

ARTIGO

Gênero e Discriminação Algorítmica: o papel da governança corporativa para mitigar vieses comportamentais em inteligência artificial

Débora Schwartz

debora.schwartz@alumni.usp.br

Mestranda em Direito Comercial
pela Universidade de São Paulo.
Pesquisadora bolsista no Grupo
Direito e Políticas da FDUSP. Integrante
da Comissão Especial de Direito da
Concorrência e Regulação Econômica
(CECORE), da OAB/SP. Advogada.

Gênero e Discriminação Algorítmica: o papel da governança corporativa para mitigar vieses comportamentais em inteligência artificial

Palavras-chave

Discriminação de gênero
Algoritmos
Inteligência artificial
Governança Corporativa

Resumo

O artigo pretende investigar se eventuais malefícios relacionados à discriminação algorítmica de gênero podem ser objeto de disciplina da governança corporativa. A pesquisa apresentará, primeiro, a ideia de que a discriminação algorítmica de gênero se mostra como uma nova roupagem de um problema antigo. Em seguida, argumenta-se que o Direito, e em especial o Direito Empresarial, deve detectar os problemas que essas novas tecnologias acarretam, mediante análise interdisciplinar da interface com a ciência da computação. Por último, levantam-se sugestões práticas de iniciativas, medidas e ferramentas de governança corporativa que podem ser utilizadas para mitigar a discriminação de gênero presente em algoritmos. Ao final, apresenta-se conclusão e breves perspectivas de agendas de pesquisa no campo aqui analisado.

Gender and algorithmic discrimination: the role of corporate governance in mitigating behavioural biases in artificial intelligence

Keywords

Gender discrimination
Algorithms
Artificial Intelligence
Corporate Governance

Abstract

The article aims to investigate whether possible harm related to algorithmic gender discrimination can be the subject of corporate governance discipline. The research will first present the idea that algorithmic gender discrimination shows itself as a new vest for an old problem. It then argues that law, and especially business law, must detect the problems that these new technologies bring, by analysing the interdisciplinary interface with computer science. Finally, practical suggestions are made for initiatives, measures and corporate governance tools that can be used to mitigate the gender discrimination present in algorithms. Finally, there is a conclusion and brief prospects for research agendas in the field analysed.

1 Introdução

No primeiro semestre de 2023, o Instituto InternetLab publicou uma pesquisa que examinava como os aplicativos de música Deezer e Spotify faziam recomendações de artistas aos usuários das plataformas, levando em conta o marcador de gênero na análise¹. Ao se escutar uma música nessas plataformas, o algoritmo automaticamente sugere uma próxima música ao usuário. Com base nesse movimento de recomendações, o resultado da pesquisa indicou que havia diferença significativa na quantidade de artistas mulheres e homens sendo recomendados, no geral, em ambos os aplicativos. Em um caso extremo, chegou-se à marca de uma sequência de 45 músicas recomendadas em que os intérpretes eram, exclusivamente, homens. Por outro lado, quando esse mesmo método foi testado com músicas de intérpretes femininas, o máximo de recomendações seguidas que se chegou de músicas de artistas mulheres foi de 17.

Esse estudo ilustra como as tecnologias de inteligência artificial, movidas por algoritmos, não são tão neutras quanto parecem. Em verdade, sua neutralidade é artificial (Mulholland, 2022). Elas aparentam revelar vieses comportamentais discriminatórios já há tempos identificados na conduta humana. Assim, em alguns casos, a inteligência artificial (IA) poderia reforçar problemas de discriminação de gênero há tanto discutidos, principalmente se considerarmos que suas escolhas virtuais e algorítmicas impactam diretamente na realidade de empresas e consumidores (Van Giffen et al., 2022, p. 93–106).

O conceito de IA é recorrentemente debatido no campo acadêmico e político (Wang, 2019). Segundo estudo recente da União Europeia (2019b, p. 6), pode-se definir IA como:

sistemas de software (e eventualmente também de hardware) concebidos por seres humanos, que, tendo recebido um objetivo complexo, atuam na dimensão física ou digital percebendo o seu ambiente mediante a aquisição de dados, interpretando os dados estruturados ou não estruturados recolhidos, raciocinando sobre o conhecimento ou processando as informações resultantes desses dados e decidindo as melhores ações a adotar para atingir o objetivo estabelecido.

A OCDE, por sua vez, adota uma definição em seu sítio eletrônico que está constantemente se alterando, vez que recebe aportes e sugestões de mudanças em tempo real dada a volatilidade do debate². A proposta mais atualizada da definição de sistemas de IA se aproxima daquela indicada pela UE, mas dá enfoque ao potencial de influência que esses mecanismos possuem, bem como na possibilidade de seus objetivos serem explícitos ou implícitos. Embora não haja um consenso unânime sobre qual deve ser a definição de IA mais adequada, argumenta-se que sua conceitualização é passo importante para se determinar uma moldura regulatória mínima a respeito dessas novas tecnologias (Wang, 2019).

O tema da discriminação algorítmica por sistemas de IA vem tomando força há alguns anos, especialmente por conta de casos em que os setores de recursos humanos de diversas empresas utilizavam IA para contratação de novos funcionários, ou contratavam os serviços de sites de vagas que fizessem essa seleção também por meio de inteligência artificial (Mulholland, 2022). Identificou-se, no entanto, que essas tecnologias acarretavam discriminação de gênero quando da seleção de currículos de possíveis candidatos às vagas disponíveis nas respectivas empresas contratantes.

Entre 2014 e 2015, a empresa Amazon desenvolveu um sistema de inteligência artificial

para ajudar a revisar currículos de candidatos e classificá-los com base em sua adequação para diferentes posições. No entanto, descobriu-se que o sistema tinha uma tendência a discriminar candidatas mulheres. Isso ocorreu porque ele teria sido treinado com currículos enviados à Amazon ao longo de uma década, cujos candidatos eram, em sua maioria, homens. Como resultado, o sistema aprendeu a associar características masculinas às melhores classificações³.

Já em 2018, a plataforma de recrutamento Gupy também apresentou problema similar, em que se questionou a transparência dos critérios de seleção da IA utilizada. Descobriu-se que, apesar de não ser programado para julgar o gênero dos possíveis candidatos ou candidatas, o algoritmo da Gupy não aprovava mulheres em suas seleções⁴.

Em relatório recente denominado *Governing AI for Humanity*, a ONU (2024) reconheceu que sistemas de inteligência artificial podem ser discriminatórios em relação à raça e ao gênero. Dessa forma, classificou tal possibilidade como risco latente a grupos específicos e potencialmente vulneráveis quando do uso de IA.

Tendo em vista esse risco, é corrente o debate sobre a necessidade de regulação das novas tecnologias de inteligência artificial. Em abril de 2021, a Comissão Europeia propôs o primeiro quadro regulatório da União Europeia (UE) para a IA. Ele estabelecia que sistemas de IA que pudessem ser usados em diferentes aplicações fossem analisados e classificados de acordo com o risco que representam aos usuários. Os diferentes níveis de risco significariam maior ou menor regulamentação⁵.

Mais recentemente, em junho de 2023, a UE anunciou que pretende elaborar legislação específica para regular inteligência artificial (também chamada de *Artificial Intelligence Act* ou *AI Act*). A ideia é que o documento estabeleça regras para provedores e usuários, a depender do nível de risco da inteligência artificial. Sistemas de IA que apresentassem risco mínimo teriam

regras e avaliação mais brandas. Já para aqueles sistemas que aparentassem riscos inaceitáveis, como manipulação cognitivo-comportamental de pessoas, scores sociais e identificação biométrica em tempo real, seriam proibidos⁶.

No Brasil, até o presente momento, tem-se notícia da existência de um Projeto de Lei (PL) voltado a esse tema, proposto pelo senador Rodrigo Pacheco (PSD-MG). O PL nº 2.338/2023 estabelece diretrizes para a disponibilização de sistemas de inteligência artificial no Brasil, com especial atenção àqueles impactados pelo seu uso. Além disso, determina critérios para a utilização desses sistemas pelo governo – incluindo sanções para possíveis infrações – e concede ao Poder Executivo a autoridade de designar a entidade responsável por supervisionar e regulamentar o setor⁷.

O debate, portanto, centra-se na perspectiva de atuação do Poder Público face aos eventuais malefícios que as tecnologias envolvendo algoritmos poderiam acarretar, inclusive aqueles relacionados à discriminação de gênero. Mas e as empresas? Que papel devem desempenhar? Devem tomar atitudes para, elas próprias, mitigarem os riscos apresentados por IAs? Como o setor privado se insere nessa questão?

A igualdade de gênero é um direito humano, reconhecido pela ONU⁸, e um direito fundamental, agasalhado pelo art. 5º Constituição de 1988⁹. As empresas, por sua vez, devem possuir papel ativo para com a redução das desigualdades de gênero. Não se trata de preocupação só do Estado ou da sociedade. Desde 2011, os Princípios Orientadores da ONU para Empresas e Direitos Humanos ressaltam essa necessidade de atuação constante e vigilante do setor corporativo (ONU, 2011).

Ademais, com o crescimento do debate envolvendo a pauta ESG (sigla para *Environmental, Social and Governance*), o tema tem se capitalizado e entrado nas estruturas de governança interna das empresas, como será visto adiante. Entende-se, portanto, que as empresas

possuem uma responsabilidade social e devem adquirir posição ativa diante dos problemas que a discriminação algorítmica de gênero pode acarretar. Trata-se de uma questão de governança corporativa.

Por fim, cabe ao ramo do Direito Empresarial endereçar as questões trazidas pelo tema da discriminação algorítmica de gênero. Conforme aponta recente revisão sistemática de literatura conduzida por Sainz et al. (2024, p. 271), o debate aqui exposto tem se concentrado nas áreas do Direito Administrativo e Direitos Humanos. Não há sequer notícia de artigo que trate do tema sob o prisma comercialista. O presente trabalho, portanto, objetiva suprir essa lacuna na literatura e dar início a essa agenda de pesquisa para que mais investigadores se debruçam futuramente sobre o tema.

Para isso, a metodologia adotada foi de levantamento não sistemático da literatura disponível sobre o tema, tendo em vista a escassez de trabalhos já apontada pelo levantamento bibliográfico preliminar no ramo do Direito ora destacado. Após, houve a escolha do referencial teórico para a análise sobre como a discriminação de gênero ocorre ao longo das etapas de aprendizado dos sistemas de IA, sendo elegido o trabalho de Van Giffen et al. (2022), por sua profundidade no tema. Por fim, foram levantados alguns casos de vieses algorítmicos de gênero e ferramentas para mitigá-los com base nas indicações feitas pela própria doutrina selecionada.

Tendo esses pontos como premissa, o presente artigo busca, primeiro, demonstrar que a discriminação algorítmica de gênero se apresenta como uma nova roupagem de um problema antigo. Em seguida, argumenta-se que o Direito Empresarial deve detectar e endereçar os problemas que essas novas tecnologias acarretam, mediante análise interdisciplinar das especificidades trazidas pela interface com a ciência da computação. O último capítulo deste trabalho, por sua vez, busca trazer sugestões

práticas de iniciativas, medidas e ferramentas de governança corporativa para tentar mitigar a discriminação de gênero presente em algoritmos. Ao final, apresenta-se conclusão e breves perspectivas de agendas de pesquisa no campo aqui analisado.

2 Uma velha conhecida em novas roupagens: discriminação algorítmica nas novas tecnologias

Falar em novas tecnologias, por vezes, implica olhar para novos problemas com os quais o Direito deva se preocupar. Embora, à primeira vista, a discriminação algorítmica se apresente como um tema novo no debate jurídico atual, ao se analisar mais detidamente, percebe-se que suas implicações e origens são há muito estudadas pela literatura. Exemplo disso é a falta de representação feminina no campo das ciências exatas, biológicas e humanas, fato que gerou distorções nas pesquisas até então empreendidas por setores que são, majoritariamente, masculinos (Latorre Ruiz & Pérez Sedeño, 2023, p. 66).

Assim, no âmbito algorítmico, essa disparidade social estaria escalada. Segundo revisão de literatura conduzida por Sainz et al. (2024, p. 274), restou pacífico na doutrina que a discriminação por meio de algoritmos reproduz discriminações sociais pré-existentes.

Os crescentes debates sobre discriminação algorítmica serviram, em verdade, para revelar como as desigualdades de gênero, raça etc se mostram e se perpetuam de forma quase “inconsciente” nos vieses comportamentais da sociedade. A particularidade, aqui, evidencia-se no discurso.

Na origem do desenvolvimento de sistemas de IA, difundiu-se a falácia de que tais tecnologias seriam neutras e, portanto, supostamente blindadas dos vieses humanos (Latorre Ruiz & Pérez Sedeño, 2023, p. 65). Vendeu-se a ideia de tecnologia com o intuito de substituição da ação humana pela ação da máquina, justificada pelo fato de a IA ser, em teoria, mais eficiente e livre de subjetividade (Tonucci & Caldeira, 2023). Todavia, isso não se verifica empiricamente (Mulholland, 2022; Lindoso, 2021).

Eroğlu & Karatepe Kaya (2022) destacam, por exemplo, que se cogitou utilizar IAs nos Conselhos de Administração de empresas, para mitigar distorções causadas pela falta de representatividade feminina nesses órgãos corporativos. Resultado: infrutífero. A suposta neutralidade da inteligência artificial não foi comprovada empiricamente e suas decisões, portanto, eram incapazes de substituir a perspectiva feminina – objetivo principal das políticas de cotas.

Esse tema foi mais aprofundado no Brasil sob o recorte da discriminação racial através de algoritmos. Novamente, na revisão de literatura elaborada por Sainz et al (2024, p. 268), a intersecção com o debate racial esteve presente em 69% dos estudos levantados. Há que se destacar também a obra de Silva (2022), a qual procura traçar um panorama sobre como o racismo se mescla e se interpenetra nas redes digitais.

O autor traz o exemplo do *chatbox* Tay, criado em 2016 pela Microsoft (Silva, 2022, p. 67). A inteligência artificial havia sido criada para interagir com o público jovem, usuário da plataforma Twitter, e tinha o objetivo de popularizar os *chatbots*. Após algumas horas de sua implementação, notou-se que a IA passou a reproduzir comentários xenófobos e racistas. Sucede que a base para a coleta e o processamento dos dados que a alimentavam era, precisamente, a interação com humanos através da rede social.

Lobacheva & Kashtanova (2022, p. 99) apontam para um caso nos EUA que possuía

facetas de discriminação ainda mais complexas. Tratava-se de um algoritmo usado em hospitais do país para determinar quais pacientes possivelmente demandariam cuidados médicos adicionais. Em 2019, foi descoberto que esse sistema de IA dava preferência para pacientes brancos em detrimento de pacientes negros. A informação sobre raça não era coletada pela IA, porém era utilizado como variável o valor dispendido pelo paciente em assistência médica. Por uma série de motivos, a parcela negra da população estadunidense aloca menos recursos em saúde do que a parcela branca, o que gerou uma discriminação indireta dos usuários.

Esses exemplos demonstram que os vieses discriminatórios podem influenciar as decisões dos algoritmos, ainda que critérios como raça ou gênero não sejam utilizados para treinamento da IA. Em verdade, os casos narrados demonstram que é justamente na ausência desses dados onde a discriminação tem mais chance de ocorrer indiretamente.

Nesse sentido, entende-se que os aportes feitos sobre a discriminação algorítmica em matéria de raça são aplicáveis e transponíveis para o debate de gênero. Embora o escopo da presente pesquisa não seja a questão racial, este trabalho utiliza esse referencial teórico para embasar a literatura que se debruça sobre discriminação e vieses em sistemas de IA.

Rechaçada a teoria da neutralidade dos sistemas de IA, resta clara a possibilidade de haver discriminação de grupos sociais no processo de aprendizado dessas tecnologias. O conceito de discriminação algorítmica ainda é amplamente debatido na doutrina (Sainz et al, 2024). Todavia, pode ser entendido como a capacidade da IA em replicar ou reforçar “preconceitos existentes na sociedade, ocasionando distinções, preferências ou exclusões capazes de afetar a igualdade de tratamento entre os indivíduos, sobretudo os grupos vulneráveis” (Requião & Costa, 2022, p. 4).

Algumas autoras também argumentam que o próprio conceito de discriminação, por si só, é de difícil determinação (Lobacheva & Kashtanova, 2022). Assim, interpretar a ideia de discriminação algorítmica também se apresenta como desafio. Por essa razão, a literatura tem buscado meios de aprofundar a conceitualização desse termo. Mendes e Mattiuzo (2019), por exemplo, entendem que a discriminação algorítmica pode ocorrer em quatro cenários: (i) por erro estatístico, (ii) por generalização, (iii) por uso de informações sensíveis ou (iv) como limitadora de direitos. Cada uma dessas situações pode ocorrer a depender da etapa de aprendizagem algorítmica que se observa, conforme aprofundado adiante. Pelo exposto, entende-se que a questão da discriminação algorítmica nada mais é do que uma nova roupagem para um problema antigo. No entanto, esse cenário traz novas implicações éticas e políticas, as quais o Direito deve endereçar. Mas, para entender como tratar essas questões, deve-se compreender as particularidades dos sistemas algorítmicos de IA para identificar como a discriminação se mostra nessas novas tecnologias.

3 Em qual momento do desenvolvimento algorítmico ocorre a discriminação?

Para conseguir mitigar o problema da discriminação algorítmica, é necessário compreender, em primeiro lugar, as etapas do chamado *machine learning*. Esse processo consiste, basicamente, na forma pela qual a inteligência artificial recebe, processa e transforma dados em decisões. A União Europeia descreve esse tipo de aprendizagem como técnicas que “permitem que um sistema de IA aprenda a resolver problemas que não podem ser especificados de forma precisa, ou cujo método de resolução não pode ser descrito por regras de raciocínio simbólico” (União Europeia, 2019b, p. 3).

O *machine learning* pode ser dividido em três modalidades: de aprendizagem supervisionada, aprendizagem não supervisionada e aprendizagem por reforço. No primeiro tipo, o algoritmo é treinado a partir de dados fornecidos por um operador humano, de modo que suas decisões serão pautadas com base nos exemplos recebidos. O segundo tipo de aprendizagem decodifica padrões e características comuns aos dados fornecidos e, portanto, não há exemplo correto a ser seguido. Por fim, o terceiro tipo de aprendizagem se desenvolve com base na experiência, pela lógica de tentativa e erro.

A literatura aponta que os problemas de discriminação de gênero são identificados, normalmente, em algoritmos com aprendizagem supervisionada, justamente por haver uma noção de “certo e errado” uma vez que o operador insere os exemplos indicando um viés. Nesse sentido, Van Giffen et al. (2022, p. 96) propõem uma revisão da literatura para compreender como a discriminação pode ocorrer em cada uma das etapas do *machine learning*.

esse fenômeno de *feedback loop* (O’Neil, 2018), já que, em determinado ponto, não é mais possível descobrir em que momento da cadeia algorítmica foi gerada a discriminação (Van Giffen et al., 2022, p. 98).

Essa ampla gama de vieses inconscientes passíveis de acarretar discriminação fez com que a União Europeia, em 2019, desenvolvesse um Guia de Orientações Éticas para uma IA de confiança (no título em inglês, *Ethics Guidelines for Trustworthy AI*), que estabelece sete requisitos para que se tenha um grau razoável de confiabilidade em um algoritmo: (i) ação e supervisão humanas; (ii) solidez técnica e segurança; (iii) privacidade e governação dos dados; (iv) transparência; (v) diversidade, não discriminação e equidade; (vi) bem-estar ambiental e societal; (vii) responsabilização (União Europeia, 2019a, p. 6). Especificamente para o presente artigo, destacam-se os requisitos (iv) e (v).

Sobre diversidade, não discriminação e equidade, o Guia pontua que a utilização de sistemas de IA deve vir acompanhada (i) da Prevenção de enviesamentos injustos, (ii) da noção de acessibilidade e conceção universal e (iii) da participação das partes interessadas. São diretrizes genéricas que resumem as preocupações de Van Giffen et al. (2022) descritas acima. O problema da discriminação algorítmica é, portanto, completamente reconhecido pela União Europeia.

A questão da transparência dos algoritmos – requisito (iv) – se dá pelo fato de que, com o avanço dessa tecnologia, os algoritmos têm ficado cada vez mais complexos e sofisticados. Sua documentação, portanto, é dificultada, bem como o seu acesso por terceiros. Tal situação fez com que as decisões tomadas por algoritmos ficassem opacas, não apenas para a sociedade e as autoridades, mas também para os próprios desenvolvedores de algoritmos (Latorre Ruiz & Pérez Sedeño, 2023, p. 71).

Ademais, é do interesse das empresas que esses algoritmos permaneçam na forma de

segredos de negócio, pois consistem em vantagem competitiva perante seus concorrentes. Assim, ainda que existentes, tais documentos não são comumente divulgados, sob a justificativa de se evitar a ocorrência de “efeito carona” (no jargão em inglês, *free-riding*). Pasquale (2015) explica que esse âmbito de opacidade pode ser chamado de sigilo legal, isto é, quando o conteúdo do algoritmo não pode ser revelado por determinações normativas.

O Guia da UE, por sua vez, propõe três pilares para se endereçar o problema da transparência: (i) rastreabilidade; (ii) explicabilidade e (iii) comunicação. No primeiro, há grande foco na necessidade de documentação dos algoritmos. A explicabilidade técnica, por sua vez, exige que as decisões tomadas por um sistema de IA possam ser compreendidas e rastreadas por seres humanos. Por fim, o pilar da comunicação estipula que não se deve entender os sistemas de IA como se humanos fossem, de forma que as pessoas têm direito de serem informadas de que estão interagindo com uma IA.

O Guia, contudo, não apresentou iniciativas práticas para se implementar os requisitos de uma IA de confiança (Hickman & Petrin, 2021). A tentativa da doutrina, desde então, foi buscar estabelecer mecanismos e instrumentos que pudessem endereçar os problemas discriminatórios ocasionados pelo uso de inteligência artificial. Acredita-se que a governança corporativa, por sua vez, possa ter um papel relevante na mitigação dessas questões.

4 Mecanismos de governança corporativa para mitigar a discriminação algorítmica de gênero

A partir do cenário até aqui narrado, doutrina, empresas, organizações internacionais e governos vêm tentando buscar alternativas para amenizar ou solucionar a discriminação identificada em tecnologias de inteligência artificial. Um exemplo recente é o kit de ferramentas da IBM chamado *AI Fairness 360*, que procura facilitar adequação dos algoritmos de uso industrial aos princípios éticos destacados nesta pesquisa, bem como para auxiliar pesquisadores a compartilhar e avaliarem algoritmos entre si (Bellamy et al., 2021). Assim como ela, outras ferramentas foram desenvolvidas por grandes empresas de tecnologias para tentar amenizar os problemas aqui discutidos, como se verá adiante.

Tais medidas nada mais são do que práticas de governança corporativa elaboradas pelo setor privado para endereçar o problema dos vieses em IA, também chamada nesse contexto de governança algorítmica (Mendes & Mattiuzzo, 2019). A ideia principal é que, ao projetar sistemas algorítmicos, devem ser seguidos os princípios de conscientização, explicabilidade, precisão, auditabilidade, justiça (*fairness*) e fiscalização e reparação de eventuais danos causados pela IA, sob a perspectiva do chamado *accountability* (Mendes & Mattiuzzo, 2019, p. 56).

Preconiza-se que as empresas devem ser responsabilizadas pelo fato de que as pessoas serão afetadas pelo processo decisório dos sistemas de IA. Como forma de mitigar eventuais efeitos negativos, recorre-se também à noção de transparência, justamente como forma de aclarar as *black boxes* mencionadas por Pasquale (2015), que consistiriam no fenômeno da opacidade algorítmica.

Todavia, a mera listagem de tais princípios não é suficiente para garantir seu efetivo *enforcement*. É necessário pensar medidas mais concretas e aplicáveis à prática do setor privado.

Nessa linha, acerca das possíveis discriminações algorítmicas descritas no capítulo anterior, Van Giffen et al. (2022) propõem um modelo com diversas sugestões para se mitigar os efeitos que vieses comportamentais podem causar no âmbito do *machine learning*, aqui entendidas como práticas de governança corporativa. Destacam-se três delas: (i) documentação e transparência do algoritmo; (ii) supervisão humana obrigatória e (iii) times diversos e multidisciplinares.

Cumprido desde logo destacar que tais sugestões em muito se aproximam dos princípios destacados pelo Guia da União Europeia para se atingir uma “IA confiável”. O elemento adicional aqui consiste no aspecto prático dessas possíveis soluções. Frisa-se, também, que essas alternativas para se lidar com a discriminação algorítmica devem incidir principalmente na etapa de preparação de dados do *machine learning* – etapa (iii) descrita no capítulo anterior –, onde mais se encontra o risco de incidência de vieses.

Quanto à solução de documentação e transparência do algoritmo, novas ferramentas vêm sendo desenvolvidas e testadas pelos cientistas da computação, cabendo destaque para duas: *a Model Cards*, da Google, e *AI System Cards*, da Meta. A tecnologia de *Model Cards* busca simplificar e automatizar a geração de documentos de *machine learning* que oferecem contexto e transparência para o desenvolvimento e o desempenho de um modelo de algoritmo¹⁰. Ela seria integrada ao processo de *machine learning*, de forma a possibilitar que os criadores de um sistema de IA compartilhem as métricas e os metadados do modelo com demais pesquisadores, desenvolvedores, informantes etc.

A *AI System Cards*, por sua vez, tem uma lógica semelhante. A diferença residiria no fato de que

esse modelo foca, principalmente, em informações de nível de sistema. Um System Card forneceria uma visão geral de vários modelos de aprendizado de máquina que compõem um sistema de aprendizado de máquina, bem como detalhes sobre esses componentes e um guia passo-a-passo com um exemplo de entrada¹¹.

Tais ferramentas podem ser úteis não apenas na fase de preparação dos dados, mas também na de entendimento desses dados e de implementação da IA – respectivamente etapas (ii) e (vii) mencionadas por Van Giffen et al. (2022). Uma vez ampliada a transparência no *machine learning*, é possível haver maior treinamento das equipes que o alimentam e que dele se utilizam para a criação de novos sistemas de IA. Assim, o risco de haver transmissão de vieses comportamentais entre diferentes IAs fica reduzido.

Contudo, ressalta-se que as empresas podem encontrar dificuldades na implementação dessa prática de transparência, tendo em vista a vantagem competitiva que reveste a opacidade algorítmica. Uma vez agasalhado pelo segredo de negócio, o algoritmo se transforma em um ativo a ser protegido do chamado “efeito carona”. Daí a necessidade de medidas regulatórias, capitaneadas pelo Poder Público, para que a disciplina do segredo empresarial não seja um empecilho para se atingir um grau adequado de confiabilidade e transparência algorítmicas.

Quanto à segunda sugestão, sobre a possibilidade de haver supervisão humana obrigatória dos sistemas de IA, tem sido frequente o debate sobre a criação de equipes de *Responsible AI* ou de *Explainable AI*, mais conhecidas pela sigla XAI (Arrieta et al., 2020). A fim de se garantir que haja redução nas discriminações que seus respectivos algoritmos poderiam causar, propõe-se a existência de equipes responsáveis por aplicar os princípios da transparência e da explicabilidade nos sistemas de IA utilizados no bojo de suas respectivas empresas.

Destaca-se que a supervisão humana também pode ser de extrema valia na etapa (v) do

machine learning, de avaliação do modelo, sugerida por Van Giffen et al. (2022). Isso porque, ao passar pelo crivo humano, há maiores chances de que eventual incorporação de vieses pela IA seja identificada e corrigida. Todavia, a supervisão acrítica não basta, razão pela qual se destaca a importância da criação de equipes diversas e inclusivas para realizar essa avaliação da modelagem algorítmica.

Chega-se, dessa forma, à sugestão de estabelecimento de times multidisciplinares e diversos para análise e supervisão da IA, mencionados por Van Giffen et al. (2022) e outros expoentes da literatura (Tonucci & Caldeira, 2023, p. 336). No bojo dessas equipes, o cientista da computação estaria imbuído de coletar evidências suficientes para se testar o modelo e garantir a correta documentação do algoritmo. Ao operador do Direito, caberia examinar se os resultados e o comportamento do algoritmo implicam discriminação ou ilegalidade.

Ademais, o problema da sub-representação feminina no bojo das ciências exatas, especialmente na ciência da computação e na ciência de dados, é fenômeno latente. Como apontado por Tonucci & Caldeira (2023, p. 330), estudos realizados em 2020 indicaram que 15% das cientistas de dados do mundo eram mulheres e, das profissionais que ocupavam funções de análise de dados, eram apenas 26%.

A respeito dos vieses inconscientes envolvidos na contratação de profissionais, Pisanelli (2022) desenvolveu interessante trabalho acerca da possibilidade de reduzir as disparidades de gênero usando, precisamente, inteligência artificial.

Como largamente demonstrado no início do artigo, a grande parte das notícias sobre o tema apontam para casos de aumento da discriminação de gênero, e não de redução. Todavia, o estudo empírico ora mencionado aponta para outras modelagens possíveis do *machine learning*.

Em suma, a autora explica que os problemas principais da discriminação de IAs como da

Gupy, da Amazon etc residiam no fato de que essas empresas faziam o *machine learning* com base em seus próprios funcionários, adotando uma estratégia de IA preditiva (*predictive AI*). Todavia, caso fosse adotada uma estratégia de triagem através da inteligência artificial (*screening AI*), isto é, que o sistema tivesse como base os próprios currículos e escolhesse os melhores a partir tão somente da amostra dos candidatos, a discriminação de gênero era mitigada.

O trabalho de Pisanelli (2022), portanto, indica que a adoção de modelos de IA baseados em triagem pode ser mais benéfica quando da contratação de profissionais, inclusive aumentando a probabilidade de haver contratadas mulheres.

Todavia, ainda que se transponha a barreira da contratação, compor equipes diversas e multidisciplinares não é uma tarefa simples, tampouco óbvia. Isso porque a própria integração dessas equipes é um desafio enfrentado pelas empresas quando da implementação dessa medida de governança.

Uma política comumente vista nesse âmbito é a criação de cotas para grupos minorizados¹² em processos seletivos de contratação ou em altos cargos de gestão, de modo a permitir que eles acessem espaços de poder. Esse é, definitivamente, um primeiro passo importante. Todavia, essa medida deve vir acompanhada de mecanismo de ascensão e permanência das mulheres em tais ambientes majoritariamente masculinos.

Abreu (2024), ao estudar as políticas de cotas de gênero nas companhias brasileiras, identifica que ações afirmativas de reservas de vagas a mulheres em altos cargos de gestão não são suficientes para endereçar o problema da sub-representação feminina nesses ambientes¹³. Primeiro, porque a grande maioria das empresas possui, no máximo, 30% de cargos de diretoria ocupados por mulheres, o que, segundo a autora, seria insuficiente para se falar em efetiva representação. Segundo, porque se identificou grande número de diretoras que atuam em

vários conselhos de administração – ou seja, acumulam cargos como membros independentes, fenômeno também conhecido como *trophy directors* ou *golden skirts*. São mulheres que possuem dupla ou tripla jornada para garantir a representação feminina nos cargos decisórios das empresas. Tendo em vista o elevado comprometimento exigido nesses cargos, além de poder acarretar sobrecarga, a autora aponta que esse poderia ser um indício de baixa participação destas conselheiras nas companhias em que atuam (Abreu, 2024). Não basta, pois, criar condições de acesso.

Assim, o que se debate hoje é que, para favorecer o surgimento de ambientes diversos, é preciso, também, criar políticas de inclusão e permanência de grupos minorizados. O tema já fora muito debatido, por exemplo, no âmbito do ingresso desses grupos às universidades públicas. As primeiras políticas de cotas no ambiente universitário foram pioneiras ao possibilitar o acesso de grupos minorizados aos espaços acadêmicos, antes destinados a uma pequeníssima parcela da sociedade. Com o tempo, percebeu-se que tais políticas não eram suficientes. Deveriam vir acompanhadas, também, de medidas que possibilitassem a permanência e a inclusão desses grupos na universidade. Assim surgiram as diversas iniciativas que vemos hoje nas instituições de ensino superior do país¹⁴.

A mesma lógica pode ser aplicada para se debater a questão da sub-representação feminina no ambiente de trabalho. Se queremos combater o problema da discriminação algorítmica com a inserção de mulheres nos mercados de tecnologia, especialmente na computação, não basta garantir o simples acesso a esses espaços (Requião & Costa, 2022). É necessário traçar estratégias de ascensão e de permanência das mulheres nesse ambiente.

A ONU Mulheres, a WE Impact e a Consultoria GEMA elaboraram uma série de iniciativas práticas possíveis de serem utilizadas

por startups que buscam diversificar seu quadro de funcionários. O relatório “Princípios de Empoderamento das Mulheres para *Startups*” (GEMA, 2021) sugere, por exemplo, (i) realização de mentorias entre as mulheres da empresa; (ii) treinamento de novos colaboradores e novas colaboradoras voltado a questões de gênero; (iii) criação de espaços inclusivos *net-working* para pessoas que estão em posições de cuidado, como mães e aquelas que são responsáveis por idosos; dentre outras políticas e programas de governança que fomentem o desenvolvimento profissional do quadro feminino em empresas de tecnologia.

Ademais, pode-se dizer que uma empresa comprometida em mitigar a discriminação algorítmica (seja de gênero, raça, classe social etc) é também uma empresa que se pretende alinhada aos princípios ESG, tão em voga nas atuais discussões de Direito Empresarial (Eroğlu & Karatepe Kaya, 2022, p. 548).

No Brasil, este debate está, em certa medida, inserido no referido PL nº 2338/2023, em trâmite perante o Senado. Os arts. 19 a 26 de seu texto original estabelecem medidas gerais de governança dos sistemas de inteligência artificial. Por sua vez, o art. 20, IV, estabelece que:

Art. 20. Além das medidas indicadas no art. 19, os agentes de inteligência artificial que forneçam ou operem sistemas de alto risco adotarão as seguintes medidas de governança e processos internos:

[...]

IV – medidas de gestão de dados para mitigar e prevenir vieses discriminatórios, incluindo:

a) avaliação dos dados com medidas apropriadas de controle de vieses cognitivos humanos que possam afetar

a coleta e organização dos dados e para evitar a geração de vieses por problemas na classificação, falhas ou falta de informação em relação a grupos afetados, falta de cobertura ou distorções em representatividade, conforme a aplicação pretendida, bem como medidas corretivas para evitar a incorporação de vieses sociais estruturais que possam ser perpetuados e ampliados pela tecnologia; e

b) composição de equipe inclusiva responsável pela concepção e desenvolvimento do sistema, orientada pela busca da diversidade.

Como se vê, o trecho acima incorpora algumas das sugestões esboçadas na literatura e largamente debatidas neste artigo. Todavia, embora o PL se esforce para prever a necessidade de governança corporativa quanto aos vieses algoritmos, não avança muito quanto à forma que essas medidas podem tomar.

A alínea ‘a’ destaca a necessidade de avaliação dos dados objeto do *machine learning*. Todavia, não detalha como poderia ser feita essa avaliação, se ela necessariamente deve passar pelo crivo humano e nem como isso se daria dentro das empresas e demais entidades do mercado. Já a alínea ‘b’ menciona a implementação de equipes diversas, mas não ressalta a questão da multidisciplinariedade para se enfrentar a complexidade do problema, por exemplo.

Dessa forma, entende-se que o debate sobre vieses discriminatórios de sistemas de IA ainda caminha a passos lentos, principalmente quando se trata de medidas a serem adotadas pelo setor privado para mitigar tal situação.

De fato, essa ausência de robustez regulatória também é um obstáculo enfrentado pelo mercado quando da implementação dessas práticas de governança. Uma vez que não há clareza sequer na definição de IA, tampouco nas formas

de tratá-la, há um vácuo normativo acerca de quais medidas mais concretas de governança corporativa devem ser incentivadas e/ou seguidas pelas empresas.

Tendo em vista essa lacuna, a grande maioria das propostas apresentadas neste capítulo perpassam o aprofundamento de medidas voltadas à governança corporativa. A urgência em se regular tecnologias de inteligência artificial não é um debate que cabe apenas ao âmbito do Direito Público, mas também ao Direito Privado. A discriminação algorítmica de gênero é um problema interseccional, multisetorial e, portanto, deve ser objeto de estudo de todos os campos do saber jurídico.

5 Conclusão

Apesar da falácia de neutralidade das tecnologias de inteligência artificial, a crença de que as IAs são eficientes e livres de subjetividade foi desafiada e desmistificada. Não há, pois, que se falar em tecnologia neutra ou livre de vieses comportamentais. Existem diversas variáveis invisíveis quando falamos da tomada de decisões por algoritmos. Uma delas é, precisamente, a questão de gênero.

Isso traz implicações éticas e políticas que o Direito deve enfrentar, exigindo uma compreensão profunda e interdisciplinar das particularidades dos sistemas algorítmicos de IA para lidar com a discriminação nesse contexto.

A literatura aponta que, em todas as etapas do chamado *machine learning*, pode haver discriminação algorítmica, com especial destaque para a fase de preparação dos dados. Isso porque é nela que os vieses comportamentais dos próprios seres humanos mais podem incidir quando do desenvolvimento do algoritmo.

Tais preocupações fizeram com que União Europeia, por exemplo, propusesse um Guia de diretrizes éticas para se atingir um maior

nível de confiabilidade na inteligência artificial, com especial enfoque para as noções de transparência e não discriminação. Embora referido Guia não tenha apresentado iniciativas práticas, a doutrina buscou mecanismos para abordar os problemas discriminatórios causados pela IA. O Direito Empresarial, por sua vez, teria o papel de propor abordagens regulatórias capazes de mitigar esse problema no âmbito do setor privado

Nesse contexto, várias instituições, empresas, organizações internacionais e governos têm buscado alternativas para atenuar ou resolver a discriminação identificada em tecnologias de inteligência artificial. Por exemplo, o desenvolvimento de ferramentas de transparência por grandes empresas de tecnologia – como IBM, Google e Meta –, na tentativa de abordar os problemas ora discutidos.

Como sugestões práticas, propõe-se um modelo que oferece alternativas para mitigar os efeitos dos vieses comportamentais no *machine learning*, em especial a documentação e transparência do algoritmo, a supervisão humana obrigatória e a criação de equipes multidisciplinares.

Especialmente quanto à última sugestão, o debate se intensifica. A composição dessas equipes não é simples e exige a criação de políticas de inclusão e permanência de grupos minorizados. Isso pode ser observado no debate sobre a sub-representação feminina no campo das ciências exatas e no ambiente de trabalho, onde estratégias de ascensão e inclusão são essenciais para combater a discriminação.

Tais estratégias de mitigação do problema, no entanto, são pouco reguladas ou disseminadas na atualidade. Para que haja um combate mais robusto dos vieses algorítmicos de gênero, faz-se necessário não apenas um papel ativo e integrado do setor privado, mas também dos próprios ramos do saber jurídico.

Referências bibliográficas

- Abreu, T. (2024). A diversidade de gênero na alta administração: tokenismo, conexões pessoais e diretoras de enfeite. *Revista de Direito Mercantil industrial, econômico e financeiro – RDM*, 180/181, 161–212.
- Arrieta, A., B. et al. (2020). Explainable Artificial Intelligence (XAI): *Concepts, taxonomies, opportunities and challenges toward responsible AI*. *Information Fusion*, 58, 82–115.
- Bellamy, R., K., E. et al. (2021). AI Fairness 360: An Extensible Toolkit for Detecting and Mitigating Algorithmic Bias. *IBM Journal of Research and Development*.
- Eroğlu, M., & Karatepe Kaya, M. (2022). Impact of Artificial Intelligence on Corporate Board Diversity Policies and Regulations. *European Business Organization Law Review*, 23(3), 541–572. <https://doi.org/10.1007/s40804-022-00251-5>
- Gema, ONU Mulheres & We Impact (2021). *Princípios de Empoderamento das Mulheres para Startups*. Recuperado de: http://www.onumulheres.org.br/wp-content/uploads/2021/09/WEPs_Startups_Toolkit_POR_BRA_2021_WinWin.pdf.
- Hickman, E., & Petrin, M. (2021). Trustworthy AI and Corporate Governance: The EU's Ethics Guidelines for Trustworthy Artificial Intelligence from a Company Law Perspective. *European Business Organization Law Review*, 22(4), 593–625. <https://doi.org/10.1007/s40804-021-00224-0>
- Latorre Ruiz, E., & Pérez Sedeño, E. (2023). Gender Bias in Artificial Intelligence. In J. Vallverdú (Org.), *Gender in AI and Robotics* (Vol. 235, p. 61–75). Springer International Publishing. https://doi.org/10.1007/978-3-031-21606-0_4
- Lobacheva, A., & Kashtanova, E. (2022). Social Discrimination in the Epoch of Artificial Intelligence. *Wisdom*, 2(1), 97–105. <https://doi.org/10.24234/wisdom.v2i1.767>
- Lindoso, M., C., B. (2021). *Discriminação de gênero no Tratamento Automatizado de Dados Pessoais: como a automatização incorpora vieses de gênero e perpetua a discriminação de mulheres*. Rio de Janeiro: Processo.
- Mulholland, C. (2022). Inteligência artificial e discriminação de gênero. In Schreiber, A.; Martins, G., M., & Carpena, H. (eds.) *Direitos Fundamentais e Sociedade Tecnológica* (pp. 169-182). São Paulo: Foco.
- O'NEIL, C. (2018). *Weapons of Math Destruction: How Big Data increases Inequality and Threatens Democracy*. Nova York: Crown.
- ONU. (2024). *Governing AI for Humanity: Final Report*. Recuperado de: https://www.un.org/sites/un2.un.org/files/governing_ai_for_humanity_final_report_en.pdf.
- ONU. (2011). *Guiding Principles on Business and Human Rights*. Disponível em: https://www.ohchr.org/sites/default/files/documents/publications/guidingprinciplesbusinesshr_en.pdf.
- Pasquale, F. (2015). *The Black Box Society: The secret Algorithms that control Money and Information*. Cambridge, Massachusetts, London, England: Harvard University Press.
- Pisanelli, E. (2022). *A new turning point for women: artificial intelligence as a tool for reducing gender discrimination in hiring*. Presentation at the European University Institute. Recuperado de: https://congress-files.s3.amazonaws.com/2022-08/Pisanelli_paper.pdf.
- Requião, M., & Costa, D., C. (2022) Discriminação algorítmica: ações afirmativas como estratégia de combate. *Civilistica*, 11(3), 1–24.

- Sainz, N., Gabardo, E., & Ongaratto, N. (2024). Discriminação algorítmica no Brasil: uma análise da pesquisa jurídica e suas perspectivas para compreensão do fenômeno. *Revista Direito Público*, 21(110), 258-289. doi: <https://doi.org/10.11117/rdp.v21i110.7295>
- Santos, R. (2020). *Maioria minorizada*. Rio de Janeiro: Telha.
- Silva, T. (2022). *Racismo algorítmico: Inteligência artificial e discriminação nas redes digitais*. São Paulo: Edições Sesc.
- Tonucci, D., C., G., & Caldeira, M., M. (2023). Feminismo de dados: uma nova perspectiva para a ciência e os vieses "inconscientes" de gênero. In Barbosa, B., Tresca, L., & Lauschner, T. (eds.) *3a Coletânea de Artigos—TIC, Governança da Internet, Gênero, Raça e Diversidade—Tendências e Desafios*. Núcleo de Informação e Coordenação do Ponto BR.
- União Europeia. (2019a). *Ethics guidelines for trustworthy AI*. Recuperado de: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.
- União Europeia. (2019b). *Uma definição de IA: principais capacidades e disciplinas científicas*. Recuperado de: <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>.
- Van Giffen, B., Herhausen, D., & Fahse, T. (2022). Overcoming the pitfalls and perils of algorithms: A classification of machine learning biases and mitigation methods. *Journal of Business Research*, 144, 93–106. <https://doi.org/10.1016/j.jbusres.2022.01.076>
- Wang, P. (2019). On Defining Artificial Intelligence. *Journal of Artificial General Intelligence*, 10(2), 1-37. doi: <http://dx.doi.org/10.2478/jagi-2019-0002>.

Notas finais

1 Projeto Algo_Ritmos. Disponível em: <https://algoritmos.internetlab.org.br>.

2 Updates to the OECD's definition of an AI system explained. Disponível em: <https://oecd.ai/en/wonk/ai-system-definition-update>.

3 Disponível em: <https://epocanegocios.globo.com/Empresa/noticia/2018/10/amazon-desiste-de-ferramenta-secreta-de-recrutamento-que-mostrou-vies-contramulheres.html>.

4 Disponível em: <https://www.intercept.com.br/2022/11/24/como-plataformas-de-inteligencia-artificial-podem-discriminar-mulheres-idosos-e-faculdades-populares-em-processos-seletivos/>.

5 Disponível em: <https://digital-strategy.ec.europa.eu/en/policies/european-approach-artificial-intelligence>.

6 Disponível em: <https://www.europarl.europa.eu/news/en/headlines/society/20230601STO93804/eu-ai-act-first-regulation-on-artificial-intelligence>.

6 Disponível em: <https://www12.senado.leg.br/noticias/materias/2023/05/12/senado-analisa-projeto-que-regulamenta-a-inteligencia-artificial>.

7 Disponível em: <https://www12.senado.leg.br/noticias/materias/2023/05/12/senado-analisa-projeto-que-regulamenta-a-inteligencia-artificial>.

8 A igualdade de gênero passou a fazer parte do direito internacional dos direitos humanos pela Declaração Universal dos Direitos Humanos, que foi adotada pela Assembleia Geral em 10 de dezembro de 1948. Disponível em: <http://www.onumulheres.org.br/noticias/igualdade-de-genero-e-assembleia-geral-da-onu-fatos-e-historia-a-saber/>.

9 “Art. 5º Todos são iguais perante a lei, sem distinção de qualquer natureza, garantindo-se aos brasileiros e aos estrangeiros residentes no País a inviolabilidade do direito à vida, à liberdade, à igualdade, à segurança e à propriedade, nos termos seguintes: I - homens e mulheres são iguais em direitos e obrigações, nos termos desta Constituição [...]”.

10 Disponível em: <https://modelcards.withgoogle.com/about>.

11 Disponível em: <https://modelcards.withgoogle.com/about>.

12 Opta-se pela expressão “grupo minorizado” por entender que ela reflete melhor a realidade e o processo que as chamadas “minorias sociais” enfrentam cotidianamente. Conforme Santos (2020), o termo passou a ser utilizado a partir dos anos 90 e deriva da ideia de que determinados segmentos sociais, independentemente da quantidade de indivíduos, têm pouca representação social, econômica e política. Estão, portanto, à margem dos interesses sociais.

13 O regulamento de emissores da B3, por exemplo, contemplou medidas relacionadas a temas Ambientais, Sociais e de Governança Corporativa (conhecidos como ESG), voltadas para a temática da representação feminina nos conselhos de administração ou nas diretorias estatutárias de companhias abertas. Disponível em: https://www.b3.com.br/data/files/3B/31/0A/CF/394798101DBF7498AC094EA8/Regulamento%20de%20Emissores%20_20.07.2023_.pdf#page30.

14 Exemplo disso é o Programa de Apoio à Permanência e Formação Estudantil (PAPFE), na USP, vide: <https://jornal.usp.br/radio-usp/programa-estudantil-na-usp-pretende-assegurar-maior-inclusao-e-permanencia-de-alunos/>.