72%. Esse foi o algoritmo que também atingiu o melhor desempenho usando a curva ROC como métrica. Em seguida, o algoritmo melhor colocado foi o Floresta Aleatória (rf), que atingiu uma acurácia média de 70,5% (+-0,05%), e máxima

de 71%. O Floresta Aleatória é também o segundo melhor colocado medindo-se a área abaixo da curva ROC.

Também tiveram bons resultados uma variante do XGBoost, que incorpora a técnica DART (xgbDART), que é usada em redes neurais artificiais para descartar aleatoriamente algumas árvores. O algoritmo de Boosting Gradiente (gbm), outra versão de um ensemble *boosting* também ficou entre os melhores resultados.

O desempenho dos algoritmos de árvores de decisão foi irregular. Apenas o C5.0 teve bom desempenho. Os demais, como as árvores de decisão condicional simples (rpart, ctree) tiveram desempenho inferior.

Dois algoritmos de redes neurais artificiais, um com apenas uma camada intermediária de neurônios (nnet), e um perceptron multicamadas do tipo *ensemble* com 4 camadas e 9 neurônios com pesos aleatórios (avNNET), tiveram desempenho praticamente idêntico, e bastante próximo dos melhores colocados. O baixo desempenho do algoritmo Perceptron de Multicamadas de Propagação Inversa (mlpML) e da máquina de aprendizado extremo (elmNN), surpreenderam. Esperava-se melhores resultados em ambos algoritmos mais sofisticados de redes neurais artificiais. O baixo desempenho talvez esteja relacionados a busca por bons hiperparâmetros. Talvez esses algoritmos tenham melhor desempenho na resolução de problemas que envolvam dados não-estruturados (imagens, texto, etc.).

Os dois algoritmos que tentam extrair regras a partir dos dados (JRip e OneR, do pacote RWeka) tiveram baixo desempenho. O *One Rule* foi o algoritmo que obteve o pior resultado entre todos os algoritmos testados. As abordagens bayesianas como o modelo generalizado (bglm) e o *Naïve Bayes* (nb) tiveram desempenho intermediário, assim como o modelo linear generalizado com penalidade (glmnet) e a regressão logística (Log) que foram testados.

Acurácia na amostra de validação

Para verificar se o modelo tem sobreajuste ou não, ou seja, se ele tem capacidade de generalização em dados novos, foram comparados os resultados da Acurácia na amostra de casos completos da base de validação ($n = 13497$). Na figura 8, no capítulo dos resultados, pode-se constatar que o XGBoost teve acurácia um pouco maior na amostra de validação. O mesmo foi notado em um dos algoritmos de Redes Neurais Artificiais (avNNET). Os

algoritmos Floresta Aleatória (rf), Boosting Gradiente (gbm) e Boosting Gradiente Extremo DART (xgbDART) tiveram praticamente o mesmo desempenho em ambas as amostras. Os demais tiveram desempenho inferior na base de validação, com destaque ao C5.0 que teve boa performance na amostra de treinamento, e um pouco menor na amostra de validação, indicando que pode haver algum sobreajuste.

A partir desses resultados, foi possível identificar os algoritmos com melhor desempenho. Entretanto, a diferença entre os primeiros colocados foi ainda pequena, dificultando a decisão. Como os algoritmos foram treinados com os mesmos parâmetros pode-se empregar métodos estatísticos para determinar se o desempenho de um modelo supera outro. Nesse caso, a hipótese nula é que a diferença observada entre os desempenhos dos algoritmos é aleatória (EUGSTER, Manuel J A; HOTHORN; LEISCH, 2008; DEMSAR, 2006; DIETTERICH, 1998). Para um nível de significância ($\alpha = 0,05$) o teste de Friedman encontrou evidências para descartar a hipótese nula de que todos os algoritmos conseguem a mesma acurácia, κ ou mesmo ROC (Teste de Friedman = 1685.77).

Entretanto, esse teste não determina qual ou quais os pares que são estatisticamente diferentes. Para poder identificá-los deve se recorrer à uma análise *post hoc*. Assumindo uma distribuição não-normal, menos restritiva, um teste não-paramétrico é recomendado. O teste de Wilcoxon-Mann-Whitney por pares revelou um *p-value* abaixo do $\alpha = 0,05$ para a maioria dos pares, indicando também que as diferenças entre os modelos são estatisticamente significantes. O *p-value* acima do $\alpha = 0,05$ indica entretanto que alguns dos pares dos estimadores melhores colocados (e.g. C5.0, gbm e xgbDART) são muito próximos.

Entre os dois melhores colocados, o XGBoost (xgb) e o Floresta Aleatória (rf), o teste revelou que a diferença entre eles é estatisticamente significativa (*p-value* = 3.88e-15). A seção seguinte apresenta os resultados do treinamento usando o XGBoost, usando esse critério.

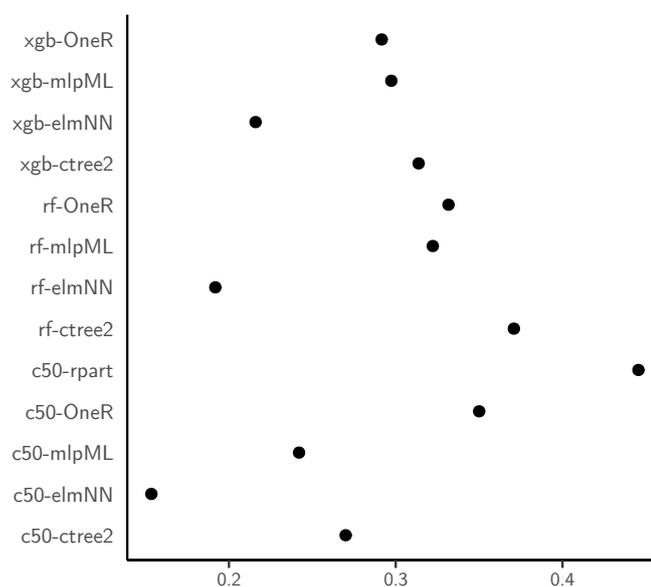
Stacking ou empilhamento

Também foram testados, sem muito sucesso, algumas formas de empilhamento de algoritmos (*stacking*), que serve para combinar a capacidade preditiva de diversos algoritmos (WOLPERT, 1992). A chave para realizar esta combinação é a diversidade entre os algoritmos. Por diversidade, entende-se a capacidade de prever casos distintos: quando um dos algoritmo não é capaz de classificar uma determinada observação, o outro o seria. Assim, quando

combinados, teriam melhor desempenho.

A possibilidade de serem combinados em um *ensemble* de empilhamento foi estudada a partir da correlação entre as predições dos diferentes algoritmos. Quanto menor a correlação entre elas, maior a chance de formarem um bom ensemble. Os testes revelaram uma correlação relativamente alta (0,722) entre o Extreme Boosting (xgbTree) e o Random Forest (rf), os dois melhores colocados, indicando que uma combinação de ambos os algoritmos não seria tão benéfica quanto a escolha de um deles como algoritmo principal. Também foram testadas a combinação desses algoritmos com outros de menor capacidade preditiva para o problema proposto, mas os resultados obtidos não foram bons e omitimos uma descrição mais detalhada aqui. Além disso, a interpretação dos resultados em um empilhamento dificultaria o trabalho de interpretação dos dados.

Figura 19: Correlação entre as predições dos algoritmos de melhor desempenho. O gráfico mostra as mais baixas correlações (<0,5) entre as predições dos três algoritmos de melhor desempenho, e famílias distintas (XgbTree, RF e C50), e os demais. No eixo y estão as combinações e no eixo x a escala de correlação de 0-1, sendo 1 a correlação total



Como pode ser visto na figura 19, a análise das correlações revelou um potencial de combinação entre os algoritmos de melhor desempenho com os algoritmos de redes neurais artificiais mlpML e elmNN, regras de decisão OneR, e árvores de de decisão CTree e CART. Tendo em vista o baixo desempenho individual do OneR e do elmNN, a inclusão desses algoritmos pode não trazer uma grande vantagem preditiva. Os algoritmos CTree, CART et mlpML ou avNNET, entretanto, podem beneficiar a capacidade de previsão se treinados

em *ensemble* de empilhamento.

O escolhido: XGBoost

Como foi visto acima, o algoritmo XGBoost teve o melhor desempenho no estudo comparativo. O *eXtreme Gradient Boosting*, proposto por Chen e Guestrin (2016), é uma evolução recente das técnicas desenvolvidas nos algoritmos *AdaBoost* e *GBM - Stochastic Gradient Boosting Machines* (FREUND; SCHAPIRE, 1997, 1999), e na *Maquina de Boosting Gradual* (FRIEDMAN, 2001; FRIEDMAN, J. H.; HASTIE; TIBSHIRANI, Robert, 1998; FRIEDMAN, 2002). O XGBoost melhora estas técnicas através de otimizações e melhorias em termos de paralelização, cortes de baixo para cima (*pruning*), uso do cache, regularização (*shrinkage* e *column subsampling*), e também a gestão de não-respostas.

O XGBoost é um algoritmo do tipo *ensemble* que combina diversas árvores de decisão e usa a técnica do *boosting* gradiente para otimização. A ideia principal dos *ensembles* é combinar diversos algoritmos fracos e simplificados. No caso do XGBoost são combinadas diversas árvores de decisão tipo CARTs (*Classification and Regression Trees*) usando a técnica de bagging (*Bootstrap AGGregatING*). Esta técnica consiste em incorporar a contribuição de cada uma das árvores por meio de um processo de votação.

O XGBoost difere do Floresta Aleatória, que também é uma combinação de árvores de decisão. Em uma Floresta Aleatória, apenas uma parte das variáveis explicativas é sorteada de maneira aleatória para cada árvore (*column subsampling*) e são computadas médias para formar o *ensemble*. A ideia principal do *boosting* é adicionar novos modelos (*i.e.* novas árvores) ao *ensemble* de forma sequencial. O aprendizado envolve ajustar/otimizar consecutivamente cada modelo de forma a encontrar estimativas mais acuradas da variável resposta. Cada nova árvores $N + 1$ do XGBoost foca no que não foi aprendido na árvore N (no Florestas Aleatórias elas são independentes).

O termo gradiente, por sua vez, vem do fato que cada nova árvore treinada é otimizada em função da precedente usando a técnica do gradiente descendente. Em termos de aprendizado isso consiste em avaliar o erro do modelo atual (testado em uma sub-amostra de validação) e modificá-lo seguindo a maior inclinação descendente da curva de erro de modo a minimizá-la¹. Esta técnica permite escapar dos mínimos locais que interromperiam

¹Pesquisadores descrevem diversos fenômenos naturais de otimização: um espermatozóide segue a direção

o aprendizado (NATEKIN; KNOLL, 2013). O hiperparâmetro *encolhimento* garante que a influência de cada árvore seja equilibrada em relação às futuras/próximas, para que todas possam continuar contribuindo para a otimização do modelo final.

Como outros algoritmos de *gradient boosting*, o XGBoost usa a função *log loss* como função custo para aprender a relação entre as variáveis explicativa x_n e resposta y . Esta função mede o desempenho do modelo em termos de classificação, sendo o resultado uma probabilidade entre 0 e 1. Um modelo perfeito teria um *log loss* igual a 0. O custo aumenta quando são atribuídas probabilidades baixas à uma observação que tem a classe verdadeira. O *log loss* penaliza tanto FP quanto FN, mas em especial as predições que tem alta probabilidade e são erradas².

de maior inclinação em termos de concentração dos marcadores químicos do óvulo; um mosquito segue a direção de maior inclinação da curva de distribuição luminosa para chegar até a fonte de emissão da luz. No caso do AM se trata do campo de parâmetros do modelo sendo treinado. Exemplos Boullier e El Mhamdi (2020)

²A *log loss* também é conhecida como entropia cruzada. A cada divisão das árvores (cada nova decisão) implica um ganho informacional na medida que reduz a entropia dos dados. Se trata assim de encontrar a combinação de variáveis que retorna o maior ganho de informação, tornando os ramos da árvore mais homogêneos e com menor entropia.

Tabela 9: Algoritmos testados no benchmark. A tabela contém a descrição de todos os algoritmos testados no experimento de benchmarking de acordo com a classificação proposta por Fernandez-Delgado, 2014 (coluna 1). As colunas 3 e 4 listam as bibliotecas que foram usadas e os hiperparâmetros de cada algoritmo. A coluna Autores se refere tanto aos autores das bibliotecas quanto das teorias subjacentes a cada um deles.

| Família | Algoritmo | Biblioteca R | Hiperparâmetros | Autores(as) |
|-------------------------|--|--------------|---|---|
| Análises Discriminantes | Regularized Discriminant Analysis (rda) | klaR | gamma, lambda | Jerome Friedman, Hastie e Rob Tibshirani (2010) |
| Abordagens Bayesianas | Bayesian Additive Regression Trees (bart) | bartMachine | num_trees, k, alpha, beta, nu | Chipman, George e McCulloch (2010a,b, 1998), Denison, Mallick e Smith (1998) |
| | Naive Bayes (nb) | klaR | fL, usekernel, adjust | Minsky (1961) |
| | Bayesian Generalized Linear Model (bglm) | arm | | Gelman, Jakulin et al. (2008) |
| Redes Neurais | Neural Network (nnet) | nnet | size, decay | Ackley, Hinton e Sejnowski (1985), Hopfield (1982) |
| | Multi-Layer Perceptron, with multiple layers (mlpML) | RSNNS | layer1, layer2, layer2 | Rosenblatt (1958) |
| | Model Averaged Neural Network (avNNET) | nnet | size, decay, bag | Ripley e Hjort (1995) |
| | Extreme Learning Machines (elmNN) | elmNN | nhid, actfun | Huang et al. (2012) |
| Árvores de decisão | Classification and Regression Tree - CART (rpart) | Rpart | cp | Breiman, Friedman et al. (1984), Therneau, Atkinson et al. (2015) |
| | C5.0 (c50) | C50 | trials, model, winnow | Quinlan (2014), Quinlan (1986) |
| | Conditional Inference Trees (ctree2) | party | maxdepth, mincriterion | Hothorn, Hornik e Zeileis (2006) |
| Ensemble Boosting | Stochastic Gradient Boosting (gbm) | gbm | n.trees, interaction.depth, shrinkage, n.minobsinnode | Freund e Schapire (1997), Freund e Schapire (1999), Friedman (2001), Jerome H. Friedman, Hastie e Robert Tibshirani (1998), Friedman (2002) |
| | eXtreme Gradient Boosting with decision Trees (xgb) | xgboost | nrounds, max_depth, eta, gamma, colsample_bytree, min_child_weight, subsample | Friedman (2001), Jerome H. Friedman, Hastie e Robert Tibshirani (1998), Chen e He (2019) |

Tabela 9: Algoritmos testados no benchmark. A tabela contém a descrição de todos os algoritmos testados no experimento de benchmarking de acordo com a classificação proposta por Fernandez-Delgado, 2014 (coluna 1). As colunas 3 e 4 listam as bibliotecas que foram usadas e os hiperparâmetros de cada algoritmo. A coluna Autores se refere tanto aos autores das bibliotecas quanto das teorias subjacentes a cada um deles. continuação

| Família | Algoritmo | Biblioteca R | Hiperparâmetros | Autores(as) |
|--------------------------------|---|--------------|---|--|
| | eXtreme Gradient Boosting with Dropout regulation (xgbDART) | xgboost | nrounds, max_depth, eta, gamma, colsample_bytree, rate_drop, , skip_drop, min_child_weight, subsample | Rashmi e Gilad-Bachrach (2015), Srivastava et al. (2014) |
| Ensemble Florestas Aleatórias | Random Forest (rf) | randomForest | mtry | Breiman (2001a) |
| Modelos lineares generalizados | GLM with elastic net (glmnet) | glmnet | alpha, lambda | Jerome Friedman, Hastie e Rob Tibshirani (2010) |
| | Regularized Logistic Regression (Log) | LiblineaR | cost, loss, epsilon | Paul (2017) |
| Baseado em Regras | Repeated Incremental Pruning to Produce Error Reduction - RIPPER (JRip) | RWeka | NumOpt, NumFolds, MinWeights | Cohen (1995) |
| | One R (OneR) | RWeka | | Holte (1993) |

^a Classificação de acordo com Fernández-Delgado et al. (2014)

^b Treinados usando os casos completos da Base de Treinamento

D. APÊNDICE D - ARTIGOS EMPÍRICOS SOBRE PATRIOTISMO

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|-----------------------------|---|------|--------|--|
| Kosterman e Feshbach (1989) | Survey com alunos de ensino médio e superior, e membros de um sindicato de empresas de construção | EUA | Não | O estudo mostrou o caráter afetivo do patriotismo, e de superioridade associado ao nacionalismo. Apesar de correlacionadas, atitudes nacionalistas se mostraram mais ligadas à opiniões positivas sobre políticas em armamentos nucleares. A conclusão dos autores é que atitudes patrióticas, ou positivas *intra*grupos, não significam necessariamente hostilidade *extra*grupos. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|---------------------------------|--|--|--------|--|
| Smith e Jarkko (1998) | Survey ISSP 1995 | 23 países | Não | Estudou os fatores que contribuem para a formação do orgulho nacional, entendido como um pré-requisito para tanto patriotismo, quanto nacionalismo. O estudo mostrou que o orgulho nacional era menor em países onde regimes democráticos foram mais longevos e estáveis, nos países da antiga União Soviética, em países onde persistem conflitos étnicos, e em países onde havia uma "culpa de guerra" arraigada na população (Japão e Alemanha). O estudo também mostrou que os mais velhos eram em geral mais orgulhosos devido ao efeito da globalização e multilateralismo nas gerações mais novas, em reação ao nacionalismo agressivo da Segunda Guerra Mundial. No mesmo ano, usando também os dados da Pesquisa Mundial de Valores para a Austrália, Alemanha, Inglaterra e Suécia, @hjern_national_1998 não encontrou correlação entre identidade cívica e orgulho nacional e xenofobia, que teria uma relação maior com um tipo de identidade étnica-nacional. |
| Hjern (1998) | ISSP 1995; WVS Wave 6 (2010-2012) | Austrália, Alemanha, Inglaterra e Suécia | Não | Não encontrou correlação entre identidade cívica e xenofobia, que teria uma relação maior com um tipo de identidade étnica-nacional. |
| Schatz, Staub e Lavine (1999) | 291 estudantes de graduação da University of Massachusetts–Amherst | EUA | Não | O autor encontrou evidências de que o patriotismo cedo estaria associado ao nacionalismo, desengajamento político, percepção de ameaças externas e da importância de comportamentos simbólicos, além de um viés de exposição mediática pró-nacional. O segundo, patriotismo construtivo por sua vez estaria associado a um tipo de virtude política. |
| Mummendey, Klink e Brown (2001) | 381 estudantes de instituições da Alemanha e Inglaterra | Inglaterra e Alemanha | Não | Encontrou uma correlação entre a hostilidade a grupos externos e a construção da identidade nacional com base em comparações intergrupos, e não intragrupo (no tempo por exemplo). O primeiro teria base comparações autônomas e a segunda relacionais, o que diferenciaria orientações patrióticas e nacionalistas. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|-------------------------------|--|-----------|--------|--|
| Jones e Smith (2001) | ISSM 1995 | 23 países | Não | Explorou as origens do pertencimento nacional e identificou formas cívicas/voluntaristas, e prescritas/étnicas de nacionalismo. No estudo foi identificada uma intersecção entre o desempenho do país e a identidade, mas que não se sobrepõe totalmente aos fatores objetivos e pessoais na identificação nacional (local de nascimento, religião). Apesar dos efeitos das migrações, globalização os autores encontraram a identificação nacional com sólidas bases tradicionais. |
| Evans e Kelley (2002) | ISSP 1995-1996 | 24 países | Não | Pesquisou a influência de fatores nacionais como a ciência, economia, artes e literatura, e esporte no orgulho nacional. Contrariamente as hipóteses da globalização, o autor encontrou correlações com o desempenho dos países nestas áreas. Também encontrou o efeito da globalização na correlação positiva entre a idade e orgulho, além de semelhanças e diferenças entre as fontes de orgulho em diferentes sociedades (*e.g.* ingleses se orgulham mais da ciência do que das artes, orgulho do desempenho econômico em todas as sociedades, orgulho no esporte em nações menores). |
| De Figueiredo e Elkins (2003) | ISSP 1995; WVS 1981, 1990-91, e 1995-97; GSS 1994 e 1996 | | Sim | Considera diferentes formas de orgulho pelo país correlacionadas em formas de patriotismo e nacionalismo. Usando dados combinados de diversas surveys realizadas entre 1990 e 1981 encontrou diferenças entre atitudes nacionalistas e patrióticas, e uma correlação entre orgulho e preconceito em relação a membros de grupos externos. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|---------------------------|--|----------------|--------|--|
| Smith e Kim (2006) | ISSP 1995/96 (Alemanhas separadas) e 2003/04 33 países (incluindo as duas Alemanhas) | 24 países | Não | Estudou as mudanças entre os padrões de orgulho nacional entre os anos 1995-96 e 2003/04 também usando o ISSP em 24 países. O estudo mostra uma estabilidade e variação entre países e grupos socio-econômicos. Pouca diferença entre gêneros, com os homens ligeiramente mais orgulhosos, e uma diferença importante em faixas etárias, sendo os mais velhos mais orgulhosos. No estudo, o orgulho nacional geral está mais relacionado com um pertencimento que o autor chama de "verdadeiro": oposição ao internacionalismo e globalização, e visão negativa da imigração e imigrantes. Este tipo de orgulho estaria ligado à culturas majoritárias de um país e coortes que foram socializadas em períodos menos globalizados. |
| Tilley e Heath (2007) | WVS 1981, 1990 e Eurobarometer 1982, 1983, 1984, 1985, 1986, 1988, 1994 e 1997 | Inglaterra | Não | Estudou o declínio do orgulho nacional na Inglaterra, e mostrou que o efeito era essencialmente geracional. |
| Duchesne e Frogner (2007) | Eurobarometer | União Européia | Não | Se perguntarm porque é tão difícil entender se o orgulho nacional impede ou facilita a identidade Européia. Usando dados do Eurobarômetro dos anos 1994 até 2000, os autores mostram uma correlação entre as duas identidades, denotando uma evolução paralela e complementar do sentimento de identificação. |
| Muñoz (2009) | WVS/EVS 1981, 1990, 1995, 1999 e 2000 | Espanha | Não | Acompanhou a mudança no orgulho nacional para estudar as mudanças nas base sociais de apoio à nação na transição do Franquismo à democracia |
| Davidov (2009/ed) | ISSP 2003 | 34 países | Não | identificou o patriotismo construtivo como distinto do nacionalismo. O patriotismo construtivo inclui o orgulho na democracia, seguridade social e igualdade no tratamento dos cidadãos, enquanto que nacionalismo se relaciona com declarações relacionadas à superioridade cultural. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|---------------------|---|--|--------|---|
| Shayo (2009) | ISSP 1995; WVS: Wave 3 (1995-1998); Wave 2 (1990-1994); Wave 1 (1981-1984) (LIS). | 20 países | Sim | Sugere que a identidade nacional na forma de orgulho do país é mais comum entre os indivíduos mais pobres, e tende a reduzir o apoio à políticas redistributivas. O trabalho mostra ainda que a prevalência em nível societal de orgulho tem correlação negativa com o nível de redistribuição. |
| Bonikowski (2010) | ISSP 1995 e 2003 | EUA | Não | Identificou quatro tipos de nacionalismos, um deles centrado no orgulho do estado, que muda de acordo com acontecimentos políticos e econômicos que exacerbam a saliência do estado-nação. O autor também identificou diferenças importantes em termos de nacionalismo entre classes mais educadas e jovens, e mais velhas e menos educadas. |
| Green et al. (2011) | ISSP 2003 | Suíça | Não | Usando dados de 2003 para a Suíça encontra uma correlação entre patriotismo, na forma de orgulho pelas instituições democráticas, e atitudes mais tolerantes. |
| Ariely (2011) | ISSP National Identity 2003 | 15 democracias ocidentais (Áustria, Inglaterra, Canadá, Dinamarca, Finlândia, Alemanha, Irlanda, França, Países Baixos, Noruega, Portugal, Espanha, Suécia, Suíça e EUA) | Não | Identificou elementos culturais (história, arte, literatura) e políticos (seguridade social, tratamento justo e igualitário, funcionamento da democracia) de forma a identificar o patriotismo constitucional e o nacionalismo liberal em 15 democracias ocidentais. Ariely encontrou evidências que nacionalismo de bases culturais tende a contrapor valores liberais, e que pelo contrário o patriotismo cultural pode atingir os objetivos da democracia liberal mesmo dentro de grande diversidade cultural. |
| Davidov (2011) | ISSP 1995, 2003 | 22 países | Não | Repetição do mesmo estudo de 2009 |
| Solt (2011) | WVS 1981, 1990, 1995-97 e 1999-2001; Eurobarometer 1986 | 22 países | Sim | Os resultados sugerem que a desigualdade tem efeitos no orgulho nacional pelo efeito de "distanciamento" (*diversion*) que os governantes exercem sobre os cidadãos de forma a criar ou reforçar "mitos" e contrapor demandas sociais. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|---------------------------------------|--|----------------------------|------------------|--|
| Ariely (2012) | ISSP 2003; World Values Survey 2005 | 63 países | Sim | Globalization (modernization and democracy) is related negatively to patriotism (pride), (BR lower than expected), willingness to fight and ethnic conceptions of membership in the nation (BR higher than expected but ancestry is very low); not impact in national identification and sense of nationalism (not measured for BR) |
| Kavetsos (2012) | Eurobarômetro | União Européia (12 países) | Não | Identificou crescimento do orgulho depois de eventos esportivos de grande porte em 12 países União Européia pesquisados pelo Eurobarômetro. |
| Han (2013) | WVS 1995, 2000 e 2005 | | Sim | Os testes que o autor realizou identificaram que pessoas mais pobres são mais afetadas pela desigualdade que aumentam o orgulho, especialmente em países com muitos migrantes e migração. |
| Fabrykant (2014) | WVS Wave 5 (2005-2009), 4 (1999-2004) e 3 (1995-1998) | 82 países | Sim (1997, 2006) | Usou dados do WVS 5 (2005-2009), 4 (1999-2004) e 3 (1995-1998) para pesquisar as diferenças entre nacionalismo e multiculturalismo em 82 países. Segundo ela existem diferentes nacionalismos contemporâneos ligados a conjuntos de valores gerais. Nacionalismo tem mais relação com a tolerância, religião, e felicidade do que aspectos políticos. Também identificou identidades híbridas leais tanto a grupos subnacionais quanto a formas de cosmopolitismo. |
| Eckel (2014) | EVS 4 (2008) e indicadores nacionais (Comparative Manifesto Project) | União Européia | Não | Confirma os resultados de Solt e o papel importante das elites na instrumentalização da desigualdade. |
| Pawlowski, Downward e Rasciute (2014) | ISSP 2007 | 33 países | Não | Ver @sullivan_collective_2014 para o Brasil |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|------------------------------|---|---|--------|--|
| Lan e Li (2015) | Chinese Political Compass; World Value Surveys 2001, 2007 | China + 15 países (Argentina, Canadá, Chile, Índia, Indonésia, Japão, Coréia do Sul, México, Marrocos, África do Sul, Espanha, Suécia, Turquia e USA) | Não | Sugerem que uma atividade econômica voltada para o comércio exterior, desvia a atenção do mercado interno reduzindo o nacionalismo (medido usando orgulho, apoio às forças armadas e confiança em algumas instituições) |
| Ha e Jang (2015) | Korean General Social Survey (KGSS) 2010 | Coréia do Sul | Não | Orgulho é uma forma de satisfação nacional, relacionado com o sentimento de felicidade |
| Ariely (2016) | World Values Survey 2005 e 2010; European Values Study 2008 | 93 países | Sim | Estudou o patriotismo como orgulho nacional, usando o WVS de 2005 e 2010, e o EVS 2008 para 93 países. O autor procurou encontrar os fatores que determinam os níveis de orgulho nacional. O estudo mostrou que países mais desenvolvidos e globalizados tem níveis mais baixos de orgulho, enquanto países mais desiguais, menos homogêneos quanto à religião, ou em conflito apresentam padrão oposto. O autor sugere que estas evidências não indicam patriotismo como uma virtude. |
| Wolak e Dawkins (2016) | American National Election Study (ANES) 2012 | EUA | Não | Os norte-americanos são mais propensos a dizer que sentem amor pelo seu país em contextos políticos congruentes com os seus valores e entendimento da cidadania. |
| Bonikowski e DiMaggio (2016) | GSS 1996, 2004 e 2012 | EUA | Não | |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|--|--------------------------|---|--------|--|
| Dimitrova-Grajzl, Eastwood e Grajzl (2016) | Expert survey | Europa e países da antiga União Soviética | Não | Identidade nacional é um artefato cultural que leva tempo para emergir. |
| Ariely (2018) | ISSP 2013 | 29 países | Não | Patriotismo é geralmente positivo; na média maiorias tem melhor impressão do que minorias; minorias avaliam o governo melhor quando existem políticas inclusivas; quanto mais inclusivas, mais as maiorias percebem o patriotismo |
| Qian e Hung (2018) | WVS 1990 à 2014 | 20 economias (Argentina, Austrália, Brasil, Chile, China, Alemanha, Índia, Japão, México, Nigéria, Peru, Polônia, Rússia, África do Sul, Coréia do Sul, Spain, Sweden, Taiwan, Turquia e EUA) | Sim | Buscou testar o efeito da desigualdade sobre o nacionalismo. Usando o WVS de 1990 a 2014 para 20 economias não encontrou uma relação neutra. A sua hipótese é que os efeitos identificados de forma empírica por Shayo são anulados pelo efeito descrito por Solt. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|-------------------------|--|--|--------|---|
| Noh (2018) | Barômetro das Américas 2016/17 | 11 países da América do Sul e Central (México, Guatemala, El Salvador, Honduras, Nicarágua, Costa Rica, Panamá, Equador, Bolívia, Peru, Paraguai, Chile, Uruguai, Brasil, Venezuela, República Dominicana, Haiti, Jamaica) | Sim | Brasil é o mais baixo em termos de muito orgulho entre países latino-americanos (61%), somente EUA e Canada ficam abaixo. Mais velhos tendem a demonstrar mais orgulho, enquanto sexo, local de residência, educação e riqueza não tem influência. Itens salientes no país influenciam as respostas. No Brasil especificamente, a situação econômica, o crime a corrupção e insegurança levam a menores respostas positivas em termos de orgulho. Desemprego não. |
| Hofstede et al. (2010) | Values Survey Module em estudos de mercado entre 2000 e 2001 | Brasil | Sim | Jeitinho como regra e forma de evitar incerteza; Individualism/collectivismo: neither modern nor traditional but both at the same time; |
| Satherley et al. (2019) | | Nova Zelândia | Não | Encontrou também diferentes perfis de patriotismo, e que nacionalismo e patriotismo geralmente ocorrem simultaneamente. Além disso, os resultados apontaram que atitudes puramente patrióticas estão relacionadas com alta tolerância a migrantes, e orientações relativamente liberais. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|----------------------------|--|----------------|--------|--|
| Leite et al. (2018) | Estudantes Universitários | Brasil | Sim | Patriotismo e nacionalismo são mediados pelo essencialismo, entendido como a crença em uma personalidade compartilhada pelos membros de um grupo que os diferencia de outros grupos, gerando relações sociais baseadas em categorias |
| Qari, Konrad e Geys (2012) | ISSP 2003 e dados econômicos da OCDE | 22 países | Não | Encontrou pelo contrário uma correlação positiva entre imposição progressiva e redistributiva e atitudes patrióticas, e uma instrumentalização dos sentimentos patrióticos para políticas fiscais. |
| Nincic e Ramos (2012) | PEW 1987-2007 e experimentos | Estados Unidos | Não | Usando dados da PEW e experimentos estudaram os efeitos da política externa no patriotismo, seja ele absoluto ou crítico, "contingente" a aprovação das políticas do país. Os resultados indicam a formação do patriotismo crítico se dá mais por fatores endógenos, mas que ambos exercem influência no patriotismo genérico. |
| Blank e Schmidt (2003) | National identity of Germans, painél realizado em 1996 | Alemanha | Não | Patriotismo tolerante e com uma certa distância crítica do estado e do regime político. |
| Green (2020) | Barômetro da África | Uganda | Não | Representação política do grupo étnico influencia as declarações de orgulho |
| Wimmer (2017) | WVS 1990 à 2014 | 123 países | Yes | Orgulho é produzido pelo poder. Membros de grupos não representados são menos orgulhosos. Em nível de país, quanto maior a parte não-representada, menor o orgulho nacional. Pré-existência de conflitos étnicos também reduz orgulho. |

Tabela 10: Artigos empíricos sobre patriotismo e orgulho nacional. A tabela contém os autores, as bases de dados usadas, se o Brasil faz parte ou não da análise e as principais conclusões dos autores (*continued*)

| Artigo | Técnicas e base de dados | País | Brasil | Principais Conclusões |
|--|--------------------------|------|--------|-----------------------|
| ^a ISSP: International Social Survey Programme | | | | |
| ^b NIS: National Identity Study | | | | |
| ^c GSS: General Social Survey | | | | |
| ^d EVS: European Values Survey | | | | |
| ^e WVS: World Values Survey | | | | |
| ^f LIS: Luxembourg Income Study | | | | |
| ^g ANES: American National Election Study | | | | |
| ^h KGSS: Korean General Social Survey | | | | |
| ⁱ CPC: Chinese Political Compass | | | | |
| ^j VSM: Values Survey Module | | | | |

BIBLIOGRAFIA

ABRAMSON, Paul R.; INGLEHART, Ronald. **Value Change in Global Perspective**. Ann Arbor, Mich: University of Michigan Press, 1995. ISBN 978-0-472-09591-9 978-0-472-06591-2.

ACKLEY, David H.; HINTON, Geoffrey E.; SEJNOWSKI, Terrence J. A Learning Algorithm for Boltzmann Machines. en. **Cognitive Science**, v. 9, n. 1, p. 147–169, 1985. ISSN 1551-6709. DOI: 10.1207/s15516709cog0901_7.

ADADI, Amina; BERRADA, Mohammed. Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence (XAI). **IEEE Access**, v. 6, p. 52138–52160, 2018. ISSN 2169-3536. DOI: 10.1109/ACCESS.2018.2870052.

ADAMCZYK, Amy; PITT, Cassidy. Shaping Attitudes about Homosexuality: The Role of Religion and Cultural Context. **Social Science Research**, v. 38, n. 2, p. 338–351, jun. 2009. ISSN 0049-089X. DOI: 10.1016/j.ssresearch.2009.01.002.

ADORNO, T. W. et al. **The Authoritarian Personality**. First Edition edition. [S.l.]: Harper & Brothers, 1950. ISBN 978-0-06-030150-7.

AGRESTI, Alan. **An Introduction to Categorical Data Analysis**. Hoboken, NJ, USA: John Wiley & Sons, Inc., mar. 2007. ISBN 978-0-470-11475-9 978-0-471-22618-5. DOI: 10.1002/0470114754.

ALENCASTRO, Luiz Felipe de. **O trato dos viventes: formação do Brasil no Atlântico Sul, séculos XVI e XVII**. [S.l.]: Companhia das Letras, 2000. ISBN 978-85-359-0008-8.

ALLAIRE, JJ et al. **Rmarkdown: Dynamic Documents for r**. [S.l.: s.n.], 2020.

ALLISON, Paul David. **Missing Data**. Thousand Oaks, Calif: Sage Publications, 2002. (Sage University Papers. Quantitative Applications in the Social Sciences, no. 07-136). ISBN 978-0-7619-1672-7.

ALMOND, Gabriel; VERBA, Sidney. **Civic Culture Study, 1959-1960: Version 2**. [S.I.]: Inter-university Consortium for Political and Social Research, jun. 1984. DOI: 10.3886/ICPSR07201.v2.

ALMOND, Gabriel A. Comparative Political Systems. **The Journal of Politics**, v. 18, n. 3, p. 391–409, ago. 1956.

_____. The Intellectual History of the Civic Culture. In: ALMOND, Gabriel A.; VERBA, Sidney (Ed.). **The Civic Culture Revisited**. Revised edition (1 mai 1989). [S.I.]: SAGE Publications Inc, 1989. p. 432. ISBN 0-8039-3560-9.

ALMOND, Gabriel A.; VERBA, Sidney. **The Civic Culture: Political Attitudes and Democracy in Five Nations**. Newbury Park, Calif: SAGE Publications Inc, 1963. ISBN 978-0-8039-3558-7.

ALPAYDIN, Ethem. **Introduction to Machine Learning**. Second. [S.I.]: The MIT Press, 2010. ISBN 978-0-262-01243-0.

_____. **Machine Learning. The New AI**. [S.I.]: The MIT Press, 2017.

ALVAREZ, R. Michael (Ed.). **Computational Social Science: Discovery and Prediction**. Reprint edition. New York, NY: Cambridge University Press, mar. 2016. ISBN 978-1-107-51841-4.

ALVAREZ-MELIS, David; JAAKKOLA, Tommi S. On the Robustness of Interpretability Methods. **arXiv:1806.08049 [cs, stat]**, jun. 2018. arXiv: 1806.08049 [cs, stat].

ANDERSON, Benedict R. O'G. Benedict Anderson e as Fronteiras (e Anomalias) Do Nacionalismo. **Jornal da Unicamp**, set. 2011.

_____. **Imagined Communities: Reflections on the Origin and Spread of Nationalism**. Rev. and extended ed. London ; New York: Verso, 1991. ISBN 978-0-86091-329-0 978-0-86091-546-1.

ANDRIDGE, Rebecca R.; LITTLE, Roderick J. A. A Review of Hot Deck Imputation for Survey Non-Response. **International statistical review = Revue internationale de statistique**, v. 78, n. 1, p. 40–64, abr. 2010. ISSN 0306-7734. DOI: 10.1111/j.1751-5823.2010.00103.x.

APLEY, Daniel W.; ZHU, Jingyu. Visualizing the Effects of Predictor Variables in Black Box Supervised Learning Models. **arXiv:1612.08468 [stat]**, dez. 2016. arXiv: 1612.08468 [stat].

APPIAH, Kwame Anthony. Patriotas Cosmopolitas. Tradução: Antonio Sérgio Alfredo Guimarães. **Revista Brasileira de Ciências Sociais**, v. 13, n. 36, 1998.

ARIELY, Gal. Constitutional Patriotism, Liberal Nationalism and Membership in the Nation: An Empirical Assessment. en. **Acta Politica**, v. 46, n. 3, p. 294–319, jul. 2011. ISSN 1741-1416. DOI: 10.1057/ap.2011.10.

_____. Evaluations of Patriotism across Countries, Groups, and Policy Domains. en. **Journal of Ethnic and Migration Studies**, v. 44, n. 3, p. 462–481, fev. 2018. ISSN 1369-183X, 1469-9451. DOI: 10.1080/1369183X.2017.1319761.

_____. Globalisation and the Decline of National Identity? An Exploration across Sixty-Three Countries: Globalisation and the Decline of National Identity. en. **Nations and Nationalism**, v. 18, n. 3, p. 461–482, jul. 2012. ISSN 13545078. DOI: 10.1111/j.1469-8129.2011.00532.x.

_____. The Nexus between Globalization and Ethnic Identity: A View from Below. **Ethnicities**, v. 0, n. 0, p. 1–22, 2019. DOI: 10.1177/1468796819834951.

_____. Why Does Patriotism Prevail? Contextual Explanations of Patriotism across Countries. en. **Identities**, p. 1–27, mar. 2016. ISSN 1070-289X, 1547-3384. DOI: 10.1080/1070289X.2016.1149069.

ATHEY, Susan. Beyond Prediction: Using Big Data for Policy Problems. en. **Science**, v. 355, n. 6324, p. 483–485, fev. 2017. ISSN 0036-8075, 1095-9203. DOI: 10.1126/science.aal4321.

BADIE, Bertrand. **Culture et politique**. [S.l.]: Economica, jan. 1993. ISBN 978-2-7178-2426-1.

BAEHRENS, David et al. How to Explain Individual Classification Decisions. en. **Journal of Machine Learning Research**, v. 11, p. 1803–1831, 2010.

- BAINBRIDGE, William Sims et al. Artificial Social Intelligence. **Annual Review of Sociology**, v. 20, n. 1, p. 407–436, 1994. DOI: 10.1146/annurev.so.20.080194.002203.
- BAQUERO, Cesar Marcello Jacome; GONZALEZ, Rodrigo Stumpf. Political Culture, Economic Changes, And Inertial Democracy: A Post-2014 Elections Analysis. Portuguese. **Scopus**, Universidade Estadual de Campinas, 2016. ISSN 0104-6276. DOI: 10.1590/1807-01912016223492.
- BAQUERO, Marcello. **Democracia Inercial: Assimetrias Entre Economia e Cultura Política Na América Latina**. Porto Alegre: UFRGS Editora, 2018. ISBN 978-85-386-0403-7.
- BAQUERO, Marcello; CASTRO, Henrique C. O. de; RANINCHESKI, Sônia M. (Des) Confiança Nas Instituições e Partidos Políticos Na Constituição de Uma Democracia Inercial No Brasil: O Caso Das Eleições de 2014. **Política & Sociedade**, v. 15, n. 32, p. 9–38, 2016.
- BAR-TAL, Daniel; STAUB, Ervin. **Patriotism in the Lives of Individuals and Nations**. [S.l.]: Nelson-Hall Publishers, 1997. ISBN 978-0-8304-1410-9.
- BARRY, Brian M. **Sociologists, Economists and Democracy**. Reprint. Chicago: Univ. of Chicago Press, 1988. ISBN 978-0-226-03824-7.
- BATISTA, Gustavo E. A. P. A.; MONARD, Maria Carolina. An Analysis of Four Missing Data Treatment Methods for Supervised Learning. en. **Applied Artificial Intelligence**, v. 17, n. 5-6, p. 519–533, mai. 2003. ISSN 0883-9514, 1087-6545. DOI: 10.1080/713827181.
- BAUMEISTER, Roy F. Self-Esteem, Self-Presentation, and Future Interaction: A Dilemma of Reputation. en. **Journal of Personality**, v. 50, n. 1, p. 29–45, 1982. ISSN 1467-6494. DOI: 10.1111/j.1467-6494.1982.tb00743.x.
- BBC. Japan Schools to Teach Patriotism. en-GB. **BBC News**, mai. 2007.
- BEAUCHAMP, Nicholas. Predicting and Interpolating State-Level Polls Using Twitter Textual Data: PREDICTING POLLS WITH TWITTER. en. **American Journal of Political Science**, v. 61, n. 2, p. 490–503, abr. 2017. ISSN 00925853. DOI: 10.1111/ajps.12274.
- BENGIO, Yoshua; COURVILLE, Aaron; VINCENT, Pascal. Representation Learning: A Review and New Perspectives. **arXiv:1206.5538 [cs]**, jun. 2012. arXiv: 1206.5538 [cs].

BERGMAN, Manfred Max. A Theoretical Note on the Differences Between Attitudes, Opinions, and Values. en. **Swiss Political Science Review**, v. 4, n. 2, p. 81–93, jun. 1998. ISSN 14247755. DOI: 10.1002/j.1662-6370.1998.tb00239.x.

BERGMEIR, Christoph et al. **RSNNS: Neural Networks Using the Stuttgart Neural Network Simulator (SNNS)**. [S.l.: s.n.], set. 2019.

BERGSTRA, James; BENGIO, Yoshua. Random Search for Hyper-Parameter Optimization. en. **Journal of Machine Learning Research**, v. 13, p. 281–305, 2012.

BERGSTRA, James S.; BARDENET, Rémi et al. Algorithms for Hyper-Parameter Optimization. In: SHAWE-TAYLOR, J. et al. (Ed.). **Advances in Neural Information Processing Systems 24**. [S.l.]: Curran Associates, Inc., 2011. p. 2546–2554.

BERK, Richard. **Criminal Justice Forecasts of Risk: A Machine Learning Approach**. New York: Springer-Verlag, 2012. (SpringerBriefs in Computer Science). ISBN 978-1-4614-3084-1.

BERNARDES, Denis. **O patriotismo constitucional: Pernambuco, 1820-1822**. [S.l.]: Editora Universitária UFPE, 2006. ISBN 978-85-60438-00-6.

BIAU, Gérard; SCORNET, Erwan. A Random Forest Guided Tour. en. **TEST**, v. 25, n. 2, p. 197–227, jun. 2016. ISSN 1863-8260. DOI: 10.1007/s11749-016-0481-7.

BILL, James A.; HARDGRAVE, Robert L. **Comparative Politics: The Quest for Theory**. Washington, D.C: University Press Of America, 1973. ISBN 978-0-8191-2090-8.

BILLIG, Michael. **Banal Nationalism**. [S.l.]: SAGE Publications, set. 1995. ISBN 978-0-8039-7525-5.

BLANK, Thomas; SCHMIDT, Peter. National Identity in a United Germany: Nationalism or Patriotism? An Empirical Test with Representative Data. **Political Psychology**, v. 24, n. 2, p. 289–312, 2003. ISSN 0162-895X. DOI: 10.1111/0162-895X.00329.

BLITZER, John; DREDZE, Mark; PEREIRA, Fernando. Biographies, Bollywood, Boom-Boxes and Blenders: Domain Adaptation for Sentiment Classification. en. **Association of Computational Linguistics**, 2007.

BOBBIO, Norberto. Política. In: MATTEUCCI, Nicola; PASQUINO, Gianfranco; BOBBIO, Norberto (Ed.). **Dicionário de Política**. Brasília, DF: Ed. Univ. de Brasília, 1993. ISBN 978-85-230-0308-1.

- BOELAERT, Julien; OLLION, Etienne. The Great Regression. Machine Learning, Econometrics, and the Future of Quantitative Social Sciences. **Revue française de sociologie**, 2018.
- BOLTON, Richard J.; HAND, David J. Statistical Fraud Detection: A Review. en. **Statistical Science**, v. 17, n. 3, p. 235–255, ago. 2002. ISSN 0883-4237, 2168-8745. DOI: 10.1214/ss/1042727940.
- BONIKOWSKI, Bart. **Varieties of Popular Nationalism in Modern Democracies: An Inductive Approach to Comparative Research on Political Culture**. Cambridge Mass., 2010.
- BONIKOWSKI, Bart; DIMAGGIO, Paul. Varieties of American Popular Nationalism. **American Sociological Review**, p. 0003122416663683, 2016.
- BORBA, Julian. Cultura Política, Ideologia e Comportamento Eleitoral: Alguns Apontamentos Teóricos Sobre o Caso Brasileiro. **Opinião Pública**, v. 11, n. 1, p. 147–168, mar. 2005. ISSN 0104-6276. DOI: 10.1590/S0104-62762005000100006.
- BOULLIER, Dominique; EL MHAMDI, El Mahdi. Le machine learning et les sciences sociales à l'épreuve des échelles de complexité algorithmique. fr. **Revue d'anthropologie des connaissances**, Société d'Anthropologie des Connaissances, v. 14, n. 1, mar. 2020. ISSN 1760-5393.
- BREHM, John O. **The Phantom Respondents: Opinion Surveys and Political Representation**. [S.l.]: University of Michigan Press, 1993. (Michigan Studies in Political Analysis). ISBN 978-0-472-09523-0.
- BREIMAN, Leo. Bagging Predictors. en. **Machine Learning**, v. 24, n. 2, p. 123–140, ago. 1996. ISSN 1573-0565. DOI: 10.1023/A:1018054314350.
- _____. Random Forests. en. **Machine Learning**, v. 45, n. 1, p. 5–32, out. 2001. ISSN 1573-0565. DOI: 10.1023/A:1010933404324.
- _____. Statistical Modeling: The Two Cultures. en. **Statistical Science**, v. 16, n. 3, p. 199–215, 2001.
- BREIMAN, Leo; CUTLER, Adele et al. **randomForest: Breiman and Cutler's Random Forests for Classification and Regression**. [S.l.: s.n.], 2018.

- BREIMAN, Leo; FRIEDMAN, Jerome et al. **Classification and Regression Trees**. [S.l.]: Taylor & Francis, jan. 1984. ISBN 978-0-412-04841-8.
- BRENT, Edward. Designing Social Science Research with Expert Systems. **Anthropological Quarterly**, The George Washington University Institute for Ethnographic Research, v. 62, n. 3, p. 121–130, 1989. ISSN 0003-5491. DOI: 10.2307/3317452.
- BRINGSJORD, Selmer; GOVINDARAJULU, Naveen Sundar. Artificial Intelligence. In: ZALTA, Edward N. (Ed.). **The Stanford Encyclopedia of Philosophy**. Fall 2018. [S.l.]: Metaphysics Research Lab, Stanford University, 2018.
- BRUBAKER, Rogers. In the Name of the Nation: Reflections on Nationalism and Patriotism. en. **Citizenship Studies**, v. 8, n. 2, p. 115–127, jun. 2004. ISSN 1362-1025, 1469-3593. DOI: 10.1080/1362102042000214705.
- BUIS, Maarten L. Logistic Regression: When Can We Do What We Think We Can Do? In: UNITED Kingdom Stata Users' Group Meetings 2015. United Kingdom: Stata Users Group, 2015. v. 8, p. 19.
- BURNS, E. B. **Nationalism in Brazil A Historical Survey**. [S.l.]: Frederick A. Praeger, 1968.
- CALHOUN, Craig; WIEVIORKA, Michel. Manifeste pour les sciences sociales. fr. **Socio. La nouvelle revue des sciences sociales**, n. 1, p. 5–39, mar. 2013. ISSN 2266-3134. DOI: 10.4000/socio.200.
- CÂNDIDO, Antônio. Literature and the Rise of Brazilian National Self-Identity. **Luso-Brazilian Review**, v. 5, n. 1, p. 27–43, 1968. ISSN 0024-7413.
- CANOVAN, Margaret. Patriotism Is Not Enough. en. **British Journal of Political Science**, v. 30, n. 3, p. 413–432, jul. 2000. ISSN 1469-2112, 0007-1234. DOI: 10.1017/S000712340000017X.
- CANTU, Francisco; SAIEGH, Sebastian M. A Supervised Machine Learning Procedure to Detect Electoral Fraud Using Digital Analysis. en. **SSRN Electronic Journal**, 2010. ISSN 1556-5068. DOI: 10.2139/ssrn.1594406.
- CARVALHO, Diogo V.; PEREIRA, Eduardo M.; CARDOSO, Jaime S. Machine Learning Interpretability: A Survey on Methods and Metrics. en. **Electronics**, Multidisciplinary Digital Publishing Institute, v. 8, n. 8, p. 832, ago. 2019. DOI: 10.3390/electronics8080832.

CARVALHO, José Murilo de. **A Construção Da Ordem e Teatro de Sombras**. [S.l.]: Civilização Brasileira, 2008.

_____. **A formação das almas - O imaginário da República no Brasil**. Edição: 1ª. [S.l.]: Companhia das Letras, 2017. ISBN 978-85-359-2895-2.

_____. A Utopia de Oliveira Viana. **Revista Estudos Históricos**, v. 4, n. 7, p. 82–99, 1991.

_____. Brazil 1870-1914. The Force of Tradition. **Journal of Latin American Studies**, v. 24, p. 145–162, 1992. ISSN 0022-216X. DOI: 10.1017/S0022216X00023816.

_____. **Cidadania No Brasil: O Longo Caminho**. Rio de Janeiro: Civilização Brasileira, 2001. ISBN 978-85-200-0565-1.

_____. Cidadania: tipos e percursos. pt. **Estudos Históricos**, n. 18, 1996.

_____. Dimensiones de la ciudadanía en el Brasil del siglo XIX. In: SÁBATO, Hilda (Ed.). **Ciudadanía política y formación de las naciones: perspectivas históricas de América Latina**. [S.l.]: Colegio de México, jan. 1999. pp. 321–344. ISBN 978-968-16-5147-3.

_____. **Elite and state-building in imperial Brazil**. [S.l.]: Stanford University, 1974.

_____. O Motivo Edênico No Imaginário Social Brasileiro. **Revista Brasileira de Ciências Sociais**, v. 13, n. 38, out. 1998. ISSN 0102-6909. DOI: 10.1590/S0102-69091998000300004.

_____. Political Elites and State Building: The Case of Nineteenth-Century Brazil. **Comparative Studies in Society and History**, v. 24, n. 3, p. 378–399, 1982. ISSN 0010-4175.

CASTRO, Henrique C. de O. de. **Cultura Política Comparada: Democracia e Mudanças Econômicas: Brasil, Argentina e Chile**. Brasília: Verbena, 2014. ISBN 978-85-64857-18-6.

_____. Cultura Política, Democracia e Hegemonia na América Latina. pt. **Revista de Estudos e Pesquisas sobre as Américas**, v. 5, n. 2, p. 79–96, 2011. ISSN 1984-1639.

_____. Cultura Política: A Tentativa de Construção de Um Conceito Adequado à América Latina. **Revista de Estudos e Pesquisas sobre as Américas**, v. 2, n. 1, 2008.

CASTRO, Henrique C. de O. de; CAPISTRANO, Daniel. Cultura Política Pós-Consenso de Washington: O Conceito de Cultura Cívica e a Mudança Política Na América Latina.

Revista Debates, v. 2, n. 1, p. 75, 2008.

CHABLO, Alexander. What Can Artificial Intelligence Do for Anthropology? **Current Anthropology**, v. 37, n. 3, p. 553–555, 1996. ISSN 0011-3204. DOI: 10.1086/204518.

CHANG, Linchiat; KROSNICK, Jon A. National Surveys Via Rdd Telephone Interviewing Versus the Internet Comparing Sample Representativeness and Response Quality. en. **Public Opinion Quarterly**, v. 73, n. 4, p. 641–678, jan. 2009. ISSN 0033-362X. DOI: 10.1093/poq/nfp075.

CHEN, Jiahua; SHAO, Jun. Nearest Neighbor Imputation for Survey Data. en. **Journal of Official Statistics**, v. 16, n. 2, p. 113–131, 2000.

CHEN, Nan-Chen et al. Using Machine Learning to Support Qualitative Coding in Social Science: Shifting The Focus to Ambiguity. en. **ACM Transactions on Interactive Intelligent Systems**, v. 9, n. 4, p. 21, mar. 2018.

CHEN, Tianqi; GUESTRIN, Carlos. XGBoost: A Scalable Tree Boosting System. en. **Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining - KDD '16**, p. 785–794, 2016. DOI: 10.1145/2939672.2939785. arXiv: 1603.02754.

CHEN, Tianqi; HE, Tong. **Xgboost: eXtreme Gradient Boosting**. [S.l.: s.n.], 2019.

CHEN, Tianqi et al. **Xgboost: Extreme Gradient Boosting**. [S.l.: s.n.], 2020.

CHILCOTE, Ronald H. Development and Nationalism in Brazil and Portuguese Africa. en. **Comparative Political Studies**, v. 1, n. 4, p. 501–525, jan. 1969. ISSN 0010-4140, 1552-3829. DOI: 10.1177/001041406900100403.

CHIPMAN, Hugh A.; GEORGE, Edward I.; MCCULLOCH, Robert E. BART: Bayesian Additive Regression Trees. EN. **The Annals of Applied Statistics**, v. 4, n. 1, p. 266–298, mar. 2010. ISSN 1932-6157, 1941-7330. DOI: 10.1214/09-AOAS285.

_____. _____. **The Annals of Applied Statistics**, v. 4, n. 1, p. 266–298, 2010. DOI: 10.1214/09-AOAS285.

_____. Bayesian CART Model Search. **Journal of the American Statistical Association**, v. 93, n. 443, p. 935–948, 1998.

- CIFTCI, Sabri. Modernization, Islam, or Social Capital: What Explains Attitudes Toward Democracy in the Muslim World? en. **Comparative Political Studies**, v. 43, n. 11, p. 1442–1470, nov. 2010. ISSN 0010-4140, 1552-3829. DOI: 10.1177/0010414010371903.
- CLARK, William Roberts; GOLDER, Matt. Big Data, Causal Inference, and Formal Theory: Contradictory Trends in Political Science?: Introduction. en. **PS: Political Science & Politics**, v. 48, n. 1, p. 65–70, jan. 2015. ISSN 1049-0965, 1537-5935. DOI: 10.1017/S1049096514001759.
- CLINTON, Joshua; JACKMAN, Simon; RIVERS, Douglas. The Statistical Analysis of Roll Call Data. en. **American Political Science Review**, v. 98, n. 2, p. 355–370, mai. 2004. ISSN 0003-0554, 1537-5943. DOI: 10.1017/S0003055404001194.
- COHEN, William W. Fast Effective Rule Induction. In: IN Proceedings of the Twelfth International Conference on Machine Learning. [S.l.]: Morgan Kaufmann, 1995. p. 115–123.
- CUMMINS, Robert; POLLOCK, John. **Philosophy and AI: Essays at the Interface**. First. [S.l.]: MIT Press, 1992. (A Bradford Book). ISBN 978-0-262-03180-6.
- DAHL, Robert Alan. **Polyarchy: Participation and Opposition**. [S.l.]: Yale University Press, 1973. ISBN 978-0-300-15357-6.
- DALTON, Russell J.; WELZEL, Christian (Ed.). **The Civic Culture Transformed: From Allegiant to Assertive Citizens**. New York, NY: Cambridge University Press, 2014. ISBN 978-1-107-03926-1 978-1-107-68272-6.
- DAVIDOV, Eldad. Measurement Equivalence of Nationalism and Constructive Patriotism in the ISSP: 34 Countries in a Comparative Perspective. en. **Political Analysis**, v. 17, n. 1, p. 64–82, 2009/ed. ISSN 1047-1987, 1476-4989. DOI: 10.1093/pan/mpn014.
- _____. Nationalism and Constructive Patriotism: A Longitudinal Test of Comparability in 22 Countries with the ISSP. en. **International Journal of Public Opinion Research**, v. 23, n. 1, p. 88–103, mar. 2011. ISSN 0954-2892. DOI: 10.1093/ijpor/edq031.
- DE FIGUEIREDO, Rui J. P.; ELKINS, Zachary. Are Patriots Bigots? An Inquiry into the Vices of In-Group Pride. en. **American Journal of Political Science**, v. 47, n. 1, p. 171–188, jan. 2003. ISSN 1540-5907. DOI: 10.1111/1540-5907.00012.

DEMPSTER, A. P.; LAIRD, N. M.; RUBIN, D. B. Maximum Likelihood from Incomplete Data Via the *EM* Algorithm. en. **Journal of the Royal Statistical Society: Series B (Methodological)**, v. 39, n. 1, p. 1–22, set. 1977. ISSN 00359246. DOI: 10.1111/j.2517-6161.1977.tb01600.x.

DEMSAR, Janez. Statistical Comparisons of Classifiers over Multiple Data Sets. **Journal of Machine Learning Research**, v. 7, p. 1–30, 2006.

DENISON, David G. T.; MALLICK, Bani K.; SMITH, Adrian F. M. A Bayesian CART Algorithm. en. **Biometrika**, v. 85, n. 2, p. 363–377, jun. 1998. ISSN 0006-3444. DOI: 10.1093/biomet/85.2.363.

DIAMOND, Larry (Ed.). **Political Culture and Democracy in Developing Countries: Textbook Edition**. Boulder: Lynne Rienner Pub, abr. 1994. ISBN 978-1-55587-515-2.

DIETTERICH, Thomas G. Approximate Statistical Tests for Comparing Supervised Classification Learning Algorithms. **Neural Computation**, v. 10, n. 7, p. 1895–1923, out. 1998. ISSN 0899-7667. DOI: 10.1162/089976698300017197.

DILLMAN, Don A.; PHELPS, Glenn et al. Response Rate and Measurement Differences in Mixed-Mode Surveys Using Mail, Telephone, Interactive Voice Response (IVR) and the Internet. **Social Science Research**, v. 38, n. 1, p. 1–18, mar. 2009. ISSN 0049-089X. DOI: 10.1016/j.ssresearch.2008.03.007.

DILLMAN, Don A.; SMYTH, Jolene D.; CHRISTIAN, Leah Melani. **Internet, Phone, Mail, and Mixed-Mode Surveys: The Tailored Design Method**. 4th edition. Hoboken: Wiley, 2014. ISBN 978-1-118-45614-9.

DIMITROVA-GRAJZL, Valentina; EASTWOOD, Jonathan; GRAJZL, Peter. The Longevity of National Identity and National Pride: Evidence from Wider Europe. en. **Research & Politics**, v. 3, n. 2, p. 2053168016653424, abr. 2016. ISSN 2053-1680. DOI: 10.1177/2053168016653424.

DOMINGOS, Pedro. A Few Useful Things to Know about Machine Learning. en. **Communications of the ACM**, v. 55, n. 10, p. 78, out. 2012. ISSN 00010782. DOI: 10.1145/2347736.2347755.

_____. Occam's Two Razors: The Sharp and the Blunt. In: PROCEEDINGS of the Fourth International Conference on Knowledge Discovery and Data Mining. [S.l.]: American Association for Artificial Intelligence, ago. 1998. p. 37–43.

DOMINGOS, Pedro. **The Master Algorithm: How the Quest for the Ultimate Learning Machine Will Remake Our World**. [S.l.]: Basic Books, 2015. ISBN 978-0-465-06570-7.

_____. The Role of Occam's Razor in Knowledge Discovery. en, p. 19, 1999.

DOOB, Leonard William. **Patriotism and Nationalism: Their Psychological Foundations**. [S.l.]: Greenwood Press, 1976. ISBN 978-0-8371-8978-9.

DOSHI-VELEZ, Finale; KIM, Been. Towards A Rigorous Science of Interpretable Machine Learning. **arXiv:1702.08608 [cs, stat]**, fev. 2017. arXiv: 1702.08608 [cs, stat].

DRECHSLER, Jörg. Multiple Imputation of Multilevel Missing Data—Rigor Versus Simplicity. en. **Journal of Educational and Behavioral Statistics**, fev. 2015. DOI: 10.3102/1076998614563393.

DU, Mengnan; LIU, Ninghao; HU, Xia. Techniques for Interpretable Machine Learning. **arXiv:1808.00033 [cs, stat]**, mai. 2019. arXiv: 1808.00033 [cs, stat].

DUCHESNE, Sophie; FROGNIER, André-Paul. Why Is It so Difficult to Know If National Pride Leads the Way to European Identity or Prevents It? en. **SSRN Electronic Journal**, 2007. ISSN 1556-5068. DOI: 10.2139/ssrn.1522904.

DURRANT, Gabriele B. Imputation Methods for Handling Item - Nonresponse in the Social Sciences: A Methodological Review. **NCRM Methods Review Papers**, NCRM/002, 2005.

EASTON, David. **A Framework for Political Analysis**. [S.l.]: Prentice-Hall, 1965.

ECKEL, Matthew C. Inequality, Elite Messaging and National Pride. en. In: EUROACADEMIA Conference: Identities and Identifications: Politicized Uses of Collective Identities. Florence: [s.n.], out. 2014. p. 24.

ECKSTEIN, Harry. A Culturalist Theory of Political Change. **American Political Science Review**, v. 82, n. 03, p. 789–804, 1988.

_____. Congruence Theory Explained. **Center for the Study of Democracy**, jul. 1997.

_____. **Division and Cohesion in Democracy: A Study of Norway**. [S.l.]: Princeton University Press, mar. 2015. ISBN 978-1-4008-6816-2.

EFRON, Bradley; HASTIE, Trevor. **Computer Age Statistical Inference: Algorithms, Evidence, and Data Science**. First. [S.l.]: Cambridge University Press, jul. 2016. ISBN 978-1-107-14989-2 978-1-316-57653-3. DOI: 10.1017/CBO9781316576533.

ELSTER, Jon. **Peças e engrenagens das ciencias sociais**. Rio de Janeiro: Relume Dumara, 1994. ISBN 978-85-85427-91-7.

_____. **The Cement of Society: A Survey of Social Order**. Cambridge England ; New York: Cambridge University Press, jul. 1989. ISBN 978-0-521-37607-5.

ENZWEILER, Markus. The Mobile Revolution – Machine Intelligence for Autonomous Vehicles. en. **it - Information Technology**, v. 57, n. 3, jan. 2015. ISSN 1611-2776, 2196-7032. DOI: 10.1515/itit-2015-0009.

EUGSTER, Manuel J A; HOTHORN, Torsten; LEISCH, Friedrich. **Exploratory and Inferential Analysis of Benchmark Experiments**. en. Munich, 2008. p. 29.

_____. Domain-Based Benchmark Experiments: Exploratory and Inferential Analysis. en. **Austrian Journal of Statistics**, v. 41, n. 1, p. 5, fev. 2016. ISSN 1026597X. DOI: 10.17713/ajs.v41i1.185.

EVANS, M. D. R.; KELLEY, Jonathan. National Pride in the Developed World: Survey Data from 24 Nations. en. **International Journal of Public Opinion Research**, v. 14, n. 3, p. 303–338, set. 2002. ISSN 0954-2892. DOI: 10.1093/ijpor/14.3.303.

FABRYKANT, Marharyta; MAGUN, Vladimir. **Grounded and Normative Dimensions of National Pride in Comparative Perspective**. en. Rochester, NY, abr. 2015.

FABRYKANT, Marharyta S. Value of (Expla)Nation: Testing Modernist Theories of Nationalism. **Journal of Siberian Federal University. Humanities & Social Sciences**, v. 1, n. 7, p. 152–168, 2014.

FAVARETTO, Maddalena et al. What Is Your Definition of Big Data? Researchers' Understanding of the Phenomenon of the Decade. en. **PLOS ONE**, Public Library of Science, v. 15, n. 2, e0228987, fev. 2020. ISSN 1932-6203. DOI: 10.1371/journal.pone.0228987.

FERNÁNDEZ-DELGADO, Manuel et al. Do We Need Hundreds of Classifiers to Solve Real World Classification Problems? **Journal of Machine Learning Research**, v. 15, p. 3133–3181, 2014.

FISHER, Aaron; RUDIN, Cynthia; DOMINICI, Francesca. All Models Are Wrong, but Many Are Useful: Learning a Variable's Importance by Studying an Entire Class of Prediction Models Simultaneously. **arXiv:1801.01489 [stat]**, jan. 2018. arXiv: 1801.01489 [stat].

FITZMAURICE, Garrett M. et al. **Handbook of Missing Data Methodology**. [S.l.]: CRC Press, 2014. (Chapman & Hall/CRC Handbooks of Modern Statistical Methods). ISBN 978-1-4398-5461-7.

FREUND, Yoav; SCHAPIRE, Robert E. A Decision-Theoretic Generalization of On-Line Learning and an Application to Boosting. **Journal of Computer and System Sciences**, v. 55, n. 1, p. 119–139, ago. 1997. ISSN 0022-0000. DOI: 10.1006/jcss.1997.1504.

_____. A Short Introduction to Boosting. en. **Journal of Japanese Society for Artificial Intelligence**, v. 14, n. 5, p. 771–780, set. 1999.

FRIEDMAN, Jerome; HASTIE, Trevor; TIBSHIRANI, Rob. Regularization Paths for Generalized Linear Models via Coordinate Descent. **Journal of statistical software**, v. 33, n. 1, p. 1–22, 2010. ISSN 1548-7660.

FRIEDMAN, Jerome; HASTIE, Trevor et al. **Glmnet: Lasso and Elastic-Net Regularized Generalized Linear Models**. [S.l.: s.n.], 2019.

FRIEDMAN, Jerome H. Greedy Function Approximation: A Gradient Boosting Machine. **The Annals of Statistics**, v. 29, n. 5, p. 1189–1232, 2001. ISSN 0090-5364.

_____. Stochastic Gradient Boosting. **Journal Computational Statistics & Data Analysis - Nonlinear methods and data mining**, v. 38, n. 4, p. 367–378, fev. 2002.

FRIEDMAN, Jerome H.; HASTIE, Trevor; TIBSHIRANI, Robert. Additive Logistic Regression: A Statistical View of Boosting. en, p. 45, 1998.

FRIEDMAN, Jerome H.; POPESCU, Bogdan E. Predictive Learning via Rule Ensembles. **arXiv:0811.1679 [stat]**, nov. 2008. DOI: 10.1214/07-AOAS148. arXiv: 0811.1679 [stat].

GARSON, G. David. Expert Systems: An Overview for Social Scientists. **Social Science Computer Review**, v. 8, n. 3, p. 387–410, 1990.

GECAS, V. The Self-Concept. **Annual Review of Sociology**, v. 8, n. 1, p. 1–33, 1982. DOI: 10.1146/annurev.so.08.080182.000245.

GELLNER, Ernest. **Nations and Nationalism**. [S.l.]: Cornell University Press, 2008. ISBN 978-0-8014-7500-9.

GELMAN, Andrew. Scaling Regression Inputs by Dividing by Two Standard Deviations. en. **Statistics in Medicine**, v. 27, n. 15, p. 2865–2873, jul. 2008. ISSN 02776715, 10970258. DOI: 10.1002/sim.3107.

- GELMAN, Andrew; JAKULIN, Aleks et al. A Weakly Informative Default Prior Distribution for Logistic and Other Regression Models. en. **The Annals of Applied Statistics**, v. 2, n. 4, p. 1360–1383, dez. 2008. ISSN 1932-6157. DOI: 10.1214/08-AOAS191.
- GELMAN, Andrew; KING, Gary; LIU, Chuanhai. Not Asked and Not Answered: Multiple Imputation for Multiple Surveys. **Journal of the American Statistical Association**, v. 93, n. 443, p. 846–857, 1998. ISSN 0162-1459. DOI: 10.2307/2669819.
- GIUFFRIDA, Angela. Nicola Zingaretti Pledges to Take on Italy's Rightwing Government. en-GB. **The Guardian**, mai. 2019. ISSN 0261-3077.
- GOLDSTEIN, Alex et al. Peeking Inside the Black Box: Visualizing Statistical Learning with Plots of Individual Conditional Expectation. **arXiv:1309.6392 [stat]**, set. 2013. arXiv: 1309.6392 [stat].
- GREEN, Elliott. Ethnicity, National Identity and the State: Evidence from Sub-Saharan Africa. en. **British Journal of Political Science**, Cambridge University Press, v. 50, n. 2, p. 757–779, abr. 2020. ISSN 0007-1234, 1469-2112. DOI: 10.1017/S0007123417000783.
- GREEN, Eva G. T. et al. Nationalism and Patriotism as Predictors of Immigration Attitudes in Switzerland: A Municipality-Level Analysis. en. **Swiss Political Science Review**, v. 17, n. 4, p. 369–393, 2011. ISSN 1424-7755. DOI: 10.1111/j.1662-6370.2011.02030.x.
- GREENFELD, Liah. **Nationalism: Five Roads to Modernity**. Cambridge, Mass: Harvard University Press, 1992. ISBN 978-0-674-60318-9.
- GREENWALD, Hal S.; OERTEL, Carsten K. Future Directions in Machine Learning. English. **Frontiers in Robotics and AI**, v. 3, 2017. ISSN 2296-9144. DOI: 10.3389/frobt.2016.00079.
- GREENWELL, Brandon et al. **Gbm: Generalized Boosted Regression Models**. [S.l.: s.n.], 2019.
- GREGORUTTI, Baptiste; MICHEL, Bertrand; SAINT-PIERRE, Philippe. Correlation and Variable Importance in Random Forests. en. **Statistics and Computing**, v. 27, n. 3, p. 659–678, mai. 2017. ISSN 1573-1375. DOI: 10.1007/s11222-016-9646-1.
- _____. Grouped Variable Importance with Random Forests and Application to Multivariate Functional Data Analysis. en, p. 32, 2013.

GRIMMER, Justin; STEWART, Brandon M. Text as Data: The Promise and Pitfalls of Automatic Content Analysis Methods for Political Texts. en. **Political Analysis**, v. 21, n. 3, p. 267–297, 2013. ISSN 1047-1987, 1476-4989. DOI: 10.1093/pan/mps028.

GUIMARÃES, Manoel Luis Lima Salgado. Nação e Civilização Nos Trópicos: O Instituto Histórico Geográfico Brasileiro e o Projeto de Uma História Nacional. **Revista Estudos Históricos**, v. 1, n. 1, p. 5–27, jan. 1988. ISSN 2178-1494.

HA, Shang E.; JANG, Seung-Jin. National Identity, National Pride, and Happiness: The Case of South Korea. en. **Social Indicators Research**, v. 121, n. 2, p. 471–482, abr. 2015. ISSN 0303-8300, 1573-0921. DOI: 10.1007/s11205-014-0641-7.

HABERMAS, Jürgen. **Citizenship and National Identity: Some Reflections on the Future of Europe**. [S.l.: s.n.], 1992.

HAERPFER, Christian W.; KIZILOVA, Kseniya. 7. Support for Democracy in Postcommunist Europe and Post-Soviet Eurasia. In: DALTON, Russell J.; WELZEL, Christian (Ed.). **The Civic Culture Transformed: From Allegiant to Assertive Citizens**. New York, NY: Cambridge University Press, 2014. ISBN 978-1-107-03926-1 978-1-107-68272-6.

HAGGARD, Stephan; KAUFMAN, Robert R.; LONG, James D. Income, Occupation, and Preferences for Redistribution in the Developing World. en. **Studies in Comparative International Development**, v. 48, n. 2, p. 113–140, jun. 2013. ISSN 1936-6167. DOI: 10.1007/s12116-013-9129-8.

HAN, Kyung Joon. Income Inequality, International Migration, and National Pride: A Test of Social Identification Theory. **International Journal of Public Opinion Research**, v. 25, n. 4, p. 502–521, mai. 2013. DOI: 10.1093/ijpor/edt011.

HARRELL, Frank. **Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression, and Survival Analysis**. New York: Springer-Verlag, 2001. (Springer Series in Statistics). ISBN 978-1-4419-2918-1.

HASTIE, Trevor; TIBSHIRANI, Robert; FRIEDMAN, Jerome H. **The Elements of Statistical Learning Data Mining, Inference, and Prediction**. Second. [S.l.]: Springer, jan. 2017. (Springer Series in Statistics).

HERDER, Johann Gottfried; FORSTER, Michael N. **Philosophical Writings**. Cambridge, UK ; New York: Cambridge University Press, 2002. (Cambridge Texts in the History of Philosophy). ISBN 978-0-521-79088-8 978-0-521-79409-1.

- HJERM, Mikael. National Identities, National Pride and Xenophobia: A Comparison of Four Western Countries. **Acta Sociologica**, v. 41, n. 4, p. 335–347, 1998.
- HOBBSAWM, E. J. **Nations and Nationalism since 1780: Programme, Myth, Reality**. 2nd ed. Cambridge [England] ; New York: Cambridge University Press, 1992. ISBN 978-0-521-43961-9.
- HOFMAN, Jake M.; SHARMA, Amit; WATTS, Duncan J. Prediction and Explanation in Social Systems. en. **Science**, American Association for the Advancement of Science, v. 355, n. 6324, p. 486–488, fev. 2017. ISSN 0036-8075, 1095-9203. DOI: 10.1126/science.aal3856.
- HOFSTEDE, Geert et al. Comparing Regional Cultures Within a Country: Lessons From Brazil. en. **Journal of Cross-Cultural Psychology**, v. 41, n. 3, p. 336–352, mai. 2010. ISSN 0022-0221, 1552-5422. DOI: 10.1177/0022022109359696.
- HOLANDA, Sérgio de. **Raízes do Brasil**. 7. impr. São Paulo: Companhia das Letras, 1999. ISBN 978-85-7164-448-9.
- HOLANDA, Sérgio Buarque de. **Visão do paraíso: os motivos edênicos no descobrimento e colonização do Brasil**. [S.l.: s.n.], 2000. ISBN 978-85-11-13109-3 978-85-7402-189-8.
- HOLTE, Robert C. Very Simple Classification Rules Perform Well on Most Commonly Used Datasets. en. **Machine Learning**, v. 11, n. 1, p. 63–90, abr. 1993. ISSN 1573-0565. DOI: 10.1023/A:1022631118932.
- HONAKER, James; KING, Gary. What to Do about Missing Values in Time-Series Cross-Section Data. en. **American Journal of Political Science**, v. 54, n. 2, p. 561–581, abr. 2010. ISSN 00925853, 15405907. DOI: 10.1111/j.1540-5907.2010.00447.x.
- HONAKER, James; KING, Gary; BLACKWELL, Matthew. Amelia II: A Program for Missing Data. en. **Journal of Statistical Software**, v. 45, n. 7, 2011. ISSN 1548-7660. DOI: 10.18637/jss.v045.i07.
- HOPFIELD, J. J. Neural Networks and Physical Systems with Emergent Collective Computational Abilities. en. **Proceedings of the National Academy of Sciences**, v. 79, n. 8, p. 2554–2558, abr. 1982. ISSN 0027-8424, 1091-6490. DOI: 10.1073/pnas.79.8.2554.
- HORNIK, Kurt et al. **RWeka Interface**. [S.l.: s.n.], fev. 2020.

HORVÁTH, Tomáš; MANTOVANI, Rafael G.; DE CARVALHO, André C. P. L. F. Effects of Random Sampling on SVM Hyper-Parameter Tuning. en. In: MADUREIRA, Ana Maria et al. (Ed.). **Intelligent Systems Design and Applications**. [S.l.]: Springer International Publishing, 2017. (Advances in Intelligent Systems and Computing), p. 268–278. ISBN 978-3-319-53480-0.

HOTHORN, Torsten; HORNIK, Kurt; STROBL, Carolin et al. **Party: A Laboratory for Recursive Partytioning**. [S.l.: s.n.], 2020.

HOTHORN, Torsten; HORNIK, Kurt; ZEILEIS, Achim. Unbiased Recursive Partitioning: A Conditional Inference Framework. en. **Journal of Computational and Graphical Statistics**, v. 15, n. 3, p. 651–674, set. 2006. ISSN 1061-8600, 1537-2715. DOI: 10.1198/106186006X133933.

HOTHORN, Torsten; LEISCH, Friedrich et al. The Design and Analysis of Benchmark Experiments. **Journal of Computational and Graphical Statistics**, p. 675–699, 2005. DOI: 10.1198/106186005X59630.

HUANG, G. et al. Extreme Learning Machine for Regression and Multiclass Classification. **IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)**, v. 42, n. 2, p. 513–529, abr. 2012. DOI: 10.1109/TSMCB.2011.2168604.

HUTTER, Frank; HOOS, Holger H.; LEYTON-BROWN, Kevin. Sequential Model-Based Optimization for General Algorithm Configuration. en. In: _____. **Learning and Intelligent Optimization**. Berlin, Heidelberg: Springer, 2011. (Lecture Notes in Computer Science), p. 507–523. ISBN 978-3-642-25566-3. DOI: 10.1007/978-3-642-25566-3_40.

IANNONE, Richard. **DiagrammeR: Graph/Network Visualization**. [S.l.: s.n.], 2020.

IHAKA, Ross; GENTLEMAN, Robert. R: A Language for Data Analysis and Graphics. **Journal of Computational and Graphical Statistics**, [American Statistical Association, Taylor & Francis, Ltd., Institute of Mathematical Statistics, Interface Foundation of America], v. 5, n. 3, p. 299–314, 1996. ISSN 1061-8600. DOI: 10.2307/1390807.

INGLEHART, Ronald. **Culture Shift in Advanced Industrial Society**. Princeton, N.J.: Princeton University Press, dez. 1989. ISBN 978-0-691-02296-3.

_____. **Modernization and Postmodernization: Cultural, Economic, and Political Change in 43 Societies**. [S.l.]: Princeton: Princeton University Press, 1997.

INGLEHART, Ronald. The Silent Revolution in Europe: Intergenerational Change in Post-Industrial Societies. **The American Political Science Review**, v. 65, n. 4, p. 991–1017, 1971. ISSN 0003-0554. DOI: 10.2307/1953494.

_____. **The Silent Revolution: Changing Values and Political Styles Among Western Publics**. [S.l.]: Princeton University Press, 1977. ISBN 978-1-4008-6958-9.

INGLEHART, Ronald; HAERPFER, Christian W. et al. (Ed.). **World Values Survey: All Rounds**. Madrid: JD Systems Institute., 2019.

INGLEHART, Ronald; WELZEL, Christian. Democratic Institutions and Political Culture: Misconceptions in Addressing the Ecological Fallacy. [S.l.], 2003.

_____. **Modernization, Cultural Change, and Democracy: The Human Development Sequence**. Cambridge, UK ; New York: Cambridge University Press, ago. 2005. ISBN 978-0-521-60971-5.

INGLEHART, Ronald F. **Cultural Evolution: People's Motivations Are Changing, and Reshaping the World**. First. [S.l.]: Cambridge University Press, mar. 2018. ISBN 978-1-108-61388-0 978-1-108-48931-7 978-1-108-46477-2. DOI: 10.1017/9781108613880.

INTELLIGENCE, Panel on Computer Science and Artificial; COUNCIL, National Research. **Computer Science and Artificial Intelligence**. [S.l.]: National Academies Press, 1997. ISBN 978-0-309-05831-5.

JAMES, Gareth et al. **An Introduction to Statistical Learning**. New York, NY: Springer New York, 2013. v. 103. (Springer Texts in Statistics). ISBN 978-1-4614-7137-0 978-1-4614-7138-7. DOI: 10.1007/978-1-4614-7138-7.

JANCSÓ, István (Ed.). **Brasil: Formação Do Estado e Da Nação**. São Paulo : Ijuí: Editora Hucitec : FAPESP ; Editora Unijuí, 2003. (Estudos Históricos, 50). ISBN 978-85-271-0613-9.

JONES, Frank L.; SMITH, Philip. Diversity and Commonality in National Identities: An Exploratory Analysis of Cross-National Patterns. en. **Journal of Sociology**, v. 37, n. 1, p. 45–63, mar. 2001. ISSN 1440-7833. DOI: 10.1177/144078301128756193.

JOSSE, Julie et al. On the Consistency of Supervised Learning with Missing Values. **arXiv:1902.06931 [cs, math, stat]**, fev. 2019. arXiv: 1902.06931 [cs, math, stat].

KASPAROV, Garry; GREENGARD, Mig. **Deep Thinking: Where Machine Intelligence Ends and Human Creativity Begins**. 1 edition. New York: PublicAffairs, mai. 2017. ISBN 978-1-61039-786-5.

KAVETSOS, Georgios. National Pride: War Minus the Shooting. en. **Social Indicators Research**, v. 106, n. 1, p. 173–185, mar. 2012. ISSN 0303-8300, 1573-0921. DOI: 10.1007/s11205-011-9801-1.

KEDOURIE, Elie. **Nationalism**. [S.l.: s.n.], 1960.

KIMBER, Richard. Artificial Intelligence and the Study of Democracy. en. **Social Science Computer Review**, v. 9, n. 3, p. 381–398, out. 1991. ISSN 0894-4393. DOI: 10.1177/089443939100900303.

KING, Gary. Restructuring the Social Sciences: Reflections from Harvard's Institute for Quantitative Social Science. **PS: Political Science and Politics**, v. 47, n. 1, p. 165–172, 2014. DOI: 10.1017/S1049096513001534.

KING, Gary; PAN, Jennifer; ROBERTS, Margaret. How Censorship in China Allows Government Criticism but Silences Collective Expression. **American Political Science Review**, v. 107, 2 (May), p. 1–18, 2013.

KING, Gary et al. Analyzing Incomplete Political Science Data: An Alternative Algorithm for Multiple Imputation. en. **American Political Science Review**, v. 95, n. 1, p. 21, 2001.

KING, Ross D. et al. The Robot Scientist Project. en. In: _____. **Discovery Science**. [S.l.]: Springer Berlin Heidelberg, 2005. (Lecture Notes in Computer Science), p. 16–25. ISBN 978-3-540-31698-5.

KOSTERMAN, Rick; FESHBACH, Seymour. Toward a Measure of Patriotic and Nationalistic Attitudes. **Political Psychology**, v. 10, n. 2, p. 257, jun. 1989. ISSN 0162895X. DOI: 10.2307/3791647.

KOTTHOFF, Lars et al. Auto-WEKA: Automatic Model Selection and Hyperparameter Optimization in WEKA. In: HUTTER, Frank; KOTTHOFF, Lars; VANSCHOREN, Joaquin (Ed.). **Automated Machine Learning**. Cham: Springer International Publishing, 2019. p. 81–95. ISBN 978-3-030-05317-8 978-3-030-05318-5. DOI: 10.1007/978-3-030-05318-5_4.

KRAUSE, Josua et al. A Workflow for Visual Diagnostics of Binary Classifiers Using Instance-Level Explanations. **arXiv:1705.01968 [cs, stat]**, out. 2017. arXiv: 1705.01968 [cs, stat].

KUHN, Max. **Caret: Classification and Regression Training**. [S.l.: s.n.], abr. 2019.

KUHN, Max; JOHNSON, Kjell. **Applied Predictive Modeling**. New York, NY: Springer New York, 2013. ISBN 978-1-4614-6848-6 978-1-4614-6849-3. DOI: 10.1007/978-1-4614-6849-3.

_____. **Feature Engineering and Selection: A Practical Approach for Predictive Models**. 1st. [S.l.]: Chapman & Hall, 2019. (CRC Data Science Series).

KUHN, Max; QUINLAN, Ross. **C50: C5.0 Decision Trees and Rule-Based Models**. [S.l.: s.n.], 2020.

LAN, Xiaohuan; LI, Ben G. The Economics of Nationalism. en. **American Economic Journal: Economic Policy**, v. 7, n. 2, p. 294–325, mai. 2015. ISSN 1945-7731, 1945-774X. DOI: 10.1257/pol.20130020.

LAUGEL, Thibault et al. Defining Locality for Surrogates in Post-Hoc Interpretability. **arXiv:1806.07498 [cs, stat]**, jun. 2018. arXiv: 1806.07498 [cs, stat].

LAZER, D. et al. SOCIAL SCIENCE: Computational Social Science. en. **Science**, v. 323, n. 5915, p. 721–723, fev. 2009. ISSN 0036-8075, 1095-9203. DOI: 10.1126/science.1167742.

LECUN, Yann; BENGIO, Yoshua; HINTON, Geoffrey. Deep Learning. en. **Nature**, v. 521, n. 7553, p. 436–444, mai. 2015. ISSN 0028-0836, 1476-4687. DOI: 10.1038/nature14539.

LEITE, Eldo Lima et al. Nationalism, Patriotism, and Essentialism in the Construction of Brazilian National Identity. en. **Trends in Psychology**, v. 26, n. 4, p. 2063–2075, out. 2018. ISSN 2358-1883. DOI: 10.9788/tp2018.4-13pt.

LESSA, Carlos. Nation and Nationalism Based on the Brazilian Experience. **Estudos Avançados**, v. 22, n. 62, p. 237–256, abr. 2008. ISSN 0103-4014. DOI: 10.1590/S0103-40142008000100016.

LEVIN, Ines; POMARES, Julia; ALVAREZ, R. Michael. Using Machine Learning Algorithms to Detect Election Fraud. In: ALVAREZ, R. Michael (Ed.). **Computational Social Science**. Cambridge: Cambridge University Press, 2016. p. 266–294. ISBN 978-1-316-25734-0. DOI: 10.1017/CBO9781316257340.012.

- LITTLE, R.; RUBIN, D. **Statistical Analysis with Missing Data**. [S.l.: s.n.], 1987. ISBN 978-1-69170-183-4.
- LITTLE, Roderick J. A.; RUBIN, Donald B. Maximum Likelihood for General Patterns of Missing Data: Introduction and Theory with Ignorable Nonresponse. In: **STATISTICAL Analysis with Missing Data**. [S.l.]: John Wiley & Sons, Ltd, 2014. p. 164–189. ISBN 978-1-119-01356-3. DOI: 10.1002/9781119013563.ch8.
- LOHUIZEN, Jan van; SAMOBYL, Robert Wayne. Method Effects and Robo-Polls. en. **Survey Practice**, v. 4, n. 1, p. 3057, fev. 2011. DOI: 10.29115/SP-2011-0005.
- LOUPPE, Gilles et al. Understanding Variable Importances in Forests of Randomized Trees. en, p. 9, 2013.
- LUCKIN, Rose et al. **Intelligence Unleashed: An Argument for AI in Education**. [S.l.: s.n.], 2016. ISBN 978-0-9924248-8-6.
- LUO, Jiaming; CAO, Yuan; BARZILAY, Regina. Neural Decipherment via Minimum-Cost Flow: From Ugaritic to Linear B. **arXiv:1906.06718 [cs]**, jun. 2019. arXiv: 1906.06718 [cs].
- MACIÀ, Núria et al. Learner Excellence Biased by Data Set Selection: A Case for Data Characterisation and Artificial Data Sets. **Pattern Recognition**, v. 46, n. 3, p. 1054–1066, mar. 2013. ISSN 0031-3203. DOI: 10.1016/j.patcog.2012.09.022.
- MACLIN, R.; OPITZ, D. Popular Ensemble Methods: An Empirical Study. **Journal of Artificial Intelligence Research**, v. 11, p. 169–198, ago. 1999. ISSN 1076-9757. DOI: 10.1613/jair.614. arXiv: 1106.0257.
- MARQUÉS, A. I.; GARCÍA, V.; SÁNCHEZ, J. S. Exploring the Behaviour of Base Classifiers in Credit Scoring Ensembles. **Expert Systems with Applications**, v. 39, n. 11, p. 10244–10250, set. 2012. ISSN 0957-4174. DOI: 10.1016/j.eswa.2012.02.092.
- MARTIN, Todd M. et al. Does Rational Selection of Training and Test Sets Improve the Outcome of QSAR Modeling? **Journal of Chemical Information and Modeling**, American Chemical Society, v. 52, n. 10, p. 2570–2578, out. 2012. ISSN 1549-9596. DOI: 10.1021/ci300338w.
- MASLOW, A. H. A Theory of Human Motivation. **Psychological Review**, American Psychological Association, US, v. 50, n. 4, p. 370–396, 1943. ISSN 1939-1471(Electronic),0033-295X(Print). DOI: 10.1037/h0054346.

MCCARTHY, J et al. A Proposal for the Dartmouth Summer Research Project on Artificial Intelligence. en, p. 13, ago. 1955.

MCCULLOCH, Warren S.; PITTS, Walter. A Logical Calculus of the Ideas Immanent in Nervous Activity. en. **The bulletin of mathematical biophysics**, v. 5, n. 4, p. 115–133, dez. 1943. ISSN 1522-9602. DOI: 10.1007/BF02478259.

MERCALDO, Sarah Fletcher; BLUME, Jeffrey D. Missing Data and Prediction. **arXiv:1704.08192 [stat]**, abr. 2017. arXiv: 1704.08192 [stat].

METZ, Cade. A Breakthrough for A.I. Technology: Passing an 8th-Grade Science Test. en-US. **The New York Times**, set. 2019. ISSN 0362-4331.

MICHELAT, Guy; THOMAS, Jean Pierre Hubert. **Dimensions du nationalisme**. [S.I.]: A. Colin., 1966.

MILES, Andrew. Obtaining Predictions from Models Fit to Multiply Imputed Data. en. **Sociological Methods & Research**, v. 45, n. 1, p. 175–185, fev. 2016. ISSN 0049-1241. DOI: 10.1177/0049124115610345.

MILLER, Tim. Explainable Artificial Intelligence: What Were You Thinking? In: ARTIFICIAL Intelligence for Better or for Worse. [S.I.]: Future Leaders, 2019.

_____. Explanation in Artificial Intelligence: Insights from the Social Sciences. **Artificial Intelligence**, v. 267, p. 1–38, fev. 2019. ISSN 0004-3702. DOI: 10.1016/j.artint.2018.07.007.

MILLS, C. Wright. **The Sociological Imagination**. [S.I.]: Oxford University Press, 1959. ISBN 13 978 0-1-513373-8.

MINSKY, Marvin. Steps toward Artificial Intelligence. en. **Proceedings of the IRE**, v. 49, n. 1, p. 8–30, jan. 1961. ISSN 0096-8390. DOI: 10.1109/JRPROC.1961.287775.

MITCHELL, Tom M. Explanation-Based Generalization: A Unifying View. en. **Machine Learning**, v. 1, p. 47–80, 1986.

_____. **Machine Learning**. New York: McGraw-Hill, 1997. (McGraw-Hill Series in Computer Science). ISBN 978-0-07-042807-2.

MOADDEL, Mansoor; TESSLER, Mark; INGLEHART, Ronald. Foreign Occupation and National Pride: The Case of Iraq. **Public Opinion Quarterly**, v. 72, n. 4, p. 677–705, 2008.

- MOISÉS, José Álvaro. A Desconfiança Nas Instituições Democráticas. **Opinião pública**, v. 11, n. 1, p. 33–63, 2005.
- MOISÉS, José Álvaro; MENEGUELLO, Rachel (Ed.). **A Desconfiança Política e Os Seus Impactos Na Qualidade Na Democracia - o Caso Do Brasil**. [S.l.]: University of São Paulo Press, 2013.
- MOLNAR, Christoph. **Interpretable Machine Learning**. Leanpub. [S.l.: s.n.], 2019.
- MOUSELIMIS, Lampros; GOSSO, Alberto. **elmNNRcpp: The Extreme Learning Machine Algorithm**. [S.l.: s.n.], jul. 2018.
- MULLER, Edward N.; SELIGSON, Mitchell A. Civic Culture and Democracy: The Question of Causal Relationships. en. **American Political Science Review**, v. 88, n. 3, p. 635–652, set. 1994. ISSN 0003-0554, 1537-5943. DOI: 10.2307/2944800.
- MÜLLER, Jan-Werner. **Constitutional Patriotism**. [S.l.: s.n.], 2009. ISBN 978-1-4008-2808-1.
- MÜLLER, Kirill; WICKHAM, Hadley. **Tibble: Simple Data Frames**. [S.l.: s.n.], 2020.
- MUMMENDEY, Amélie; KLINK, Andreas; BROWN, Rupert. Nationalism and Patriotism: National Identification and out-Group Rejection. en. **British Journal of Social Psychology**, v. 40, n. 2, p. 159–172, jun. 2001. ISSN 2044-8309. DOI: 10.1348/014466601164740.
- MUÑOZ, Jordi. From National-Catholicism to Democratic Patriotism? Democratization and Reconstruction of National Pride: The Case of Spain (1981–2000). en. **Ethnic and Racial Studies**, v. 32, n. 4, p. 616–639, mai. 2009. ISSN 0141-9870, 1466-4356. DOI: 10.1080/01419870701710906.
- NASCIMENTO, Francielle M.; BARONE, Dante A. C.; CASTRO, Henrique C. O. de. Social Activism Analysis: An Application of Machine Learning in the World Values Survey. en-US, ago. 2019. ISSN 2516-2314.
- NATEKIN, Alexey; KNOLL, Alois. Gradient Boosting Machines, a Tutorial. English. **Frontiers in Neurobotics**, v. 7, 2013. ISSN 1662-5218. DOI: 10.3389/fnbot.2013.00021.
- NATHANSON, Stephen. **Patriotism, Morality, and Peace**. [S.l.]: Rowman & Littlefield Publishers, 1993.

NEWELL, Allen. Remarks on the Relationship Between Artificial Intelligence and Cognitive Psychology. en. In: BANERJI, R. B.; MESAROVIC, M. D. (Ed.). **Theoretical Approaches to Non-Numerical Problem Solving**. [S.l.]: Springer Berlin Heidelberg, 1970. (Lecture Notes in Operations Research and Mathematical Systems), p. 363–400. ISBN 978-3-642-99976-5.

_____. You Can ' t Play 20 Questions with Nature and Win : Projective Comments on the Papers Of. In:

NEWELL, Allen; SIMON, Herbert A. Computer Science As Empirical Inquiry: Symbols and Search. **Commun. ACM**, v. 19, n. 3, p. 113–126, mar. 1976. ISSN 0001-0782. DOI: 10.1145/360018.360022.

_____. **Human Problem Solving**. [S.l.: s.n.], 1972.

NILSSON, Nils J. **Artificial Intelligence: A New Synthesis**. 1st. [S.l.]: Morgan Kaufmann Publishers, Inc., 1998. ISBN 978-1-55860-467-4.

NINCIC, Miroslav; RAMOS, Jennifer M. The Sources of Patriotism: Survey and Experimental Evidence. **Foreign Policy Analysis**, v. 8, n. 4, p. 373–388, 2012. ISSN 1743-8586. DOI: 10.1111/j.1743-8594.2011.00175.x.

NOH, Emily Euiyoung. Across Most of the Americas, National Pride Is High and Stable, While It Has Plummeted in the U.S. en, p. 16, ago. 2018.

NORRIS, Pippa. Confidence in the United Nations. In: WORLD Values Conference: Society, Politics and Values: 1981- 2006. Istanbul, Turkey: [s.n.], nov. 2006.

_____. **Critical Citizens: Global Support for Democratic Government**. [S.l.]: Oxford University Press, 1999. ISBN 978-0-19-829568-6.

_____. Global Governance & Cosmopolitan Citizens. In: NYE, Joseph S.; DONAHUE, John D. (Ed.). **Governance in a Globalizing World**. Cambridge, Mass. : Washington, D.C: Visions of Governance for the 21st Century ; Brookings Institution Press, 2000. ISBN 978-0-8157-6408-3 978-0-8157-6407-6.

NORRIS, Pippa; INGLEHART, Ronald. **Cosmopolitan Communications: Cultural Diversity in a Globalized World**. 1 edition. New York: Cambridge University Press, ago. 2009. ISBN 978-0-521-73838-5.

NUSSBAUM, Martha. Patriotism and Cosmopolitanism, Martha Nussbaum. en. **The Boston Review**, p. 8, 1994.

NUSSBAUM, Martha; COHEN, Joshua. **For Love of Country?** 1st edition. Boston: Beacon Press, jun. 2002. ISBN 978-0-8070-4329-5.

OLIVEIRA, Lúcia Lippi. **A Questão Nacional Na Primeira República**. Primeira Edição. São Paulo, Brasília: Brasiliense, MCT/CNPq, 1990.

OLLION, Etienne. Machine learning et sciences sociales : savoirs liés? Français. In: SÉMINAIRE DU LADHUL, AUTOMNE 2018 : FAIRE DES SHS AVEC LE NUMÉRIQUE. Unil Lausanne: [s.n.], out. 2018.

OLLION, Étienne; BOELAERT, Julien. Au delà des big data. Les sciences sociales et la multiplication des données numériques. fr. **Sociologie**, v. 6, n. 3, p. 295–310, nov. 2015. ISSN 2108-8845.

PACKAGE), Adam Kapelner and Justin Bleich (R. **bartMachine: Bayesian Additive Regression Trees**. [S.l.: s.n.], mai. 2018.

PAMPLONA, Marco A. State Building and Nation Formation in Iberian America: Independence and the Development of a Political Language of Patriotism in Brazil. English. In: PAPERS and Comments by Former Students. New York: [s.n.], out. 2011.

PATEMAN, Carole. Political Culture, Political Structure and Political Change. **British Journal of Political Science**, v. 1, n. 03, p. 291–305, jul. 1971. ISSN 1469-2112. DOI: 10.1017/S0007123400009133.

PAUL, Thibault Helleputte; Pierre Gramme; Jerome. **LiblineaR: Linear Predictive Models Based on the 'LIBLINEAR' C/C++ Library**. [S.l.: s.n.], fev. 2017.

PAWLOWSKI, Tim; DOWNWARD, Paul; RASCIUTE, Simona. Does National Pride from International Sporting Success Contribute to Well-Being? An International Investigation. **Sport Management Review**, v. 17, n. 2, p. 121–132, mai. 2014. ISSN 1441-3523. DOI: 10.1016/j.smr.2013.06.007.

PAYROVNAZIRI, Seyedeh Neelufar et al. Explainable Artificial Intelligence Models Using Real-World Electronic Health Record Data: A Systematic Scoping Review. en. **Journal of the American Medical Informatics Association**, ocaa053, 2020. DOI: 10.1093/jamia/ocaa053.

- PÉREZ-AGOTE, Alfonso. Las paradojas de la nación. spa. **REIS: Revista Española de Investigaciones Sociológicas**, Centro de Investigaciones Sociológicas (CIS), n. 61, p. 7–22, 1993. ISSN 0210-5233.
- PERKINS, Neil J.; SCHISTERMAN, Enrique F. The Inconsistency of “Optimal” Cutpoints Obtained Using Two Criteria Based on the Receiver Operating Characteristic Curve. en. **American Journal of Epidemiology**, Oxford Academic, v. 163, n. 7, p. 670–675, abr. 2006. ISSN 0002-9262. DOI: 10.1093/aje/kwj063.
- PHAM, Vu. **Bayesian Optimization for Hyperparameter Tuning**. en-US. [S.l.: s.n.], abr. 2016. Data Science.
- PIAGET, Jean. **La Psychologie de l’intelligence**. Paris: Armand Colin, set. 2012. ISBN 978-2-200-27919-6.
- PLATÃO. **A República**. Edição: Vitor Ramos. 2°. São Paulo: Difusão Européia do Livro, 1991. (Clássicos Garnier). ISBN 978-0-465-06934-7.
- PRADO JR, Caio. **Formação Do Brasil Contemporâneo**. [S.l.]: Editora Companhia das Letras, 1953.
- PRIMORATZ, Igor. Patriotism. In: ZALTA, Edward N. (Ed.). **The Stanford Encyclopedia of Philosophy**. Summer 2017. [S.l.]: Metaphysics Research Lab, Stanford University, 2017.
- PURANEN, Bi. Allegiance Eroding People’s Dwindling Willingness to Fight in Wars. In: DALTON, Russell J.; WELZEL, Christian (Ed.). **The Civic Culture Transformed: From Allegiant to Assertive Citizens**. New York, NY: Cambridge University Press, 2014. ISBN 978-1-107-03926-1 978-1-107-68272-6.
- QARI, Salmai; KONRAD, Kai A.; GEYS, Benny. Patriotism, Taxation and International Mobility. en. **Public Choice**, v. 151, n. 3, p. 695–717, jun. 2012. ISSN 1573-7101. DOI: 10.1007/s11127-011-9765-3.
- QIAN, Rui; HUNG, Juann H. Does Economic Inequality Matter for Nationalism? In: HUNG, Juann H.; CHEN, Yang (Ed.). **The State of China’s State Capitalism: Evidence of Its Successes and Pitfalls**. Singapore: Springer Singapore, 2018. p. 197–218. ISBN 9789811309830. DOI: 10.1007/978-981-13-0983-0_8.
- QUINLAN, J. Ross. **C4.5: Programs for Machine Learning**. [S.l.]: Elsevier, jun. 2014. ISBN 978-0-08-050058-4.

- QUINLAN, J. Ross. Simplifying Decision Trees. **AI Memo MIT**, n. 930, 1986.
- R CORE TEAM. **R: A Language and Environment for Statistical Computing**. Vienna, Austria: [s.n.], 2016.
- _____. **R: A Language and Environment for Statistical Computing**. Vienna, Austria: [s.n.], 2018.
- R DEVELOPMENT CORE TEAM. **R: A Language and Environment for Statistical Computing**. Vienna, Austria: [s.n.], 2004. R Foundation for Statistical Computing. ISBN 3-900051-00-3.
- RAGIN, Charles C.; AMOROSO, Lisa M. **Constructing Social Research: The Unity and Diversity of Method**. 2nd Revised edition. Los Angeles: SAGE Publications Inc, set. 2010. ISBN 978-1-4129-6018-2.
- RAMOS, Alberto Guerreiro. **O problema nacional do Brasil**. [S.l.]: Editôra Saga, 1960.
- RASHMI, K V; GILAD-BACHRACH, Ran. DART: Dropouts Meet Multiple Additive Regression Trees. en. In: **ARTIFICIAL Intelligence and Statistics**. San Diego, California, USA: [s.n.], 2015. v. 38, p. 9.
- RENNÓ, Lúcio. Teoria Da Cultura Política: Vícios e Virtudes. **BIB**, n. 45, p. 71–92, 1998.
- REPUBLICA, Instituto. **Projeto Brasilidade. Identidade Nacional e Autoestima**. [S.l.: s.n.], abr. 2010.
- RIBEIRO, Marco Tulio; SINGH, Sameer; GUESTRIN, Carlos. "Why Should I Trust You?": Explaining the Predictions of Any Classifier. **arXiv:1602.04938 [cs, stat]**, fev. 2016. arXiv: 1602.04938 [cs, stat].
- RICCI, Magda. Cabanagem, cidadania e identidade revolucionária: o problema do patriotismo na Amazônia entre 1835 e 1840. pt. **Tempo**, v. 11, n. 22, p. 5–30, 2007. ISSN 1413-7704. DOI: 10.1590/S1413-77042007000100002.
- RIPLEY, Brian. **Nnet: Feed-Forward Neural Networks and Multinomial Log-Linear Models**. [S.l.: s.n.], 2020.
- RIPLEY, Brian D.; HJORT, N. L. **Pattern Recognition and Neural Networks**. 1st. New York, NY, USA: Cambridge University Press, 1995. ISBN 978-0-521-46086-6.

- ROBILA, Mihaela; ROBILA, Stefan A. Applications of Artificial Intelligence Methodologies to Behavioral and Social Sciences. en. **Journal of Child and Family Studies**, dez. 2019. ISSN 1573-2843. DOI: 10.1007/s10826-019-01689-x.
- ROBINS, James M.; WANG, Naisyin. Inference for Imputation Estimators. In: DOI: 10.1093/biomet/87.1.113.
- ROBINSON, John P.; SHAVER, Phillip R.; WRIGHTSMAN, Lawrence S. **Measures of Personality and Social Psychological Attitudes**. [S.l.]: Gulf Professional Publishing, 1991. ISBN 978-0-12-590244-1.
- ROBNIK-SIKONJA, M.; KONONENKO, I. Explaining Classifications For Individual Instances. en. **IEEE Transactions on Knowledge and Data Engineering**, v. 20, n. 5, p. 589–600, mai. 2008. ISSN 1041-4347. DOI: 10.1109/TKDE.2007.190734.
- ROEVER, Christian et al. **klaR: Classification and Visualization**. [S.l.: s.n.], 2020.
- ROSENBLATT, F. The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain. en. **Psychological Review**, v. 65, n. 6, p. 386–408, 1958. ISSN 1939-1471, 0033-295X. DOI: 10.1037/h0042519.
- ROSSET, Saharon. **Topics in Regularization and Boosting**. 2003. PhD Dissertation – Stanford University, Department of Statistics, Stanford.
- ROUSSEAU, Jean-Jacques. **Du contrat social ou Principes du droit politique et autres écrits du contrat social**. Paris: Librairie Generale Francaise, 2003. (Le livre de Poche Classiques de la philosophie, 4644). ISBN 978-2-253-06725-2.
- ROWLAND, Robert. Patriotismo, Povo e Ódio Aos Portugueses: Notas Sobre a Construção Da Identidade Nacional No Brasil Independente. In: JANCSÓ, István (Ed.). **Brasil: Formação Do Estado e Da Nação**. São Paulo: Hucitec. [S.l.: s.n.], 2003. p. 365–388.
- RSTUDIO TEAM. **RStudio: Integrated Development Environment for r**. Boston, MA: [s.n.], 2019.
- RUBIN, Donald B. Multiple Imputations in Sample Surveys - A Phenomenological Bayesian Approach to Nonresponse. en. **Educational Testing Service**, p. 9, 1978.
- _____. Inference and Missing Data. en. **Biometrika**, v. 63, n. 3, p. 581–592, dez. 1976. ISSN 0006-3444. DOI: 10.1093/biomet/63.3.581.

RUBIN, Donald B. **Multiple Imputation for Nonresponse in Surveys**. [S.l.]: John Wiley & Sons, jun. 2004. ISBN 978-0-471-65574-9.

RUSSELL, Stuart; NORVIG, Peter. **Artificial Intelligence. A Modern Approach [Global Edition]**. 3rd. [S.l.]: Pearson, 2016. ISBN 978-1-292-15396-4.

SALANT, Priscilla; DILLMAN, Don A. **How to Conduct Your Own Survey**. New York: Wiley, 1994. ISBN 978-0-471-01273-3 978-0-471-01267-2.

SATHERLEY, Nicole et al. Differentiating between Pure Patriots and Nationalistic Patriots: A Model of National Attachment Profiles and Their Socio-Political Attitudes. **International Journal of Intercultural Relations**, v. 72, p. 13–24, set. 2019. ISSN 0147-1767. DOI: 10.1016/j.ijintrel.2019.06.005.

SCHAFER, J. L. **Analysis of Incomplete Multivariate Data**. First. [S.l.]: Chapman and Hall/CRC, 1997. (Chapman & Hall/CRC Monographs on Statistics & Applied Probability). ISBN 978-0-412-04061-0.

SCHAFER, Joseph L.; GRAHAM, John W. Missing Data: Our View of the State of the Art. en. **Psychological Methods**, v. 7, n. 2, p. 147–177, 2002. ISSN 1082-989X. DOI: 10.1037//1082-989X.7.2.147.

SCHAFER, Joseph L.; OLSEN, Maren K. Multiple Imputation for Multivariate Missing-Data Problems: A Data Analyst's Perspective. en. **Multivariate Behavioral Research**, v. 33, n. 4, p. 545–571, out. 1998. ISSN 0027-3171, 1532-7906. DOI: 10.1207/s15327906mbr3304_5.

SCHATZ, Robert T.; STAUB, Ervin; LAVINE, Howard. On the Varieties of National Attachment: Blind versus Constructive Patriotism. **Political Psychology**, v. 20, n. 1, p. 151–174, 1999. ISSN 0162-895X. DOI: 10.1111/0162-895X.00140.

SCHERER, Klaus R. What Are Emotions? And How Can They Be Measured? en. **Social Science Information**, v. 44, n. 4, p. 695–729, dez. 2005. ISSN 0539-0184, 1461-7412. DOI: 10.1177/0539018405058216.

SCHWARCZ, Lilia Moritz et al. **História Do Brasil Nação (1808-1830)**. [S.l.: s.n.], 2011. Volume I - Crise colonial.

SELIGSON, Mitchell A. The Renaissance of Political Culture or the Renaissance of the Ecological Fallacy? **Comparative Politics**, v. 34, n. 3, p. 273, abr. 2002. ISSN 00104159. DOI: 10.2307/4146954.

- SETTLES, Burr; DOW, Steven. Let's Get Together: The Formation and Success of Online Creative Collaborations. en. In: PROCEEDINGS of the SIGCHI Conference on Human Factors in Computing Systems - CHI '13. Paris, France: ACM Press, 2013. p. 2009. ISBN 978-1-4503-1899-0. DOI: 10.1145/2470654.2466266.
- SHAYO, Moses. A Model of Social Identity with an Application to Political Economy: Nation, Class, and Redistribution. en. **American Political Science Review**, v. 103, n. 02, p. 147–174, mai. 2009. ISSN 0003-0554, 1537-5943. DOI: 10.1017/S0003055409090194.
- SHEN, Haipeng. AlphaGo Seals 4-1 Victory over Go Grandmaster Lee Sedol | Technology | The Guardian. en. **The Guardian**, p. 2, 2016.
- SHMUELI, Galit. To Explain or to Predict? en. **Statistical Science**, v. 25, n. 3, p. 289–310, ago. 2010. ISSN 0883-4237. DOI: 10.1214/10-STS330.
- SIEG, Linda. Japan PM Abe's Base Aims to Restore Past Religious, Patriotic Values. en. **Reuters**, dez. 2014.
- SILVER, David et al. Mastering the Game of Go with Deep Neural Networks and Tree Search. en. **Nature**, v. 529, n. 7587, p. 484–489, jan. 2016. ISSN 1476-4687. DOI: 10.1038/nature16961.
- SIMON, Herbert A. Artificial Intelligence Systems That Understand. In: PROCEEDINGS of the 5th International Joint Conference on Artificial Intelligence - Volume 2. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1977. (IJCAI'77), p. 1059–1073.
- _____. Artificial Intelligence: An Empirical Science. en. **Artificial Intelligence**, v. 77, n. 1, p. 95–127, ago. 1995. ISSN 00043702. DOI: 10.1016/0004-3702(95)00039-H.
- SKIDMORE, Thomas E. **Brazil: Five Centuries of Change**. New York: Oxford University Press, 1999. (Latin American Histories). ISBN 978-0-19-505809-3 978-0-19-505810-9.
- SMITH, Anthony D. **Nationalism and Modernism: A Critical Survey of Recent Theories of Nations and Nationalism**. London; New York: Routledge, 1998. ISBN 978-0-203-16796-0 978-0-415-06341-8 978-0-415-06340-1.
- SMITH, Tom W.; KIM, Seokho. National Pride in Comparative Perspective: 1995/96 and 2003/04. **International Journal of Public Opinion Research**, v. 18, n. 1, p. 127–136, mar. 2006. ISSN 0954-2892. DOI: 10.1093/ijpor/edk007.

- SMITH, Tom William; JARKKO, Lars. **National Pride: A Cross-National Analysis**. [S.l.]: National Opinion Research Center, University of Chicago Chicago, IL, 1998.
- SOLT, Frederick. Diversionary Nationalism: Economic Inequality and the Formation of National Pride. **The Journal of Politics**, v. 73, n. 3, p. 821–830, 2011.
- SOROKA, George; KRAWATZEK, Félix. Nationalism, Democracy, and Memory Laws. en. **Journal of Democracy**, v. 30, n. 2, p. 157–171, 2019. ISSN 1086-3214. DOI: 10.1353/jod.2019.0032.
- SOUZA, Ewerton Pacheco de et al. Aplicações do Deep Learning para diagnóstico de doenças e identificação de insetos vetores. pt. **Saúde em Debate**, Centro Brasileiro de Estudos de Saúde, v. 43, p. 147–154, fev. 2020. ISSN 0103-1104, 0103-1104, 2358-2898. DOI: 10.1590/0103-11042019s211.
- SRIVASTAVA, Nitish et al. Dropout: A Simple Way to Prevent Neural Networks from Overfitting. en. **Journal of Machine Learning Research**, v. 15, p. 1929–1958, 2014.
- STEINRUECKEN, Christian et al. The Automatic Statistician. In: HUTTER, Frank; KOTTHOFF, Lars; VANSCHOREN, Joaquin (Ed.). **Automated Machine Learning**. Cham: Springer International Publishing, 2019. p. 161–173. ISBN 978-3-030-05317-8 978-3-030-05318-5. DOI: 10.1007/978-3-030-05318-5_9.
- STROBL, Carolin et al. Bias in Random Forest Variable Importance Measures: Illustrations, Sources and a Solution. en. **BMC Bioinformatics**, v. 8, n. 1, p. 25, jan. 2007. ISSN 1471-2105. DOI: 10.1186/1471-2105-8-25.
- ŠTRUMBELJ, Erik; KONONENKO, Igor. Explaining Prediction Models and Individual Predictions with Feature Contributions. en. **Knowledge and Information Systems**, v. 41, n. 3, p. 647–665, dez. 2014. ISSN 0219-3116. DOI: 10.1007/s10115-013-0679-x.
- SUNDHARARAJAN, Srinivasan; PAHWA, Anil; KRISHNASWAMI, Prakash. A Comparative Analysis of Genetic Algorithms and Directed Grid Search for Parametric Optimization. en. **Engineering with Computers**, v. 14, n. 3, p. 197–205, set. 1998. ISSN 1435-5663. DOI: 10.1007/BF01215973.
- SWERSKY, Kevin; SNOEK, Jasper; ADAMS, Ryan P. Multi-Task Bayesian Optimization. In: BURGESS, C. J. C. et al. (Ed.). **Advances in Neural Information Processing Systems 26**. [S.l.]: Curran Associates, Inc., 2013. p. 2004–2012.

TAJFEL, Henri (Ed.). **Differentiation Between Social Groups: Studies in the Social Psychology of Intergroup Relations**. London ; New York: Academic Pr, jan. 1979. ISBN 978-0-12-682550-3.

TAJFEL, Henri; TURNER, John C. The Social Identity Theory of Intergroup Behavior. In: JOST, John T.; SIDANIUS, Jim (Ed.). **Political Psychology**. Zeroth. [S.l.]: Psychology Press, jan. 2004. p. 276–293. ISBN 978-0-203-50598-4. DOI: 10.4324/9780203505984-16.

TEDESCHI, James T.; SCHLENKER, Barry R.; BONOMA, Thomas V. Cognitive Dissonance: Private Ratiocination or Public Spectacle? **American Psychologist**, v. 26, n. 8, p. 685–695, 1971. ISSN 1935-990X(Electronic),0003-066X(Print). DOI: 10.1037/h0032110.

THERNEAU, Terry; ATKINSON, Beth. **Rpart: Recursive Partitioning and Regression Trees**. [S.l.: s.n.], 2019.

THERNEAU, Terry; ATKINSON, Beth et al. Package 'Rpart'. **Available online: cran. ma. ic. ac. uk/web/packages/rpart/rpart. pdf (accessed on 20 April 2016)**, 2015.

TILLEY, James; HEATH, Anthony. The Decline of British National Pride. en. **The British Journal of Sociology**, v. 58, n. 4, p. 661–678, dez. 2007. ISSN 1468-4446. DOI: 10.1111/j.1468-4446.2007.00170.x.

TILLY, Charles. **Coercion, Capital, and European States, AD 990-1990**. [S.l.: s.n.], 1990.

TOCQUEVILLE, Alexis de. **A Democracia Na America**. [S.l.]: Martins Fontes, 2001.

TAI, Chih-Fong; HSU, Yu-Feng; YEN, David C. A Comparative Study of Classifier Ensembles for Bankruptcy Prediction. **Applied Soft Computing**, v. 24, p. 977–984, nov. 2014. ISSN 1568-4946. DOI: 10.1016/j.asoc.2014.08.047.

TSYRLINA-SPADY, Tatyana; LOVORN, Michael. Patriotism, History Teaching, and History Textbooks in Russia: What Was Old Is New Again. In: ZAJDA, Joseph (Ed.). **Globalisation, Ideology and Politics of Education Reforms**. Cham: Springer International Publishing, 2015. p. 41–57. ISBN 978-3-319-19505-6 978-3-319-19506-3. DOI: 10.1007/978-3-319-19506-3_4.

TURING, A. M. Computing Machinery and Intelligence. **Mind**, v. LIX, n. 236, p. 433–460, 1950.

- VAN BUUREN, S. et al. Fully Conditional Specification in Multivariate Imputation. en. **Journal of Statistical Computation and Simulation**, v. 76, n. 12, p. 1049–1064, dez. 2006. ISSN 0094-9655, 1563-5163. DOI: 10.1080/10629360600810434.
- VAN_BUUREN, Stef van. **Flexible Imputation of Missing Data**. Boca Raton: CRC Press, 2012. ISBN 978-1-4398-6825-6.
- VARNHAGEN, Francisco Adolfo de. **História Geral Do Brazil (1854-1857)**. [S.l.: s.n.], 1857.
- VERBA, Sidney. Patriotism and Nationalism: Their Psychological Foundations. By Leonard W. Doob. (New Haven: Yale University Press, 1964. Pp. 297. \$6.75.) en. **American Political Science Review**, v. 59, n. 2, p. 468–469, jun. 1965. ISSN 0003-0554, 1537-5943. DOI: 10.1017/S0003055400140821.
- VIROLI, Maurizio. **For Love of Country: An Essay on Patriotism and Nationalism**. Oxford: Clarendon Press, 1997. ISBN 978-0-19-829358-3.
- WALLACH, Hanna. Computational Social Science: Toward a Collaborative Future. In: ALVAREZ, R. Michael (Ed.). **Computational Social Science**. Cambridge: Cambridge University Press, 2016. p. 307–316. ISBN 978-1-316-25734-0. DOI: 10.1017/CBO9781316257340.014.
- WASSERSTEIN, Ronald L.; LAZAR, Nicole A. The ASA's Statement on p-Values: Context, Process, and Purpose. **The American Statistician**, v. 70, n. 2, p. 129–133, abr. 2016. ISSN 0003-1305. DOI: 10.1080/00031305.2016.1154108.
- WEISS, Sholom M.; KULIKOWSKI, Casimir A. **Computer Systems That Learn: Classification and Prediction Methods from Statistics, Neural Nets, Machine Learning, and Expert Systems**. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1991. ISBN 978-1-55860-065-2.
- WELZEL, Christian; INGLEHART, Ronald. Agency, Values, and Well-Being: A Human Development Model. en. **Social Indicators Research**, v. 97, n. 1, p. 43–63, mai. 2010. ISSN 0303-8300, 1573-0921. DOI: 10.1007/s11205-009-9557-z.
- _____. **Political Culture**. In Caramani, D. (Ed.) Oxford: Oxford University Press. [S.l.: s.n.], jan. 2010.

WICKHAM, Hadley. **Forcats: Tools for Working with Categorical Variables (Factors)**. [S.l.: s.n.], 2020.

_____. **Stringr: Simple, Consistent Wrappers for Common String Operations**. [S.l.: s.n.], 2019.

_____. **Tidyverse: Easily Install and Load the 'Tidyverse'**. [S.l.: s.n.], 2019.

WICKHAM, Hadley; BRYAN, Jennifer. **Readxl: Read Excel Files**. [S.l.: s.n.], 2019.

WICKHAM, Hadley; CHANG, Winston et al. **Ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics**. [S.l.: s.n.], 2020.

WICKHAM, Hadley; FRANÇOIS, Romain et al. **Dplyr: A Grammar of Data Manipulation**. [S.l.: s.n.], 2020.

WICKHAM, Hadley; HENRY, Lionel. **Tidyr: Tidy Messy Data**. [S.l.: s.n.], 2020.

WICKHAM, Hadley; HESTER, Jim; FRANCOIS, Romain. **Readr: Read Rectangular Text Data**. [S.l.: s.n.], 2018.

WICKHAM, Hadley; SEIDEL, Dana. **Scales: Scale Functions for Visualization**. [S.l.: s.n.], 2019.

WILSON, Woodrow. Spurious versus Real Patriotism in Education. **The School Review**, v. 7, n. 10, p. 599–620, 1899. ISSN 0036-6773. DOI: 10.1086/434094.

WIMMER, Andreas. Power and Pride: National Identity and Ethnopolitical Inequality around the World. en. **World Politics**, v. 69, n. 4, p. 605–639, out. 2017. ISSN 0043-8871, 1086-3338. DOI: 10.1017/S0043887117000120.

WITTEN, Ian H.; FRANK, Eibe. **Data Mining: Practical Machine Learning Tools and Techniques**. Second Edition. [S.l.]: Elsevier, 2011. ISBN 978-0-12-374856-0. DOI: 10.1016/C2009-0-19715-5.

WOLAK, Jennifer; DAWKINS, Ryan. The Roots of Patriotism Across Political Contexts: Roots of Patriotism Across Political Contexts. en. **Political Psychology**, set. 2016. ISSN 0162895X. DOI: 10.1111/pops.12363.

WOLPERT, David H. Stacked Generalization. **Neural Networks**, v. 5, n. 2, p. 241–259, jan. 1992. ISSN 0893-6080. DOI: 10.1016/S0893-6080(05)80023-1.

_____. The Lack of A Priori Distinctions Between Learning Algorithms. **Neural Computation**, v. 8, n. 7, p. 1341–1390, out. 1996.

- WWF-IBOPE. **Pesquisa WWF Brasil - Ibope Inteligência**. [S.l.: s.n.], 2018.
- XIE, Yihui. **Bookdown: Authoring Books and Technical Documents with r Markdown**. [S.l.: s.n.], 2020.
- _____. **Knitr: A General-Purpose Package for Dynamic Report Generation in r**. [S.l.: s.n.], 2020.
- YANG, Yezhou et al. Robot Learning Manipulation Action Plans by "Watching" Unconstrained Videos from the World Wide Web. In: PROCEEDINGS of the Twenty-Ninth AAAI Conference on Artificial Intelligence. Austin, Texas: AAAI Press, 2015. (AAAI'15), p. 3686–3692. ISBN 978-0-262-51129-2.
- YOU DEN, W. J. Index for Rating Diagnostic Tests. **Cancer**, John Wiley & Sons, Ltd, v. 3, n. 1, p. 32–35, jan. 1950. ISSN 0008-543X. DOI: 10.1002/1097-0142(1950)3:1<32::AID-CNCR2820030106>3.0.CO;2-3.
- ZALLER, John R. **The Nature and Origins of Mass Opinion**. 1st edition. Cambridge England ; New York, NY, USA: Cambridge University Press, ago. 1992. ISBN 978-0-521-40786-1.
- ZHU, Hao. **kableExtra: Construct Complex Table with 'kable' and Pipe Syntax**. [S.l.: s.n.], 2019.

ÍNDICE REMISSIVO

- aprendizado não supervisionado, 29
- aprendizado reforçado, 29
- aprendizado supervisionado, 29
- aprendizagem de máquina, 28

- big-data, 31
- boosting, 65

- ciência política, 31
- ciências sociais, 31
- Ciências Sociais Computacionais, 32
- comportamento, 31

- ensemble, 65
- European Value Survey, 54
- EVS, 54

- Floresta Aleatória, 65

- gradiente, 65

- IA forte, 25
- IA fraca, 25
- imputação múltipla, 129
- imputação única, 129
- inteligência, 24
- inteligência artificial, 15, 24

- não-respostas, 58, 62, 65, 90, 91, 121

- Pesquisa Mundial de Valores, 53

- redes neurais, 26

- sistemas especialistas, 26

- Teste de Turing, 26

- WVS, 32, 53

- árvores de decisão, 65