

# PROBLEMATIZANDO “INEFICIÊNCIAS”: REFLEXÕES NA FRONTEIRA ENTRE INTELIGÊNCIA ARTIFICIAL, MODERAÇÃO DE CONTEÚDO E DISCRIMINAÇÃO ALGORÍTMICA, EM PERSPECTIVA DECOLONIAL

## SCRUTINIZING “INEFFICIENCIES”: DELVING INTO THE CROSSROADS OF ARTIFICIAL INTELLIGENCE, CONTENT MODERATION, AND ALGORITHMIC DISCRIMINATION FROM A DECOLONIAL PERSPECTIVE

**Amanda Chami<sup>1</sup>**

Mestranda em Teoria do Estado e Direito Constitucional na PUC-Rio. Pesquisadora no grupo de Reescrita de Sentenças, organizado pela Profa. Márcia Nina Bernardes (PUC-Rio). Graduada em Direito pela PUC-Rio. Advogada no escritório de advocacia Terra, Tavares, Ferrari, Schenk, Elias Rosa. Orcid: <https://orcid.org/0009-0007-3129-8689> E-mail: [amandafchami@gmail.com](mailto:amandafchami@gmail.com)

---

**Resumo:** A atividade de moderação de conteúdo em redes sociais vem sendo majoritariamente protagonizada pela inteligência artificial. Criadas, alimentadas e treinadas por humanos, essas complexas linhas de código estão sujeitas à reprodução, na esfera digital, de vieses discriminatórios preexistentes no mundo físico. Partindo dessa premissa, este ensaio problematiza supostas “ineficiências” da inteligência artificial na atividade de moderação de conteúdo, especificamente no que tange à prática de discriminação algorítmica contra grupos vulnerabilizados. Apesar de frequentemente reconhecidos pelas próprias aplicações como “erros” da máquina, meu argumento é de que reiterados episódios de discriminação algorítmica revelam, antes, uma seletividade programada.

**Palavras-chave:** Inteligência artificial. Moderação de conteúdo. Discriminação algorítmica. Ineficiência programada.

**Abstract:** Content moderation on social media platforms has been performed mainly by artificial intelligence. Devised, fed, and trained by humans, artificial intelligence used for content moderation purposes is subject to reproducing pre-existing biases from the physical world into the digital realm. This paper aims to question alleged “shortcomings” of artificial intelligence in content moderation,

---

<sup>1</sup> Agradeço à Profa. Dra. Caitlin Mulholland pelos instigantes debates sobre o tema objeto deste ensaio e pela primeira leitura do trabalho.

namely concerning the concrete risks of algorithmic discrimination against marginalized groups. Often acknowledged by the platforms themselves as issues related to machine “inefficiency”, my argument posits that recurring episodes of algorithmic discrimination unveil, in reality, a programmed bias or selectivity.

**Keywords:** Artificial intelligence. Content moderation. Algorithmic discrimination. Programmed inefficiency.

**Sumário:** **1** Introdução: por que a perspectiva decolonial – **2** Premissas para mapear a conversa – **3** Problematizando “ineficiências”, quando a ineficiência é programada – **4** Evidências empíricas da seletividade virtual e da ineficiência programada – **5** Notas finais: por uma agenda de reforma direcionada e deliberada – Referências

---

## 1 Introdução: por que a perspectiva decolonial

Muitos hoje já se referem ao decolonial como um termo de “moda”, dessas que vêm e vão na academia, atraindo a atenção de pesquisadores e leitores enquanto o brilho da novidade se mantém. O termo “moda” não necessariamente se reveste de teor pejorativo, uma vez que, em geral, o interesse por determinados temas ou lentes teóricas, ainda que venha em ondas, se justifica pela conjuntura social vivida em dado momento. A moda pode passar, novos conceitos podem substituir o brilho dos antigos, mas isso não significa que o valor de uma determinada lente teórica tenha se desbotado. Neste ensaio, pretendo abordar a relação entre inteligência artificial, moderação de conteúdo e discriminação algorítmica a partir de uma lente decolonial.

O pensamento decolonial nos fornece ferramentas, teóricas ou empíricas, para perceber continuidades entre a forma como se estruturam as relações interpessoais e institucionais hoje e o modelo colonial, escravocrata, patriarcal que fundou nossa sociedade brasileira. Nos fornece uma lente ótima para perceber como aquele sujeito eleito pela modernidade como o sujeito-modelo-hegemônico – o homem branco, cis, hétero, cristão, capitalista e sem deficiência – permanece, ainda hoje, o detentor do capital humano, constituindo-se pela desumanização dos demais não sujeitos ou menos sujeitos que compõem a malha social.<sup>2</sup>

---

<sup>2</sup> “O padrão de normalização da condição humana eleito pela modernidade relaciona-se ao modelo de sujeito soberano de origem europeia, masculino, branco, cristão, heteronormativo, detentor dos meios de produção e sem deficiências. A aposta na universalidade para desarmar o relativismo de valores e interesses (dramatizados por conflitos sociais, políticos, econômicos, culturais, religiosos, etc.) teve como uma de suas consequências a fixação de uma lógica binária dentro da qual o universal e o relativo são mutuamente excludentes. Para além de reforçar a necessidade de proteção de determinados sujeitos e sua forma de vida, tal concepção, porque incapaz de absorver outros perfis, (reproduz hierarquizações

Algumas formas de violência e dominação características do período colonial se reproduzem mesmo nos dias de hoje, animadas pelas hierarquias fabricadas a partir das imbricações entre diversas categorias humanas – gênero, raça, classe, sexualidade e assim por diante. Há também rupturas, é claro, e, igualmente, *reinvenções* nas maneiras de reproduzir a lógica colonial na contemporaneidade.

A tecnologia e, neste caso, mais especificamente a atividade de moderação de conteúdo desenvolvida pela inteligência artificial é uma dessas ferramentas que, hoje, reinventa, produz e reproduz formas de discriminação entre (não)sujeitos. Neste escrito, portanto, pretendo tratar de alguns obstáculos atualmente enfrentados na atividade de moderação em rede sociais, quando seu objeto é o conteúdo produzido por grupos *outrificados*.

## 2 Premissas para mapear a conversa

A primeira premissa para este diálogo está no dado de que o racismo – ao lado de outras formas de violência, seja ela física, psicológica ou epistêmica, que se espriem dessa mesma matriz de poder colonial –<sup>3</sup> é estrutural e estruturante. Isso significa, em outras palavras, que:

As sociedades que vieram a constituir a chamada América Latina, foram as herdeiras históricas das ideologias de classificação social (racial e sexual) e das técnicas jurídico-administrativas das metrópoles ibéricas. Racialmente estratificadas, dispensaram formas abertas de segregação, uma vez que as hierarquias garantem a superioridade dos brancos enquanto grupo dominante [...]. Veiculada pelos meios de comunicação de massa e pelos aparelhos ideológicos tradicionais, [a ideologia do branqueamento] reproduz e perpetua a crença de que as classificações e os valores do Ocidente branco são os únicos verdadeiros e universais. Uma vez estabelecido, o mito da superioridade branca [...] é internalizado [...].<sup>4</sup>

---

entre seres humanos, saberes e cosmovisões que terão que ser sufocadas e invisibilizadas para que não ponham em risco o desenvolvimento do projeto de dominação colonial que a sustenta” (PIRES, Thula. Direitos humanos traduzidos em pretuguês. *Seminário Internacional Fazendo Gênero – 11 & 13th Women’s Worlds Congress*, Florianópolis, 2017. p. 2).

<sup>3</sup> COLLINS, Patricia Hill. *Black feminist thought: knowledge, consciousness, and the politics of empowerment*. [s.l.]: [s.n.], 2000. p. 134-135.

<sup>4</sup> GÓNZALEZ, Léila. A categoria político-cultural de amefricanidade. *Revista Tempo Brasileiro*, n. 92-93, p. 69-82, jan./jun. 1988. p. 73.

E ainda:

No Brasil e na América Latina, a violação colonial perpetrada pelos senhores brancos contra as mulheres negras e indígenas e a miscigenação daí resultante está na origem de todas as construções de nossa identidade nacional, estruturando o decantado mito da democracia racial latino-americana, que no Brasil chegou até as últimas consequências. [...] O que poderia ser considerado como história ou reminiscências do período colonial permanece, entretanto, vivo no imaginário social e adquire novos contornos e funções em uma ordem social supostamente democrática, que mantém intactas as relações de gênero segundo a cor ou a raça instituídas no período da escravidão.<sup>5</sup>

Uma segunda premissa está no consenso – foram anos de ativismo e produção acadêmica para que hoje possamos empregar aqui, com relativa confiança, o termo “consenso” – de que a inteligência artificial não é uma força neutra.<sup>6</sup> Não se trata de uma inteligência amórfica e impalpável que paira sobre um éter de ideias, com pulso e pensamento próprios. Ao menos com base no atual estado da arte, a inteligência artificial é, com o perdão da simplificação, um programa formado por linhas de código – linhas essas escritas por pessoas humanas, de carne e osso, detrás de telas de computadores. Assim é que a inteligência artificial não é neutra, o código não é neutro, uma vez que reproduzem os vieses do seu idealizador humano que, por sua vez, reproduz os vieses da sua sociedade.

Assim como o código, a linguagem (e aqui já não mais me refiro à computacional, mas à propriamente dita, que permite a comunicação entre as pessoas) tampouco é neutra. Nos usos que fazemos da linguagem, repousam disputas sobre sentidos, expressões, contexto e cultura. Linguagem é, ao fim e ao cabo, sinônimo de poder. Nós conhecemos, apreendemos e simbolizamos o mundo por meio da linguagem; imprimimos significado. O nomeado torna-se real por meio da nomeação

<sup>5</sup> CARNEIRO, Sueli. *Enegrecer o feminismo: a situação da mulher negra na América Latina a partir de uma perspectiva de gênero*. Portal Geledés, 2011. Disponível em: <https://www.geledes.org.br/enegrecer-o-feminismo-situacao-da-mulher-negra-na-america-latina-partir-de-uma-perspectiva-de-genero/>. Acesso em: 9 jan. 2022.

<sup>6</sup> Apenas a título ilustrativo: “é importante frisar que, ainda que complexos, são já conhecidos e amplamente explorados pela literatura os problemas associados à incorporação de vieses culturais e preconceitos raciais, de gênero e outros em sistemas de aprendizado por máquinas (*machine learning*), que levam a situações em que pessoas integrantes de determinados grupos sociais e étnicos sejam sistematicamente prejudicados por sistemas automatizados de decisão” (WIMMER, Miriam; DONEDA, Danilo. “Falhas de IA” e a intervenção humana em decisões automatizadas: parâmetros para a legitimação pela humanização. *Direito Público*, Brasília, v. 18, n. 100, p. 374-406, out./dez. 2021. p. 380).

e assim nossos “mundos” são construídos.<sup>7</sup> Assim, a linguagem e, sobretudo, os usos que fazemos dela, são o palco do político, um espaço de luta.<sup>8</sup>

Ainda que a internet tenha nascido sob clamores de absoluta igualdade e neutralidade, essa promessa não se concretizou. Seja pela desigualdade no seu acesso,<sup>9</sup> seja pelas regras do jogo veladamente aplicadas aos seus usuários, a tal “terra de ninguém”, na verdade, tem agendas bastante específicas, que favorecem o *ethos* hegemônico do sujeito moderno. É sobre essas regras do jogo veladamente aplicadas que passo a falar.

### 3 Problematizando “inefiências”, quando a ineficiência é programada

O campo de estudo na intersecção entre a moderação de conteúdo e teorias críticas ainda carece de desenvolvimento no Brasil. Refiro-me – vale destacar – especificamente ao estudo dos desafios da *moderação de conteúdo*, uma vez que a fronteira entre direito e tecnologia, envolvendo discriminações algorítmicas e, mais especificamente, racismo algorítmico, tem sido muito – e muito bem – explorada pela nossa literatura e ativismos<sup>10</sup> em *outras frentes*, que não necessariamente a da moderação de conteúdo.<sup>11</sup>

Dito isso, com exceção do notável trabalho desenvolvido por Carolina Bouchardet,<sup>12</sup> a maior parte dos dados *empíricos* a respeito da forma como as

<sup>7</sup> “Um homem que possui a linguagem possui, em contrapartida, o mundo que essa linguagem expressa e que lhe é implícito” (FANON, Frantz. *Pele negra, máscaras brancas*. [1952]. Salvador: EDUFBA, 2008. p. 34).

<sup>8</sup> “Antes de mais nada, é preciso esclarecer que quando utilizarmos a palavra discurso no decorrer do livro e a importância de se interromper com o regime de autorização discursiva, estamos nos referindo à noção foucaultiana de discurso. Ou seja, de não pensar discurso como amontoado de palavras ou concatenação de frases que pretendem um significado em si, mas como um sistema que estrutura determinado imaginário social, pois estaremos falando de poder e controle” (RIBEIRO, Djamilia. *O que é lugar de fala?* Belo Horizonte: Letramento, Justificando, 2017. p. 32).

<sup>9</sup> “A maior parte dos usuários de Internet brasileiros (62%) acessa a rede exclusivamente pelo celular, realidade de mais de 92 milhões de indivíduos. [...] 36 milhões de brasileiros não são usuários da rede. Esse grupo é maior entre habitantes de áreas urbanas (29 milhões); com grau de instrução até o Ensino Fundamental (29 milhões); pretos e pardos (21 milhões); das classes DE (19 milhões); e com 60 anos ou mais (18 milhões)” (92 MILHÕES de brasileiros acessam a Internet apenas pelo telefone celular, aponta TIC Domicílios 2022. *N/C.br*, 16 maio 2023. Disponível em: <https://nic.br/noticia/releases/92-milhoes-de-brasileiros-acessam-a-internet-apesar-pelo-telefone-celular-aponta-tic-domicilios-2022/>. Acesso em: 10 dez. 2023).

<sup>10</sup> Na ponte entre os dois, confira-se o trabalho desenvolvido por Nina Da Hora, disponível em: <https://www.ninadahora.dev/>; <http://lattes.cnpq.br/7768702790843282>. Acesso em: 10 dez. 2023.

<sup>11</sup> Como a proteção de dados, *cybersegurança*, policiamento com reconhecimento facial e o uso da inteligência artificial em contextos de conflito armado, apenas para citar alguns campos em que vêm se intensificando relevantes debates.

<sup>12</sup> DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuáries do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022.

hierarquias de humanidade influenciam, em concreto, as práticas de *moderação de conteúdo* foram produzidas em solo estadunidense. Apesar das sensíveis diferenças históricas na constituição dos dois países e até mesmo nos modos como discriminações estruturais se fazem presentes lá e cá, os dados levantados nas pesquisas realizadas por nossos “vizinhos do norte” chamam atenção para uma importante lacuna ainda a ser explorada empiricamente por pesquisadoras e pesquisadores no Brasil.

Em artigo publicado no *Yale Law Journal*, Anita Allen<sup>13</sup> introduz o conceito do *Black-Opticon*, que ela cunha para descrever o complexo de discriminações de que são tornadas alvo especificamente pessoas negras.<sup>14</sup> O termo é fruto de um jogo de palavras entre *black* (negro, em inglês) e o conceito de panóptico, em aceno à teorização foucaultiana sobre a invenção de Jeremy Bentham. O argumento de Anita Allen é de que o *Black-Opticon* encerra uma continuidade: ela enxerga, nos atuais mecanismos de discriminação virtual, padrões de tratamento dispensados aos grupos historicamente escravizados e marginalizados nos Estados Unidos. É o virtual encontrando novas formas de reinventar, com novo e nefasto sopro criativo, dinâmicas de poder que o mundo físico desenvolve há séculos.

Nesse sentido, a autora identifica três pilares da experiência da população negra estadunidense sob o *Black-Opticon*: a supervigilância discriminatória (*panopticon*); a exclusão discriminatória (*ban-opticon*) e a predação discriminatória (*con-opticon*). Além dos mais debatidos efeitos da supervigilância (como exemplo, o uso de ferramentas de reconhecimento facial em protestantes nas ruas), a autora traz exemplos concretos de fraudes e golpes *on-line* desenhados para atrair especificamente pessoas negras (supostos financiamentos, empregos falsos, esquemas para ganhar dinheiro), bem como a possibilidade de que anúncios pagos de acesso à moradia fossem contratados junto ao Facebook com a deliberada exclusão do disparo a pessoas negras ou pessoas que tivessem demonstrado interesse em rampas para cadeira de rodas.<sup>15</sup>

Quando se adiciona a inteligência artificial à essa equação, tem-se uma camada extra de complexidades decorrentes da falta de transparência que hoje ainda se enfrenta em relação a essas tão poderosas linhas de código. Poderosas,

<sup>13</sup> Anita L. Allen foi a primeira mulher negra a se tornar membro da American Philosophical Association. Possui formação em Direito (Harvard Law) e Filosofia (PhD pela Universidade de Michigan), especializando-se nas dimensões jurídicas e filosóficas da privacidade e proteção de dados.

<sup>14</sup> ALLEN, Anita L. Dismantling the “black opticon”: privacy, race, equity, and online data-protection reform. *Yale Law Journal*, v. 131, 20 fev. 2022. Disponível em: <https://www.yalelawjournal.org/forum/dismantling-the-black-opticon>. Acesso em: 15 out. 2023.

<sup>15</sup> ALLEN, Anita L. Dismantling the “black opticon”: privacy, race, equity, and online data-protection reform. *Yale Law Journal*, v. 131, 20 fev. 2022. Disponível em: <https://www.yalelawjournal.org/forum/dismantling-the-black-opticon>. Acesso em: 15 out. 2023.

sim, porque são capazes de decidir, em grande medida, quem fala e quem é silenciado nas redes sociais.

A inteligência artificial está sujeita àquilo que a literatura especializada chama de “Paradoxo de Moravec”. O paradoxo repousa no fato de que, por um lado, as máquinas são capazes de desempenhar determinadas tarefas com excelência sobre-humana, como processar contas matemáticas ultracomplexas em questão de milésimos de segundo, ao passo que carecem de habilidade para atividades básicas que determinada criança seria perfeitamente capaz de desempenhar, como identificar tom de voz e até mesmo andar.

Poderíamos, em tese, dividir os problemas a serem enfrentados com a inteligência artificial em dois grandes grupos: (i) ineficiências e (ii) injustiças.<sup>16</sup> Como exemplo caricato e bem-humorado de ineficiência, podemos citar o célebre caso das “câmeras inteligentes” que, programadas para seguir a bola em um jogo de futebol escocês, confundiam a esfera branco e preta com o topo calvo da cabeça do árbitro, para desilusão dos telespectadores.<sup>17</sup> Já as injustiças diriam respeito a problemas éticos oriundos da atuação da inteligência artificial, com o potencial de ferir direitos fundamentais – como seria o caso da discriminação.

No entanto, a insuficiência dessa dicotomia – *ineficiências/injustiças* – está na linha tênue e embaçada que separa uma hipótese da outra. A prática tem mostrado que grande parte das supostas ineficiências também têm origem em discriminações, ainda que os grandes provedores de aplicações não estejam – voltaremos a esse ponto adiante – abertos a admiti-lo. E nem se diga que isso é imputável ao Paradoxo de Moravec.

A suposta “ineficiência” é programada, e há ao menos três grandes razões que nos permitem traçar as suas causas a partir de discriminações humanas. As três estão intimamente relacionadas entre si: *primeiro*, a “inteligência” da inteligência artificial vem de moderadores humanos – e não apenas daqueles que escrevem as linhas de código. Especificamente no campo da moderação de conteúdo, o próprio aprendizado da inteligência artificial é alimentado por humanos. Em geral, são moderadores humanos que “treinam” a máquina para que sejam, na sequência, substituídos pela IA. Moderadores humanos produzem alguma quantidade de *training data* que a máquina então passará a replicar sozinha.<sup>18</sup> Ao fim e

<sup>16</sup> ZARSKY, Tal. The trouble with algorithmic decisions: an analytic road map to examine efficiency and fairness in automated and opaque decision making. *Science, Technology & Human Values*, v. 41, n. 1, p. 118-132, 2016. Disponível em: <https://law.haifa.ac.il/images/documents/0162243915605575.pdf>. Acesso em: 12 out. 2023.

<sup>17</sup> WIMMER, Miriam; DONEDA, Danilo. “Falhas de IA” e a intervenção humana em decisões automatizadas: parâmetros para a legitimação pela humanização. *Direito Público*, Brasília, v. 18, n. 100, p. 374-406, out./dez. 2021. p. 375.

<sup>18</sup> SIAPERA, Eugenia. AI content moderation, racism and (de)coloniality. *International Journal of Bullying Prevention*, n. 4, p. 55-66, 2022. p. 61.

ao cabo, ela estará, literalmente, reproduzindo em larguíssima escala aquilo que aprendeu dos humanos.

*Segundo* – e agora retornamos à escrita do código em si –, está na representatividade dos grupos incumbidos de escrevê-los. Pessoas racializadas, por exemplo, por muito tempo estiveram e ainda estão de fora da cúpula de tomada de decisões quanto à confecção algorítmica.<sup>19</sup> Se nos preocupamos com a diversidade de perspectivas no Judiciário, por exemplo, defendendo que ele esteja “permeado e composto por pessoas sobre as quais ele vai decidir a vida”,<sup>20</sup> o mesmo deve valer para a instância imediata de decisão sobre quem tem poder de fala no vasto e fecundo terreno das redes.

*Terceiro*, está no fato de que, muitas vezes, sistemas de inteligência artificial são alimentados de outras fontes primárias – bancos de dados preexistentes – deles importando vieses igualmente predefinidos. Se a base de dado – por exemplo, do sistema penal ou carcerário – já está impregnada de um determinado viés, essa mesma seletividade será reciclada para o novo sistema, reproduzindo-se a dinâmica de desigualdade, tal como um ciclo vicioso.<sup>21</sup> A isso soma-se o problema da transparência –<sup>22</sup> se as linhas de código, isto é, o temido “algoritmo”, são entendidas como propriedade intelectual, como determinar sua publicização?

<sup>19</sup> Confirmam-se os relatórios desenvolvidos pelo Laboratório PretaLab, disponíveis em <https://www.pretalab.com/report-quem-coda>. Acesso em: 10 dez. 2023.

<sup>20</sup> O trecho é da Professora Adriana Cruz, professora, juíza federal e primeira mulher negra a ocupar o posto de secretária-geral do Conselho Nacional de Justiça (CNJ), que, em entrevista recente, cita a Professora Jane Reis ao mencionar a importância de que o Judiciário esteja “permeado e composto por pessoas sobre as quais ele vai decidir a vida”. Disponível em: <https://www1.folha.uol.com.br/poder/2022/06/diversidade-e-crucial-para-tornar-democracia-mais-legitima-diz-cientista-politica.shtml>. Acesso em: 10 dez. 2023.

<sup>21</sup> “AI systems do not emerge in a vacuum but are already part and parcel of existing relations of power. This has been identified by research into the development and application of these systems. For instance, Ruha Benjamin (2019) has shown how the application of AI has resulted in the deeper embedding of racial codes that already permeate society. To use an example, both predictive policing and algorithms predicting recidivism have been trained on historical data and make use of demographic data, post codes, insurance rates and so on. While these systems and their designers assume that these data are neutral and represent reality in an accurate manner, in fact, they re-encode the systemic inequality that is already found in the criminal justice system. It has therefore been shown that such algorithms classify black people as more likely to offend or reoffend (Angwin et al., 2016). The use of these systems in policing, housing, risk assessments, health and so on evidently works in ways that reinforce the systemic oppression of racialised communities” (SIAPERA, Eugenia. AI content moderation, racism and (de)coloniality. *International Journal of Bullying Prevention*, n. 4, p. 55-66, 2022. p. 62).

<sup>22</sup> “Gostamos de pensar dados e números como neutros e indiscutíveis, mas, a verdade é que cada informação é resultado de um contexto, pensamentos e comportamentos em que foi construído. Nem sempre estas construções são visíveis nos resultados, os vieses embutidos se tornam o produto e fica difícil diferenciar. Por exemplo, o uso da imagem do ator negro americano, Michael B Jordan, no sistema de reconhecimento facial da segurança pública do Ceará. Como a foto dele foi parar ali? Qual a metodologia por trás desta implementação? O algoritmo que não tem transparência e que estava ajudando a determinar “culpados” foi treinado com um conjunto de dados com vieses e intenções desconhecidas. Parece que a máquina está “pensando” mas há mãos humanas por trás de todo o processo” (HORA, Nina da. Não há neutralidade, e agora IA? *Futura*. Disponível em: <https://futura.fm.org.br/conteudo/midias-educativas/artigo/nao-ha-neutralidade-e-agora-ia>. Acesso em: 12 dez. 2023).

A seguir, pretendo abordar alguns exemplos práticos de casos em que supostas “falhas” da inteligência artificial mereceram ser problematizadas.

## 4 Evidências empíricas da seletividade virtual e da ineficiência programada

Conforme tive a oportunidade de introduzir acima, a moderação de conteúdo atua, sobretudo, sobre a linguagem, seja ela verbal ou não verbal, podendo incidir igualmente sobre conteúdos gráficos (imagens, vídeos etc.) – mas que não deixam de ser uma manifestação de um usuário da rede social. No espectro do *público x privado*, as redes ocupam um lugar um tanto *sui generis*. São detidas por (gigantescas) corporações privadas, mas, ao mesmo tempo, têm se mostrado um relevante espaço de manifestação pública. Se, como falamos acima, a linguagem (quem fala e quem cala) está intimamente ligada a relações de poder, um olhar atento à forma como tem se desenvolvido a atividade de moderação das redes sociais é capaz de nos mostrar como as redes reproduzem padrões sistêmicos de silenciamento do mundo analógico. Passemos aos exemplos.

Um grupo de pesquisadores baseado nos estados de Seattle e Pittsburgh, nos Estados Unidos, identificou que postagens escritas no denominado *African American English dialect*, isto é, contendo marcas linguísticas (como gírias, expressões e formas de conjugação) características de populações negras nos Estados Unidos, tinham duas vezes mais chances de ser taxadas como ofensivas do que postagens, com o mesmo teor, mas escritas por pessoas brancas, sem essas marcas linguísticas.<sup>23</sup> A pesquisa foi desenvolvida junto a moderadores (*annotators*), especificamente sobre *tweets* que, em essência, tendem a ser postagens diretas e com número limitado de caracteres.

Vale a repetição: os moderadores eram *duas vezes* mais propensos a classificar como ofensiva uma postagem que se utilizasse de marcas linguísticas características da comunidade negra estadunidense – quando comparado a um *tweet* que transmitia a mesma mensagem, mas escrito com outras palavras, sem essas marcas linguísticas. Os pesquisadores puderam induzir que esse viés era decorrente da ausência de sensibilidade dos moderadores às diferenças linguísticas entre grupos sociais hegemônicos e minorizados e ao fato de que a ofensividade de determinados termos depende da identidade de quem o emprega (um mesmo vocábulo pode ser ofensivo se proferido por uma pessoa branca a uma pessoa negra, mas não por uma pessoa negra a outra).

<sup>23</sup> SAP, Maarten *et al.* The Risk of Racial Bias in Hate Speech Detection. *57th Annual Meeting of the Association for Computational Linguistics*, Florença, p. 1668-1678, jul./ago. 2019.

Dito de outro modo, quando lhes era fornecida a informação de que o falante era uma pessoa negra, a taxa de ofensividade do *tweet* caía drasticamente – deixava-se de tomá-lo por ofensivo. Por outro lado, se a informação sobre a raça do usuário não era fornecida, os moderadores tomavam como *parâmetro* para a avaliação do grau de ofensividade do *tweet*, ainda que inconscientemente, o registro linguístico branco-hegemônico.

A pesquisa acima foi realizada com moderadores humanos. Mas há evidências de que padrões semelhantes se repetem quando a detecção do conteúdo supostamente ofensivo é realizada por inteligência artificial – e com graves consequências.

Nos últimos anos, ativistas têm chamado a atenção para o fato de que sistemas de moderação de conteúdo de redes sociais como o Facebook têm derrubado postagens que falam sobre racismo ou denunciam a ocorrência de um episódio racista específico que o usuário tenha sofrido. Postagens que denunciam racismo são censuradas pela rede, supostamente confundidas com a prática do discurso de ódio em si. A moderação não distingue o *post* racista daquele que denuncia racismo.<sup>24</sup>

A consequência, na prática, é o literal silenciamento do usuário que fez a postagem denunciando a violência sofrida no mundo analógico. Com a derrubada arbitrária da sua manifestação, a mesma estrutura de poder opera sobre ele uma segunda vez, agora no meio digital.

Quando o usuário confronta o Facebook, recebe o retorno de que foi cometido um “erro” – *quando* recebe retorno. No entanto, e pegando emprestadas as palavras da ativista Ijeoma Oluo, “não é um erro, se ele continua se repetindo”.<sup>25</sup>

<sup>24</sup> “It’s not just black people who have their posts removed. Andy Marra, executive director of the Transgender Legal Defense & Education Fund, says allies of black people run into trouble, too. Marra’s Facebook post in late January calling on Asian Americans to protect ‘black and brown who face the brunt of white supremacy’ was removed by Facebook. Twice Marra appealed the decision to take down her Facebook post, which shared an article from a popular blog showing an Asian man throwing up ‘white power’ signs to antagonize Black Lives Matter protesters. ‘This post is expressing condemnation to anti-black racism. The post also articulates critical feedback about how other people of color – specifically those in the Asian community, including myself as an Asian person – should oppose racism in all of its forms,’ she wrote in one appeal that Facebook denied. It was only when friends reached out to Facebook to plead her case that Marra’s Facebook post was reinstated. Critics say having those kinds of connections is the only way that Facebook corrects content moderation errors, but it’s not a channel available to just anyone seeking redress” (GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023).

<sup>25</sup> “Every black person I know who has been suspended for confronting racism on Facebook has gotten the same ‘this was a mistake’ response. It is not a mistake if it keeps happening,’ Oluo said. ‘The only reason my ban was reversed was because of the outrage it generated, but so many other marginalized people in similar situations are simply forced out’” (GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023).

Os casos que ganham repercussão midiática<sup>26</sup> tendem a ser endereçados com brevidade pela aplicação. Por outro lado, quando o mesmo ocorre com pessoas “anônimas”, estas estão sujeitas até mesmo ao banimento da rede social caso haja reiteradas “ofensas”.

Por isso, muitas vezes, a via que se impõe é a do silenciamento. O usuário que já passou por essa experiência provavelmente pensará duas vezes antes de fazer uma nova postagem denunciando racismo. Afinal, as repercussões de se ter uma conta banida vão desde a impossibilidade de se comunicar com amigos e familiares por dias, semanas, até a impossibilidade de gerir negócios que dependam da rede social para existir, pondo em risco sua fonte de sustento.<sup>27</sup>

Outro estudo, este desenvolvido por pesquisadores de Michigan, EUA, se propôs a investigar quais grupos de usuários de redes sociais tinham suas postagens e suas contas derrubadas com maior frequência e qual o tipo de conteúdo frequentemente removido.<sup>28</sup> Os três grupos que tinham seu conteúdo derrubados com maior frequência eram de pessoas negras, pessoas trans e conservadores políticos. Contudo, o que a pesquisa demonstrou é que o tipo de conteúdo removido era drasticamente diferente.

O conteúdo dos grupos conservadores, quando derrubado, incluía manifestações ofensivas, desinformação, conteúdo “adulto” e discurso de ódio – isto é, todas hipóteses em que, de fato, havia violação aos termos de uso da plataforma. Já o conteúdo das pessoas trans era mais frequentemente removido quando criticavam um grupo dominante (como homens ou pessoa brancas), quando tratavam

<sup>26</sup> Por exemplo, o caso de Tanya Faison, organizadora do movimento Black Lives Matter, que, do mesmo modo, teve um de seus *posts* removidos, enquadrado como “discurso de ódio” (GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023).

<sup>27</sup> “Black activists say hate speech policies and content moderation systems formulated by a company built by and dominated by white men fail the very people Facebook claims it’s trying to protect. Not only are the voices of marginalized groups disproportionately stifled, Facebook rarely takes action on repeated reports of racial slurs, violent threats and harassment campaigns targeting black users, they say. Many of these users now think twice before posting updates on Facebook or they limit how widely their posts are shared. So to avoid being flagged, they use digital slang such as ‘wypipo,’ emojis or hashtags to elude Facebook’s computer algorithms and content moderators. They operate under aliases and maintain back-up accounts to avoid losing content and access to their community. And they’ve developed a buddy system to alert friends and followers when a fellow black activist has been sent to Facebook jail, sharing the news of the suspension and the posts that put them there. They call it getting ‘Zucked’ and black activists say these bans have serious repercussions, not just cutting people off from their friends and family for hours, days or weeks at a time, but often from the Facebook pages they operate for their small businesses and nonprofits” (GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023).

<sup>28</sup> HAIMSON, Oliver *et al.* Disproportionate removals and differing content moderation experiences for conservative, transgender, and black social media users: marginalization and moderation gray areas. *Proc. ACM Hum.-Comput. Interact.*, v. 5, n. 466, p. 466:35, 2021.

de assuntos relativos à comunidade LGBTQIA+, ou sob fundamento de se tratar de conteúdo “adulto”, a despeito de, na realidade, os pesquisadores terem constatado que eles não feriam os termos de uso da plataforma.<sup>29</sup> Por fim, o conteúdo das pessoas negras removido era relacionado à justiça racial ou antirracismo.

Em suma, apesar de grupos conservadores também serem alvo de grande quantidade de derrubada de conteúdo, quando seu conteúdo era derrubado, isso era feito de acordo com as políticas da plataforma. Quanto aos grupos de pessoas negras e trans, a derrubada era igualmente massiva, mas em absoluto descompasso com as diretrizes da própria rede social. Ineficiência, será?

É nesse sentido que Carolina Bouchardet trabalha com o conceito de “governamentabilidade algorítmica”, originalmente cunhado por Antoinette Rouvroy e Thomas Berns.<sup>30</sup> Olhando especificamente para a rede social Instagram, a autora se refere à moderação de conteúdo como uma nova forma de exercício de poder. Trata-se de *governamentabilidade*, na medida em que não apenas prevê, mas igualmente molda e condiciona condutas. Na sua pesquisa, Carolina Bouchardet parte do estudo empírico de dois casos reais no Instagram, para demonstrar como a rede “premia” as condutas, os corpos, as formas de sociabilidade que lhe são rentáveis e reprime as condutas, corpos<sup>31</sup> e formas de sociabilidade que fogem do padrão hegemônico, no limite, impedindo a sua disseminação na rede.<sup>32</sup> Nesse

<sup>29</sup> Outras pesquisas corroboram os mesmos achados: termos ligados ao universo LGBTQIA+, como “gay” ou “lésbica” estão sujeitos a ser taxados como conteúdo “adulto” ou “maduro”, o que não ocorreria do mesmo modo com termos ligados à sexualidade cis-hétero. No caso do YouTube, são considerados termos “gatilho” para o YouTube AdSense, de modo que o emprego desses termos poderia levar à perda da monetização dos vídeos produzidos. Os vídeos são desmonetizados e sofrem restrições etárias, ainda que não contenham qualquer tipo de conteúdo sexual, nudez ou linguagem imprópria. Isso gera um efeito resfriador e uma potencial invisibilização da comunidade (RODRIGUES, Gustavo; KURTZ, Lahis. *Transparência sobre moderação de conteúdo em políticas de comunidade*. Belo Horizonte: Instituto de Referência em Internet e Sociedade (IRIS), 2020. p. 67. Disponível em: <https://bit.ly/3nUbXYh>. Acesso em: 10 dez. 2023).

<sup>30</sup> ROUVROY, Antoinette; BERNIS, Thomas. Governamentalidade algorítmica e perspectivas de emancipação: o díspar como condição de individualização pela relação? Tradução de Pedro Henrique Andrade. *Revista Eco Pós, dossiê Tecnopolíticas e Vigilância*, v. 18, n. 2, p. 36-56, 2015 *apud* DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuárias do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022.

<sup>31</sup> Corpos trans são desproporcionalmente sujeitos ao controle da moderação, seja por terem suas imagens consideradas “obscenas” mesmo quando obedecem aos termos de uso, seja porque seus corpos sequer são “reconhecidos como humanos pela moderação algorítmica da rede social” (DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuárias do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022. p. 31). Também sobre o modelo estético colonial, ainda que fora da perspectiva da moderação digital, *vide* CARNEIRO, Sueli. *Enegrecer o feminismo: a situação da mulher negra na América Latina a partir de uma perspectiva de gênero*. *E-Disciplinas USP*, São Paulo, 2011.

<sup>32</sup> “No final de 2020, influenciadoras negras do Brasil realizaram experimentos com os algoritmos empregados na curadoria e recomendação de conteúdo no Instagram a fim de testar seus vieses de raça. Luana Carvalho, Sá Ollebar e Triscilla Oliveira passaram a postar imagens de mulheres brancas para verificar eventuais alterações na amplitude de distribuição do conteúdo e de engajamento de outros usuários,

processo, a autora evidencia uma série de contradições, como a discrepância entre a ferocidade com que a plataforma protege suas diretrizes antinudez (mais especificamente, antimamilos femininos), mas pouco se movimenta quando o assunto é coibir violências sofridas por usuárias mulheres na própria rede social.<sup>33</sup>

Ineficiência, *será?*

## 5 Notas finais: por uma agenda de reforma direcionada e deliberada

Entre 2017 e 2018, o Facebook cogitou implementar políticas de proteção específicas para grupos mais vulnerabilizados. No entanto, à época, o então diretor de políticas públicas do Facebook Neil Potts (hoje vice-presidente de políticas públicas da Meta Inc.) divulgou a decisão da companhia de que “todos os grupos raciais e étnicos continuariam sendo protegidos igualmente, mesmo aqueles que não sofrem opressão ou marginalização”.<sup>34</sup> As consequências disso vieram sendo sentidas nos últimos anos, como se viu acima. A problemática, naturalmente, não se resume à atual Meta. Nos Estados Unidos, também vêm sendo ajuizadas ações (*class actions*, espécie de ação coletiva, *mutatis mutandi*) em face da família Google-YouTube, em decorrência de discriminação racial<sup>35</sup> e cis-heteronormativa.<sup>36</sup>

---

verificando que tais publicações eram distribuídas a usuários da plataforma muito mais amplamente. A influenciadora digital Sá Ollebar, por exemplo, durante uma semana, publicou fotos de mulheres brancas promovendo o mesmo conteúdo de autocuidado que normalmente compartilha na plataforma. O alcance de suas publicações aumentou em 6.000% e os seguidores relataram que, apesar de não receberem suas postagens há meses, o Instagram havia circulado para eles todas as publicações da blogueira em que figuravam mulheres brancas” (DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuárias do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022. p. 46).

<sup>33</sup> “Em tese, as categorias priorizadas pela exclusão automática de conteúdo são a nudez e o discurso de ódio, mas, enquanto a moderação da rede social protege com garras as diretrizes que vedam a nudez e os mamilos femininos, falha em combater ou reconhecer sua responsabilidade no combate às violências virtuais sofridas por usuárias mulheres” (DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuárias do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022. p. 40-41).

<sup>34</sup> GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023.

<sup>35</sup> ALBERGOTTI, Reed. Black creators sue YouTube, alleging racial discrimination. *The Washington Post*, 18 jun. 2020. Disponível em: <https://www.washingtonpost.com/technology/2020/06/18/black-creators-sue-youtube-alleged-race-discrimination/>. Acesso em: 12 dez. 2023.

<sup>36</sup> KLEEMAN, Jenny. SNL producer and film-maker are latest to accuse YouTube of anti-LGBT bias. *The Guardian*, 22 nov. 2019. Disponível em: <https://www.theguardian.com/technology/2019/nov/22/youtube-lgbt-content-lawsuit-discrimination-algorithm>. Acesso em: 12 dez. 2023.

Um ponto de contato entre grande parte dos trabalhos citados é a bandeira comum por uma reforma *direcionada* e *consciente*. Em outras palavras, trata-se da constatação de que não bastam agendas neutras para combater um problema de aneutralidade (ou de uma neutralidade excludente). Se há grupos desproporcionalmente afetados pelas políticas de moderação de conteúdo no mundo *on-line*, a agenda de reforma deve igualmente se pautar por um olhar atento a essas desigualdades algorítmicas. Por isso, *direcionada*.

A transparência definitivamente é uma das pedras de toque nessa agenda de reforma, como alerta Nina da Hora.<sup>37</sup> Os algoritmos privados das redes sociais inibem qualquer possibilidade de que sejam prestadas contas da sua utilização à sociedade, de modo que, hoje, nem mesmo cientistas da computação ou desenvolvedores têm acesso às linhas de código que definem quem fala e quem cala, quem tem palco e quem é invisibilizado.

Além da transparência, contudo, é necessário que as vozes de grupos marginalizados sejam levadas a sério nesse debate e estejam representadas nas cúpulas de desenvolvimento e de decisão administrativo-institucional acerca dos algoritmos. A mera *representatividade numérica* de pessoas negras ou LGBTQIA+ ocupando cargos nessas grandes empresas – embora seja passo fundamental – não será medida suficiente caso, na estrutura institucional, não haja condições de possibilidade para a sua efetiva atuação e para o fomento de uma cultura contra-hegemônica dentro dessas instituições. O mesmo vale para a ocupação de assentos nas instâncias governamentais de controle da atividade.

Daí por que *deliberada*. É necessário que naturalizemos, dentro do nosso vocabulário tecnológico-jurídico, a necessidade da adoção de perspectivas intencional e conscientemente atentas às discriminações estruturais. A reforma eficaz virá “de dentro para fora”. Não se retira a importância de ações de natureza “reativa”, como exigir a responsabilização desses grandes *players* pelos vieses discriminatórios identificados. Mas o verdadeiro reequilíbrio da balança dependerá, ainda, de soluções “ativas” – como a aposta de Caitlin Mulholland em “ações afirmativas algorítmicas” – que nivelem as regras do jogo de dentro para fora.

## Referências

92 MILHÕES de brasileiros acessam a Internet apenas pelo telefone celular, aponta TIC Domicílios 2022. *NIC.br*, 16 maio 2023. Disponível em: <https://nic.br/noticia/releases/92-milhoes-de-brasileiros-acessam-a-internet-apenas-pelo-telefone-celular-aponta-tic-domicilios-2022/>. Acesso em: 10 dez. 2023.

<sup>37</sup> HORA, Nina da. Aula aberta: racismo e algoritmos com Nina da Hora. *YouTube*. Disponível em: <https://www.youtube.com/watch?v=7An2UVAqSW8&list=PL2dRUoriDfmvYiHyDwXN2jflyuE1i1rQf&index=14>. Acesso em: 12 dez. 2023.

ALBERGOTTI, Reed. Black creators sue YouTube, alleging racial discrimination. *The Washington Post*, 18 jun. 2020. Disponível em: <https://www.washingtonpost.com/technology/2020/06/18/black-creators-sue-youtube-alleged-race-discrimination/>. Acesso em: 12 dez. 2023.

ALLEN, Anita L. Dismantling the “black opticon”: privacy, race, equity, and online data-protection reform. *Yale Law Journal*, v. 131, 20 fev. 2022. Disponível em: <https://www.yalelawjournal.org/forum/dismantling-the-black-opticon>. Acesso em: 15 out. 2023.

CARNEIRO, Sueli. *Enegrecer o feminismo: a situação da mulher negra na América Latina a partir de uma perspectiva de gênero*. Portal Geledés, 2011. Disponível em: <https://www.geledes.org.br/enegrecer-o-feminismo-situacao-da-mulher-negra-na-america-latina-partir-de-uma-perspectiva-de-genero/>. Acesso em: 9 jan. 2022.

COLLINS, Patricia Hill. *Black feminist thought: knowledge, consciousness, and the politics of empowerment*. [s.l.]: [s.n.], 2000.

DIAS, Carolina Bouchardet. *Moderação algorítmica e autodefesa digital: a autoexposição de usuárias do Instagram como resistência ao governo algorítmico das condutas*. Dissertação (Mestrado) – Departamento de Direito, Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2022.

FANON, Frantz. *Pele negra, máscaras brancas*. [1952]. Salvador: EDUFBA, 2008.

GONZALEZ, Lélia. A categoria político-cultural de amefricanidade. *Revista Tempo Brasileiro*, n. 92-93, p. 69-82, jan./jun. 1988.

GUYNN, Jessica. Facebook while black: Users call it getting ‘Zucked,’ say talking about racism is censored as hate speech. *USA Today*, 24 abr. 2019. Disponível em: <https://www.usatoday.com/story/news/2019/04/24/facebook-while-black-zucked-users-say-they-get-blocked-racism-discussion/2859593002/>. Acesso em: 15 dez. 2023.

HAIMSON, Oliver *et al.* Disproportionate removals and differing content moderation experiences for conservative, transgender, and black social media users: marginalization and moderation gray areas. *Proc. ACM Hum.-Comput. Interact*, v. 5, n. 466, p. 466:35, 2021.

HORA, Nina da. Aula aberta: racismo e algoritmos com Nina da Hora. *YouTube*. Disponível em: <https://www.youtube.com/watch?v=7An2UVAqSW8&list=PL2dRUoriDfmvYiHyDwXN2jflyuE1i1rQf&index=14>. Acesso em: 12 dez. 2023.

HORA, Nina da. Não há neutralidade, e agora IA? *Futura*. Disponível em: <https://futura.fm.org.br/conteudo/midias-educativas/artigo/nao-ha-neutralidade-e-agora-ia>. Acesso em: 12 dez. 2023.

KLEEMAN, Jenny. SNL producer and film-maker are latest to accuse YouTube of anti-LGBT bias. *The Guardian*, 22 nov. 2019. Disponível em: <https://www.theguardian.com/technology/2019/nov/22/youtube-lgbt-content-lawsuit-discrimination-algorithm>. Acesso em: 12 dez. 2023.

PIRES, Thula. Direitos humanos traduzidos em pretuguês. *Seminário Internacional Fazendo Gênero – 11 & 13th Women’s Worlds Congress*, Florianópolis, 2017.

RIBEIRO, Djamilia. *O que é lugar de fala?* Belo Horizonte: Letramento, Justificando, 2017.

RODRIGUES, Gustavo; KURTZ, Lahis. *Transparência sobre moderação de conteúdo em políticas de comunidade*. Belo Horizonte: Instituto de Referência em Internet e Sociedade (IRIS), 2020. Disponível em: <https://bit.ly/3nUbXYh>. Acesso em: 10 dez. 2023.

SAP, Maarten *et al.* The Risk of Racial Bias in Hate Speech Detection. *57th Annual Meeting of the Association for Computational Linguistics*, Florença, p. 1668-1678, jul./ago. 2019.

SIAPER, Eugenia. AI content moderation, racism and (de)coloniality. *International Journal of Bullying Prevention*, n. 4, p. 55-66, 2022.

WIMMER, Miriam; DONEDA, Danilo. “Falhas de IA” e a intervenção humana em decisões automatizadas: parâmetros para a legitimação pela humanização. *Direito Público*, Brasília, v. 18, n. 100, p. 374-406, out./dez. 2021.

ZARSKY, Tal. The trouble with algorithmic decisions: an analytic road map to examine efficiency and fairness in automated and opaque decision making. *Science, Technology & Human Values*, v. 41, n. 1, p. 118-132, 2016. Disponível em: <https://law.haifa.ac.il/images/documents/0162243915605575.pdf>. Acesso em: 12 out. 2023.

---

Informação bibliográfica deste texto, conforme a NBR 6023:2018 da Associação Brasileira de Normas Técnicas (ABNT):

CHAMI, Amanda. Problematizando “ineficiências”: reflexões na fronteira entre inteligência artificial, moderação de conteúdo e discriminação algorítmica, em perspectiva decolonial. *Revista Brasileira de Direito Civil – RBDCivil*, Belo Horizonte, v. 33, n. 1, p. 281-296, jan./mar. 2024. DOI: 10.33242/rbdc.2024.01.012.

---

Recebido em: 05.01.2024

Aprovado em: 07.01.2024