

Resenha:

Ética na inteligência artificial¹

Review:

AI Ethics

Edgar Lyra

Professor Associado do Departamento de Filosofia da PUC-Rio. Graduou-se em Engenharia Química (UERJ) e depois seguiu caminho pela Filosofia Contemporânea. Concentra seu interesse na Filosofia da Tecnologia e em suas injunções éticas, políticas, retóricas e pedagógicas.

PUC-Rio, Programa de Pós-graduação em Filosofia, Rio de Janeiro (RJ), Brasil.

I. Introdução

O filósofo belga Mark Coeckelbergh é um dos nomes mais expressivos no atual debate sobre o desenvolvimento tecnológico e suas implicações éticas. Sua vasta obra se espalha pelo campo da inteligência artificial, da robótica, do pós-humanismo, do transumanismo, em suas consequências ambientais, políticas e psicossociais.

Coeckelbergh integra várias comissões europeias e conselhos editoriais importantes na área ampla da inteligências artificiais e tem vários livros e artigos publicados sobre o tema, dentre eles *Ética na IA*, hoje traduzido em mais de uma dezena de línguas e objeto desta resenha.

¹ COECKELBERGH, Mark. **AI Ethics**, Massachusetts, MIT Press, MIT Press, 2020, 229 p. Tradução brasileira: **Ética na inteligência artificial**, por Clarisse de Souza, Edgar Lyra, Matheus Barros e Waldyr Delgado, UBU/PUC-Rio, 2023.

Uma curiosidade é que o livro foi originalmente publicado em setembro de 2020 e, por isso, não se ocupa ainda da atual onda de preocupações em torno das IAs gerativas, de que é exemplo o ChatGPT, aberto ao público pela OpenAI no final de 2022. Solicitado pela Editora Ubu que fizesse um prefácio à edição brasileira, o autor delegou a tarefa ao ChatGPT e a validou com pouquíssimas modificações. Aproveitando a menção, toda referência e paginação aqui usadas respeita a edição brasileira.

O livro é importante por múltiplas razões. Oferece ao leitor panorama bastante completo do presente debate ético em torno da inteligência artificial, cobrindo desde as preocupações com um possível “apocalipse AI” até assuntos mais palpáveis e urgentes, como uso de dados e processos digitais de subjetivação com consequências já bem visíveis em âmbitos políticos, jurídicos e sociais. Tanto a resenha como o livro permitirão ao leitor constatar com que velocidade as coisas vêm acontecendo nos anos recentes. Apesar do lapso temporal que separa a publicação original deste ano de 2024, a leitura continua valendo como introdução ampla ao debate em curso sobre inteligência artificial.

Do ponto de vista estilístico, embora o filósofo assuma algumas posições e destile certos temores e simpatias, o tom geral do livro é sóbrio e preferencialmente interrogativo, convidando o leitor a entrar no debate de forma tão livre, aprofundada e abrangente quanto possível. As perguntas que permeiam o texto são em geral recursos heurísticos para a condução da leitura, mas há também casos que se revelam efetivamente abertos, sem resposta momentânea ou conclusiva à vista. Como seja, a linguagem usada pelo autor é clara e a exposição, bastante bem documentada, com adequado e proveitoso conjunto de notas. O livro apresenta ainda um glossário, um índice remissivo (com nomes e conceitos), uma atualizada bibliografia e algumas sugestões de leitura.

II. Sumário detalhado da obra

Dada a importância da obra, a opção aqui feita foi por uma resenha expandida, organizada segundo seus capítulos. O sumário a seguir apresentado é mais detalhado que o do livro, incluindo títulos de subseções presentes apenas no corpo do texto, mas não no índice inicial original. Visa facilitar a visão geral do tratado e proporcionar apoio às sínteses seguintes. É por certo recomendada a leitura integral do texto, mas os capítulos são separadamente inteligíveis, sendo como se segue a sua estrutura:

Prefácio à edição brasileira

1. Espelho, espelho meu - *Mirror, mirror on the wall*

a. IA – euforia e medos: espelho, espelho meu, existe alguém mais inteligente do que eu?

b. O impacto real e pervasivo da IA

c. A necessidade de discutir problemas éticos e sociais

2. Superinteligência, monstros e o apocalipse da IA - *Superintelligence, monsters, and the AI apocalypse*

a. Superinteligência e transhumanismo

b. O novo monstro de Frankenstein

c. Transcendência e apocalipse da IA

d. Como ir além das narrativas de competição e além da hipervalorização

3. Tudo sobre os humanos (*All about human*)

a. A IA geral é possível? Há diferenças fundamentais entre humanos e máquinas?

b. Modernidade, (pós-) humanismo e pós-fenomenologia

4. Simplesmente máquinas? - *Just machines?*

a. Questionando o status moral da IA: agência moral e suscetibilidade moral

b. Agentes morais

c. Suscetibilidade moral

d. Rumo a questões éticas mais práticas

5. A tecnologia - *The technology*

a. O que é inteligência artificial?

b. Diferentes abordagens e subcampos

c. Aplicações e impacto

6. Não se esqueça da ciência de dados - *Don't forget the data (science)*

a. Aprendizado de máquina

b. Ciência de dados

c. Aplicações

7. A privacidade e outros suspeitos habituais - *Privacy and other usual suspects*

a. Privacidade e proteção de dados

b. Manipulação, exploração e usuários vulneráveis

c. Fake news, o perigo do totalitarismo e o impacto nas relações pessoais

d. Segurança e Proteção

8. Máquinas irresponsáveis e decisões inexplicáveis - *A-responsible machines and unexplainable decisions*

- a. Como podemos e devemos atribuir responsabilidade moral?
- b. Transparência e explicabilidade

9. Enviesamento e sentido da vida - *Bias and the meaning of life*

- a. Viés
- b. O futuro do trabalho e o sentido da vida

10. Propostas de políticas - *Policy proposals*

- a. O que precisa ser feito e outras questões que os formuladores de políticas têm que responder
- b. Princípios éticos e justificativas
- c. Soluções tecnológicas e a questão de métodos e operacionalização

11. Desafios para os formuladores de políticas - *Challenges for the policymakers*

- a. Ética proativa: inovação responsável e incorporação de valores ao design
- b. Abordagens orientadas de baixo para cima: como traduzi-las para a prática?
- c. Rumo a uma ética positiva
- d. Interdisciplinaridade e transdisciplinaridade
- e. O risco de um inverno da IA e o perigo de um uso alienado da IA

12. É o clima, estúpido! Sobre prioridades, o Antropoceno e o carro de Elon Musk no espaço - *It's the climate, stupid! On priorities, the Anthropocene, and on Elon Musk's car in space*

- a. A ética na IA deve ser centrada no ser humano?
- b. Acertando nossas prioridades
- c. IA, mudanças climáticas e antropoceno
- d. A nova loucura espacial e a tentação platônica
- e. Retorno à Terra: rumo à IA sustentável
- f. Procura-se: inteligência e sabedoria

III. Resenha dos capítulos

O **CAPÍTULO 1 (Espelho, espelho meu)** cumpre o propósito de caracterizar a ubiquidade e o caráter ao mesmo tempo pervasivo e invisível das atuais tecnologias inteligentes. O autor parte das emblemáticas derrotas para máquinas inteligentes de Gary Kasparov e Lee Sedoll para, em seguida, evidenciar a presença tecnológica cada vez mais decisiva nas áreas dos transportes, marketing, saúde, finanças e seguros, segurança privada e militar, ciência, educação, trabalho de escritório e assistência pessoal (por exemplo, Google Duplex), entretenimento, artes (por exemplo, recuperação de música e composição), agricultura e,

é claro, manufatura.” (p. 14) Deriva dessa onipresença, enfim, “a necessidade de discutir problemas éticos e sociais” (p. 15).

O **CAPÍTULO 2 (Superinteligência, monstros e o apocalipse AI)** apresenta ao leitor o arco do debate sobre as novas tecnologias: prioridades, motivações e preocupações que marcam as posições em competição. Grande parte dele é dedicada a narrativas que flertam com um possível ponto de irreversibilidade a que nos levaria o atual crescimento exponencial das inteligências de máquina, via “auto aprimoramento recursivo” ou “emulação integral do cérebro”, como discutidos por Nick Bostrom no seu livro *Superinteligência* (2014), já traduzido para o português (2018). Não só Bostrom, outros autores de renome fazem especulações ou prognósticos nessa linha, com liberdade crescente e motivações diversas, que vão da denúncia de “risco existencial” a promessas de um futuro redimido dos males que presentemente nos afligem. Max Tegmark, Murray Shanahan, Yuval Harari, Stephen Hawking, Elon Musk e Ray Kurzweil integram assimetricamente esse elenco especulativo.

Coeckelbergh alerta, em contrapartida, para a falta de consenso em torno de conceitos como “inteligência artificial geral” e “singularidade tecnológica”, sobretudo no que concerne à sua exequibilidade e ao tempo que nos separa dessas realizações. Margareth Boden² é uma das pensadoras que questionam os prognósticos de celeridade na consumação de uma inteligência artificial de fato “geral”, entenda-se: não apenas capaz de superar campeões de xadrez, mas de guardar as peças do jogo na sua caixa, dar entrevistas à imprensa sobre o feito, administrar a agenda dele decorrente e instrumentalizar-se para seguir caminho. Outras preocupações dizem respeito à possibilidade de esse frenesi especulativo nos distrair de questões mais urgentes, mais concretas e pervasivas como as listadas no capítulo anterior.

O autor faz um longo excuro pela história do imaginário ocidental sobre a relação homens e máquinas, caracterizando desejos e medos subjacentes aos atuais desenvolvimentos, para finalizar o capítulo indagando como superar os exageros especulativos e conquistar a sobriedade e o distanciamento necessários a uma reflexão ao mesmo tempo rigorosa e ética sobre esses temas. Elenca duas linhas de escape: a primeira recomenda olhar para o Oriente, de modo a irrigar as atuais narrativas de competição

² O livro de referência é *AI – Its Nature and Future*. Oxford University Press, 2016.

com outras formas de pensar as relações entre homens, natureza, transcendência e artifício. A segunda linha é mais programática e precisa ser citada em seus 5 pontos (p. 35):

1. Usar filosofia e ciência para examinar criticamente e discutir as assunções sobre inteligência artificial e sobre a noção de homem envolvida nesses cenários e debates (Por exemplo: É de fato possível uma inteligência artificial geral? Qual é afinal a diferença entre homens e máquinas? Qual é a relação entre humanos e tecnologia? Qual é o status moral da inteligência artificial?);
2. Olhar em mais detalhe para o que são as inteligências artificiais existentes e o que efetivamente hoje fazem essas inteligências em suas várias aplicações;
3. Discutir mais concretamente os urgentes problemas éticos e sociais levantado pela inteligência artificial em suas presentes aplicações;
4. Investigar a política para futuro próximo no que concerne ao desenvolvimento da inteligência artificial;
5. Questionar se a atenção à inteligência artificial no discurso público corrente é justificável e útil à luz dos outros problemas que enfrentamos.

O **CAPÍTULO 3 (Tudo sobre o humano)** é um capítulo substancialmente filosófico, ainda que não se aprofunde nas questões levantadas. Começa com debate sobre a efetiva viabilidade de construção de uma inteligência artificial forte (*strong AI*) com características semelhantes às da cognição humana. Hubert Dreyfus, autor de *What Computers Can't Do?* (1972), texto de cariz fenomenológico, responde pela refutação principal: máquinas não podem emular o senso humano de contexto: em outras palavras, a experiência de ser-no-mundo tematizada por Martin Heidegger, em quem se inspira. A posição inversa, de defesa dessa viabilidade, é referida a filósofos analíticos como Paul Churchland e Daniel Dennett, este último postulando que nós mesmos somos “máquinas conscientes”. O debate é decerto mais matizado, havendo entre os “analíticos” pensadores como John Searle e mesmo o segundo Ludwig Wittgenstein, que não respaldam semelhante parentesco entre homens e máquinas.

Coeckelbergh identifica logo a seguir um problema contíguo ao da possibilidade técnica de emulação (e superação) da inteligência humana: a questão de quão razoável, desejável ou prudente seria esse projeto. É primeiro preciso perguntar se, do ponto de vista dos progressos puramente técnicos, não

seria melhor desatrelar o desenvolvimento da inteligência artificial da matriz humana. De todo modo, nossa humanidade precisa ser hoje novamente interrogada: ela é de fato um bem inalienável ou podemos mais abertamente flertar com perspectivas transumanistas, de expansão cognitiva, de imortalidade e aventuras para além da terra? Autores como Donna Haraway, Bruno Latour e Tim Ingold são mobilizados para indicar a complexidade dessa discussão, que não se restringe à produção de ciborgues de ficção, mas se espalha pela questão das possíveis formas de simbiose entre homens e máquinas, passadas, presentes e futuras, e pela necessidade de colocar em questão o próprio privilégio do “humano” em relação às outras espécies de seres vivos, privilégio que, afinal, nos teria conduzido aos atuais impasses civilizacionais. Tal é o nicho das abordagens pós-humanistas, sempre confrontadas com a possibilidade de se pensar em novas possíveis formas de humanismo, sobretudo não antropocêntricas.

O último ponto elencado pelo autor retoma os interditos postos por Dreyfus na década de 1970, que teriam sido revogados pelo advento das redes neurais e do *machine learning*, na sua capacidade de aprender a partir de “contextos”, agora restituídos a partir de quantidades massivas de dados. A perspectiva conhecida como pós-fenomenológica, que tem entre seus nomes mais notáveis Don Ihde e Paul Verbeek, recusa a perspectiva heideggeriana ainda central para Dreyfus e se concentra na análise da relação dos humanos com tecnologias específicas, retomando a “constituição mútua de humanos e tecnologia”, enfim sugerindo que “a batalha humanista para defender o humano contra a tecnologia é mal direcionada” (p. 47).

O **CAPÍTULO 4 (Apenas Máquinas?)** introduz o problema do status moral das máquinas dividindo-o em dois: o da possibilidade de responsabilização e o da atribuição de direitos às máquinas inteligentes. Coeckelbergh alega que o problema não é tão abstruso ou ficcional quanto possa parecer, bastando pensar nos algoritmos jurídicos que já avaliam a chance de reincidência de presos ou nas tomadas de decisões por carros autônomos; do ponto de vista dos direitos, considere-se o fato de que pessoas efetivamente se apegam a máquinas e que isso costuma depender mais de “como elas estiverem inseridas na nossa vida social, na linguagem e cultura” (p. 60) do que de quaisquer atributos essenciais. O autor se atém a problemas presentes e, por isso, não discute o possível advento das inteligências artificiais gerais,

conscientes e sencientes, a ponto de perguntar, como o faz Murray Shanahan, em *The Technological Singularity*: - Can They suffer?³ (Elas podem sofrer?).

No que concerne à agência moral, o capítulo antecipa a questão da responsabilidade por eventuais atos nocivos perpetrados por inteligências artificiais, a ser tratada mais adiante. Muitos autores e matrizes éticas são arrolados, clássicos e contemporâneos, sendo central a discussão sobre a possibilidade de inserir no desenvolvimento dessas inteligências regras éticas que garantam condutas moralmente adequadas, inclusive, como defendem alguns pós-humanistas e roboeticistas, com vantagens em relação ao julgamento humano, sempre sujeito a paixões e outros fatores heterônomos.

Coeckelberg limita-se nesse capítulo a lembrar que a questão ética não se resume à definição de regras e a sua obediência, pois, nesse caso, poderíamos definir como agente moral um algoritmo que barrasse competentemente os e-mails maliciosos a nós enviados; sem falar no problema de eventuais conflitos entre regras (vide Isaac Asimov) e da dificuldade de sua aplicação a casos particulares, problema este último que nada tem de trivial e remete às objeções de Dreyfus da década de 1970. A problemática é de todo modo retomada a partir do capítulo 8.

O **CAPÍTULO 5 (A tecnologia)** deixa momentaneamente de lado as questões éticas (somente retomadas ao seu final) para ganhar vista sobre o que a inteligência artificial hoje é e faz. O problema começa com a definição de *inteligência* e com a questão já abordada em capítulos anteriores: se o parâmetro tem que ser a inteligência humana ou se é possível falar de inteligência em sentido mais genérico. Coeckelbergh acrescenta que, paralelamente à discussão sobre ter a mente humana como modelo, há ainda o problema da relação mente-corpo. Diz ele: “transumanistas sonham com mentes futuras não mais presas a suportes biológicos”.

Uma pequena história da busca da inteligência artificial é reorganizada, história que retroage aos primórdios civilizacionais e passa por nomes como Alan Turing, Norbert Wiener, John McCarthy e Marvin Minsky, concluindo que o que hoje temos são inteligências artificiais restritas (*narrow or weak AIs*), e não gerais (*strong AIs*), e que é duvidoso se e quando teremos inteligências artificiais no último sentido.

³ MIT Press, 2015. Trata-se de uma reedição da questão posta no século XVIII por Francis Bacon, que muitos consideram como marco inaugural do que hoje é a “ética animal”.

Outro ponto importante levantado no capítulo é a possibilidade de olhar para a inteligência artificial duplamente: de modo científico ou técnico. A primeira perspectiva visa o progresso no entendimento da inteligência humana e se liga às ciências cognitivas, à psicologia, à ciência de dados e às neurociências. A segunda perspectiva visa a realização de tarefas úteis ou relevantes. Esclarece que essas perspectivas não raro se misturam e que, num plano mais geral, também a matemática, a linguística e a filosofia estão ligadas à discussão sobre inteligência artificial. Por óbvio, informática, engenharia e ciência da computação completam o caldo multidisciplinar. O autor pondera, por fim, que, embora a ciência também comporte questões éticas, o foco do livro está nos impressionantes impactos gerados pelo desenvolvimento técnico da inteligência artificial.

Mas sobre o quê, concretamente, estamos falando quando nos referimos a inteligências artificiais: algoritmos, sistemas de informação, *machine learning*, robôs com forma humana, *bots* de redes sociais, internet das coisas? É problemática mesmo a distinção entre software e hardware – sem falar das contiguidades com outros desenvolvimentos tecnológicos como nanotecnologia e biotecnologia. O autor se posiciona: seu foco é a instância que considera mais geral, os algoritmos e suas combinações. Diz que, para compreender uma ética da inteligência artificial, precisamos compreender como os algoritmos trabalham e o que fazem ou são capazes de fazer.

Na esteira dessa necessária compreensão, vem a distinção entre as inteligências artificiais simbólicas, suas “árvores de decisão” ou cadeias se-então-senão que caracterizavam as GOFAI (*Good Old-fashioned Artificial Intelligences*) e o novo “paradigma conexionista” (*connectionism*), surgido nos anos 1980 e hoje mais conhecido como tecnologia de *redes neurais*, que Coeckelbergh se ocupará de explicar em detalhe no capítulo seguinte em sua correlação com os *big data*, *machine learning* e *deep learning*. Limita-se, no capítulo 5, a dizer que muitos dos sistemas hoje conhecidos como inteligentes são híbridos, ou seja, misturam essas duas formas de computação. Vale citar uma passagem na íntegra para dar amplitude a esse elenco:

Outro paradigma importante na IA é aquele que usa abordagens mais corporificadas e situacionais, focando em tarefas motoras e na interação, em vez das chamadas tarefas cognitivas superiores. Os robôs construídos por pesquisadores de IA, como Rodney Brooks, do MIT, não resolvem problemas usando representações simbólicas, mas interagindo com o ambiente ao redor. Por exemplo, o robô humanoide Cog, desenvolvido na década de 1990, foi construído para aprender interagindo com o mundo – como os bebês fazem. (p. 74)

O autor termina o capítulo evidenciando dois fatos: que a pesquisa em inteligência artificial hoje se divide por campos de problemas, como a identificação visual, o processamento de linguagem natural e a análise de massas de dados; e que, por quaisquer caminhos, seu desenvolvimento é hoje socialmente onipresente, na indústria, na agricultura, no transporte, na saúde, no entretenimento, nos mercados financeiros, com efeitos positivos ou negativos. Entre as questões éticas que assim se põem, está aquela associada às perguntas: “Quem terá acesso à tecnologia e será capaz de usufruir de seus benefícios? Quem será capaz de aumentar seu poder por meio da inteligência artificial? Quem estará excluído desses ganhos?” (p. 76)

O **CAPÍTULO 6 – Não esqueça dos dados (Ciência)** – se debruça, na esteira do anterior, sobre *machine learning* e *data science*. O autor questiona a propriedade do uso do termo “aprendizado” e diagnostica que o aprendizado de máquina é na verdade estatística aplicada ao reconhecimento de padrões em grandes massas de dados, visando a previsões baseadas nessas regularidades. O autor dá exemplos bastante didáticos:

[...] para construir um algoritmo que reconhece imagem de gatos, os programadores não fornecem um conjunto de regras que definem para o computador o que são gatos, mas, em vez disso, o algoritmo constrói seu próprio modelo de imagens de gatos. Ele será otimizado para alcançar a mais alta precisão de previsão em um conjunto de imagens de gatos e não-gatos. Visa, assim, aprender o que são imagens de gatos. Seres humanos fornecem retornos para a máquina sobre seu desempenho, mas não as abastecem com instruções específicas ou regras. (p. 82)

Autores como Ethem Alpaydin e seu livro *Machine Learning* (2016)⁴ são mencionados, insistindo o autor na ideia de que o que se extrai nesses processos estatísticos são padrões, e não dados propriamente ditos.

Três tipos de *machine learning* são elencados: supervisionado, não supervisionado e com reforço de aprendizado. Os nomes são já em si mesmos eloquentes: no aprendizado supervisionado, há indicação ostensiva de conjuntos de dados e variáveis-alvo, o que não acontece no aprendizado não supervisionado, em que, por vezes, o algoritmo descobre regularidades até então desconhecidas dos experts humanos. O reforço de aprendizado, por sua vez, caracteriza-se pela presença decisiva de alguma “função de recompensa”, que informa ao sistema sobre estar ou não progredindo ou fazendo um bom trabalho.

⁴ Cambridge: MIT Press, 2016.

Coeckelbergh esclarece que nos três casos há humanos envolvidos, não obstante o grau distinto de autonomia da máquina em cada um dos tipos.

O *machine learning* está, por tudo isso, muito ligado à *data science*. Os conjuntos de dados precisam ser selecionados e preparados, segundo critérios estatísticos de adequação e suficiência; os algoritmos de mineração precisam ser construídos ou escolhidos; enfim, os resultados precisam ser interpretados. Trata-se de estatisticamente partir de dados particulares em busca de generalizações confiáveis. A referência bibliográfica principal é nesse momento o livro *Data Science* (2018)⁵, publicado por John Kelleher e Brendan Tierney.

Coeckelbergh evoca novamente Margareth Boden para enfatizar que a inteligência artificial, no seu presente estágio, carece de compreensão de *relevância* e, ainda, de *experiência, sensibilidade e sabedoria*, donde a centralidade da agência humana e a conseqüente reponsabilidade associada a essa agência. Um problema estatístico clássico é evocado: o fato de que correlações não são necessariamente relações causais. Algumas ilustrações são extraídas do livro *Spurious Correlations* (2015), de Tyler Vigen⁶: por exemplo, a correlação, num dado período do estado do Maine, entre o aumento da taxa de divórcios e o consumo de margarina.

Outro problema central denunciado por Kelleher e Tierney é o da abstração dos conjuntos de dados em relação às realidades que pretendem representar. Acrescente-se ao problema da abstração, por fim, o dos fatores subjetivos e enviesamentos (não explicitamente) presentes nas escolhas feitas.

O capítulo termina com uma longa enumeração de aplicações já decisivamente em vigência no atual estágio de desenvolvimento da inteligência artificial, conjunto que, num sentido usual dessas palavras, não tem nada de fraco ou restrito. Sem o cuidado de aqui identificar as aplicações específicas associadas a cada uma dessas empresas ou projetos – na suposição de que o leitor saberá fazê-lo –, encontramos no texto menções a Amazon, Walmart, Experian, American Express, BMW, IBM's Watson, Associated Press, Hello Barbie (ToyTalk), Instagram e Netflix. O autor termina dizendo que a *estatística*, outrora nada sexy, é agora “a nova mágica” (p. 91).

⁵ Cambridge: MIT Press, 2015.

⁶ Ver tylervigen.com/spurious-correlations.

O **CAPÍTULO 7 (Privacidade e outras suspeições habituais)** desloca novamente o foco de questões técnicas para as questões éticas. Não me ocorre melhor começo para resenhá-lo que a transcrição abaixo:

Um uso ético da IA requer que os dados sejam coletados, processados e compartilhados de maneira a respeitar a privacidade dos indivíduos e seu direito de saber o que acontece seus dados, de acessá-los, de se opor à sua coleta ou processamento, de saber que seus dados estão sendo coletados e processados e (se aplicável) que estão sujeitos a uma decisão tomada por uma IA. (p. 93)

O autor se apressa em explicar que “é relativamente fácil respeitar esses valores e direitos ao fazer uma pesquisa como cientista social” (p. 94), mas o que fora desse âmbito acontece está geralmente muito distante desses preceitos, com naturalização de práticas de manipulação, exploração, trabalho digital, vigilância, perfazendo uma espécie de totalitarismo democrático. Esclarece, em seguida, que não somente os usuários são objeto de exploração, mas também os trabalhadores que produzem o hardware e aqueles que treinam os algoritmos. Essa “ordenha de dados” atropela sobretudo usuários mais vulneráveis, como crianças e idosos.

Não fosse o suficiente, o uso assim naturalizado da inteligência artificial tem produzido efeitos políticos dos quais é paradigmática a ação eleitoral da *Cambridge Analytica*. As possibilidades abertas têm como efeito colateral o advento de “um mundo onde não é mais claro o que é verdade e o que é falso, onde fatos e ficção se misturam” (p. 98), com diversos e consideráveis danos para o tecido social.

O capítulo termina levantando mais uma vez o problema da responsabilização por consequências danosas do uso das inteligências artificiais, estruturais ou acidentais. Não bastassem *riscos de segurança* intrínsecos a práticas estatais, como armas letais autônomas, o autor considera que, “em um mundo em rede, todo dispositivo eletrônico ou software pode ser hackeado, invadido e manipulado por pessoas com intenções maliciosas” (p. 100), o que cria insegurança e imbróglio ainda maior.

O **CAPÍTULO 8 (Máquinas não-responsáveis e decisões inexplicáveis)** retoma e evidencia o problema da *atribuição de responsabilidade* por eventuais efeitos nocivos das atuais tecnologias, num escopo na verdade mais amplo que o das inteligências artificiais. A quem responsabilizar se um carro autodirigido se envolve num acidente com mortes, como aconteceu no estado do Arizona em 2018?

Coeckelbergh se vale de matriz ética aristotélica para descartar a possibilidade de *responsabilização de máquinas*. Elas não preencheriam, pelo menos no seu presente estágio de desenvolvimento, as duas condições de responsabilização definidas pelo filósofo grego: agência efetiva e consciência das

consequências das ações. Embora se possa conceder que máquinas tenham agência, que possam originar ou interromper ações com consequências louváveis ou condenáveis, o autor pontua que “falta a elas consciência, livre arbítrio, emoções, aptidão para formar intenções e outras condições similares” (p. 104).

Resta assim atribuir responsabilidade aos humanos que as constroem, selecionam, associam, operam e zelam pelo seu bom funcionamento. Mas como exatamente atribuí-la com equidade, quando as cadeias de *agências* que suportam os respectivos ecossistemas tecnológicos são cada vez mais longas e complexas, caracterizando o que o autor chama de problema das *many hands* e das *many things*? No caso de acidente com carro autodirigido, a cadeia se distribui pelos desenvolvedores da inteligência artificial, pelo pessoal de manutenção, pelos donos do carro, envolve as condições do trânsito, a autoridade e a regulação local, sem contar com a dificuldade de separar claramente onde começa o campo de agência da inteligência artificial e onde termina o dos sensores e demais equipamentos a ela ligados. A coisa fica ainda mais difícil se considerarmos que esses dispositivos “têm história”, que preveem uso de partes desenvolvidas alhures e, muitas vezes, deslocadas do seu contexto original de produção.

Tudo isso para dizer que, já no que concerne à definição da agência, o problema de atribuição de responsabilidade não é trivial. A coisa se complica ainda mais quando se passa ao segundo ponto da definição aristotélica, ou seja, ao real conhecimento prévio das consequências das ações que envolvem as referidas tecnologias. É possível auditar sistemas baseados em árvores de decisão, mas sistemas baseados em *deep learning* tornam o rastreamento dos caminhos que levam a uma decisão simplesmente impossível. O autor dá como exemplo dessa impossibilidade a reconstrução do caminho que leva uma inteligência artificial a selecionar seu próximo lance num jogo de xadrez.

Ainda assim, é preciso encontrar meios de contemplar o direito das pessoas afetadas a explicações plausíveis para a decisões subjacentes às “engenharias” que as afetaram ou afetam. E, neste ponto, vale resenhar um outro artigo do autor – “Artificial Intelligence, Responsibility Attribution, and a Relational Justification of Explainability”⁷ –, preferencialmente à finalização do respectivo capítulo, por ser mais claro e propositivo que o texto do livro no tratamento do tema. Particularmente importante é que o autor defende um *direcionamento ético relacional*, sugerindo que essas explicações sejam produzidas no próprio projeto e implementação da inovação tecnológica, não apenas *a posteriori*, quando da ocorrência de um

⁷ In *Science and Engineering Ethics*, 2020, n. 26, p. 2051–2068.

acidente ou constatação de dano perene. Seria esse o imperativo da *agência responsável* e da possibilidade de se continuar falando de *ética* após o advento e disseminação das inteligências artificiais.

A exigida capacidade de responder por atos e escolhas (*answerability*) traz, por sua vez, inúmeros desdobramentos práticos: a começar pela questão do que seja uma boa explicação ou uma *explicação convincente*.⁸ Ainda no mesmo artigo, Coeckelbergh esclarece que, mesmo que a responsabilidade não possa ser pensada em termos absolutos, visto que muitas vezes estão envolvidos aspectos sociais, históricos e institucionais “que impedem um *completo* controle individual, e mesmo coletivo”⁹, isso não pode significar a suspensão da busca por *explicabilidade responsável*. Pelo contrário, e com ainda mais razão, deve fomentar o debate amplo e público das concomitantes questões, incluindo a *tragicidade* que muitas vezes assombra a dimensão tecnológica.

Fica por discutir, tanto no livro quanto no artigo, como esse *compromisso prévio com a explicabilidade* haveria de encontrar acolhida numa sociedade acelerada e tecnicista, pouco reflexiva e nada afeita a exercícios hermenêuticos dessa monta, sobretudo nos nichos mais duros de produção tecnológica. Parece insofismável que a esse acolhimento deva corresponder uma substancial reforma educacional e cultural, formal e informal.

O **CAPÍTULO 9 (Enviesamento e sentido da vida)** traz para o primeiro plano um problema diretamente associado ao *machine learning* e à sua relação com as bases de dados de que se serve: o da não neutralidade ou do “enviesamento” dos resultados em função de escolhas feitas, em geral inconscientemente, pelos desenvolvedores. O autor entende que os enviesamentos são mais vezes não intencionais, ou seja, não há ação deliberada de produzir efeitos discriminatórios. Não obstante, as consequências da falta de percepção desses vieses podem ser severas, individual e coletivamente falando. O primeiro exemplo sai do sistema jurídico: o algoritmo COMPAS, que prediz a chance de reincidência em delitos de réus norte-americanos. Estudos indicam que os “falsos positivos” (réus considerados reincidentes potenciais que, todavia, não o fazem) ocorrem principalmente com negros; e que os “falsos negativos” (réus considerados não reincidentes potenciais e que, todavia, voltam a cometer crimes) ocorrem principalmente com réus brancos. Sua fonte é o livro de Hannah Fry (2018): *Hello World – Being Human in the Age of Algorithms*.

⁸ Ver Lyra, Edgar. **O esquecimento de uma arte** – retórica, educação e filosofia no século 21. São Paulo, Almedina, 2021.

⁹ Coeckelbergh, op. cit., p. 2065.

O autor esclarece que o enviesamento pode ocorrer em qualquer das etapas do desenvolvimento dos programas e ter naturezas diversas: de gênero, cor da pele, cultura, expressão corporal e por aí fora. Pode reforçar tendências nocivas existentes nos tecidos sociais, prejudicando pessoas pertencentes a grupos historicamente marginalizados.

O ponto mais sensível e menos técnico do enviesamento é o fato de ele não poder ser inteiramente evitado, sendo em certas situações até desejável, no sentido de corrigir tendências contextuais injustas. A questão aí subjacente é filosoficamente espinhosa, obrigando à discussão e à explicitação das noções de justiça e equidade presentes, conforme o caso, na defesa da busca da possível neutralidade ou da validade de enviesamentos afirmativos.

A camada mais profunda da questão anterior tem, na verdade, a ver com o mundo que desejamos construir. Se isso é verdade, o desenvolvimento das inteligências artificiais precisa frontalmente considerar questões de justiça e equidade – e não apenas tecnicidades – na medida em que transformará radicalmente, por exemplo, a economia e a divisão do trabalho.

Como a segunda parte do título do capítulo indica, o texto assume a partir desse ponto fôlego decisivamente filosófico, com questões levantadas no amplo arco que vai da filosofia política e social ao existencialismo. Como repensar as relações entre trabalho e renda, considerando a crescente capacidade de as máquinas se ocuparem não apenas do trabalho operário, mas também do intelectual? Como evitar a concentração da riqueza nas mãos dos proprietários das novas máquinas inteligentes? Que tipo de trabalho reservar para os humanos? O que significaria falar, afinal, de uma “inteligência artificial sustentável”? Haveria uma direção válida para todas as culturas hoje postas frente a frente no mundo tecnicamente globalizado?

Por abertas que possam ser essas questões, sua lembrança deve no mínimo nos precaver contra fantasias de paraísos pós-industriais e utopias do gênero, sobretudo capazes de nos desviar da tarefa que origina o livro. Fato é que o acolhimento dessas questões traz consigo não apenas problemas metafísicos, mas desafios prático-políticos nada triviais. Por exemplo: Como digerir e definir planos de ação na velocidade em que as coisas vêm ultimamente acontecendo? Como realisticamente administrar dissensos e viabilizar a implementação de diretrizes despojadas, em meio ao espírito de competição e sanha de

poder característica do tecnocapitalismo? Tais questões levam o autor aos capítulos seguintes, voltados prioritariamente para as concernentes questões políticas, legislativas e jurídicas (*policymaking*).

O **CAPÍTULO 10 (Propostas políticas)** começa afirmando que, dadas as questões éticas que envolvem a inteligência artificial, “algo preciso ser feito”, ainda que não se saiba bem o quê. O autor nomeia instâncias mobilizáveis: providências legislativas, medidas tecnológicas e educacionais, definição de standards e códigos de ética. A sabedoria prática envolveria, enfim, saber não apenas o que fazer, mas por quê, quando, como e por quem, tendo sempre em vista a natureza, a extensão e a urgência dos problemas. Coeckelbergh enumera seis necessidades seminais:

1. Justificar as medidas propostas;
2. Agir preventivamente, ou seja, não apenas quando as novas tecnologias já se encontram incrustadas no tecido social, com todos os seus efeitos colaterais disseminados;
3. Trabalhar pela otimização institucional e pela implementação das boas legislações já existentes;
4. Atentar para a pluralidade de atores que podem e necessitam ser conjuntamente mobilizados para que as medidas tenham eficácia;
5. Evitar o foco na novidade e na urgência, considerando a cauda longa e a correta articulação dos problemas com suas vizinhanças;
6. Não perder de vista a complexidade e o escopo amplo da questão ética em epígrafe.

O texto segue mapeando “a avalanche de iniciativas e documentos existentes” nos vários escopos: estados, organizações supraestatais, universidades e empreendimentos privados. Esforça-se por mostrar que há recorrência de princípios nessas iniciativas, a despeito das suas naturezas e prioridades plurais. Entre os governos, cita o dos EUA (Obama), Reino Unido, França, Áustria e China; entre as instituições supra ou multinacionais, faz referência à ICDPPC,¹⁰ à AI HLEG¹¹ e ao também europeu EGE;¹² entre as legislações, a principal é a GDPR,¹³ de 2018; e, entre as universidades com esforços notórios no campo, lista Montreal, Cambridge, Stanford e o Markkula Centre for Applied Ethics, da Santa Clara University. Faz,

¹⁰ International Conference of Data Protection and Privacy Commissioners

¹¹ European Commission High-Level Expert Group on Artificial Intelligence

¹² Group of Ethics in Science and New Technologies

¹³ General Data Protection Regulation. Por razões expostas na introdução, não há menção ao *AI Act*.

por fim, referência ao mundo corporativo, identificando iniciativas da parte de DeepMind, IBM, Amazon, Apple, Sony, Facebook, Google e Microsoft, dando, ao final da listagem, voz ao CEO da Apple, Tim Cook, que teria reputado como imperativa a necessidade de regulação técnica, na medida em que “o livre mercado não funciona” nesse sentido. Obviamente, há discussão sobre quem deve regular o quê e em que medida, num arco que vai da autorregulação à reivindicação de políticas mais centralizadas e coercitivas. Coeckelbergh arrola ainda outros atores, como a Digital Europe, que representa a indústria digital na Europa, a Campaign to Stop Killer Robots, o encontro de Asilomar, o Future of Life Institute e o IEEE,¹⁴ que em 2017 avançou uma Global Initiative on Ethics of Autonomous and Intelligent Systems.

A despeito dos dissensos a respeito da sua implementação, os princípios são, como já foi dito, mais ou menos recorrentes. Desfilam pelo capítulo menções aos direitos humanos; à promoção de bem-estar e benefício social; à inclusão; ao combate à discriminação e aos viesamentos injustos; ao respeito à privacidade e segurança dos dados; à promoção de autonomia; à proteção contra a propaganda abusiva e a manipulação; ao resguardo do direito à informação relevante sobre a ação dos algoritmos e suas tomadas de decisão; ao compromisso com práticas democráticas; à excelência científica; à confiabilidade; à explicabilidade; à transparência; à prestação de contas; ao cuidado ambiental; à sustentabilidade e à responsabilidade; enfim, à observação de princípios gerais de justiça, equidade, beneficência e não maleficência.

O autor comenta muito brevemente cada uma dessas iniciativas, apontando seus limites. Tem claro que “uma coisa é citar alguns princípios éticos e outra é descobrir como implementá-los na prática.” (p. 154) Apenas para dar corpo ao desafio, faz referência (e problematiza) as intenções de incorporação de princípios éticos ao design das novas tecnologias; ao advento de “caixas pretas éticas”, capazes, à semelhança das usadas em aviões, de apontar a causa dos erros em caso de acidentes e prejuízos; à instituição de “licenças de direção” para veículos autônomos; e à criação de “máquinas morais”, ainda que apenas funcionalmente falando. O campo dessa discussão é naturalmente amplo, incluindo a adoção de protocolos, práticas educativas e legislações, não sendo nada óbvio como articular e dar peso às várias espécies de medidas que reivindicam o rótulo de “éticas”.

¹⁴ Institute of Electrical and Electronic Engineers

O autor começa o **CAPÍTULO 11 (Desafios para os elaboradores de políticas)** retomando a ideia de que é preciso considerar as demandas éticas ainda nas fases iniciais do desenvolvimento das tecnologias inteligentes. Entendendo que não é fácil vislumbrar de antemão possíveis consequências indesejadas desses desenvolvimentos, encoraja a construção de cenários que ajudem nesse sentido. Mas não chega a delinear esse processo, limitando-se a recomendar o artigo escrito por Wessel Reijers et al. *Methods for Practising Ethics in Research and Innovation* (2018)¹⁵.

Sua constatação seguinte é a de que, em se tratando de dar materialidade às ações éticas, os procedimentos *top-down* se revelam abstratos e pouco adequados à realidade plural das partes interessadas (*stakeholders*). Recomenda a observação de práticas democráticas, *bottom-up*, baseadas no diálogo e na escuta, com maior atenção aos processos que aos ditames.

Grande obstáculo a essas práticas é a enorme concentração de poder nas mãos de poucas corporações, sendo necessária a intervenção estatal para resguardo do interesse público. A referência aí é o escrito de Paul Nemitz *Constitutional Democracy and Technology in the Age of Artificial Intelligence* (2018)¹⁶. Fica por saber em que medida os estados contemporâneos têm, de fato, autonomia em relação ao gigantismo corporativo.

Outro ponto da maior importância é o fato de que, concretamente falando, se trata de *modificar costumes, hábitos*, não apenas de enunciar *valores*. Coeckelbergh chega a se referir ao clássico conceito de *formas de vida*, de Ludwig Wittgenstein, para enfatizar que não se trata exatamente de definir listas de princípios e regras de conduta, mas de conseguir que elas sejam acolhidas pelas várias partes concernidas. Essa tarefa é ainda mais difícil na medida em que os contrapontos éticos são muitas vezes percebidos como entraves à lubrificação dos processos produtivos, sobretudo quando não fazem concessões a lavagens de imagem e expedientes do gênero. O que precisa se transformar é, portanto, a própria percepção do papel e da importância da ética na construção “do futuro”, o que, convenhamos, nada tem de trivial em sociedades imediatistas como o são sociedades ocidentais contemporâneas. O autor incentiva voltar o olhar para outras culturas políticas em busca de aprendizado, nomeadamente para culturas orientais, mas fica aqui a pergunta sobre se essas culturas não são historicamente epigonais em relação ao Ocidente, bem entendido, no que concerne ao desenvolvimento tecnológico hoje hegemônico.

¹⁵ In *Science and Engineering Ethics*, v. 24, n. 5, p. 1437-81, 2018.

¹⁶ In *Philosophical Transactions of the Royal Society*, v. 376, n. 2144, 2018.

Direção mais promissora parece ser o investimento em práticas de inter e transdisciplinaridade.

Vale citar uma passagem eloquente:

Precisamos garantir que, por um lado, pessoas com uma formação em humanidades tenham ciência da importância de pensar sobre as novas tecnologias, como a IA, e possam assim adquirir algum conhecimento dessas tecnologias e o que elas podem fazer. Do outro lado, cientistas e engenheiros precisam se tornar mais sensíveis aos aspectos éticos e sociais do desenvolvimento da tecnologia e seu uso. (p. 165)

Esse é um caminho que precisa igualmente vencer resistências, mas que pode, efetivamente, contribuir para que a ética deixe de ser uma tópica marginal em relação ao desenvolvimento tecnológico e se torne parte essencial desse desenvolvimento. Como seja, o autor encerra o capítulo chamando atenção para o perigo de uma lida banal e imprudente com os crescentes poderes e capacidades de transformação da natureza e das sociedades disponibilizados pelas novas inteligências artificiais.

O CAPÍTULO 12 (É o clima, estúpido! Sobre prioridades, antropoceno e o carro de Elon Musk no espaço) é o capítulo de encerramento do livro. Retoma questões ligadas ao “humano-centrismo” que vinham sendo tangenciadas desde o capítulo 2. De fato, o debate e as concomitantes legislações se organizam inercialmente em torno do que o autor chama de “narrativas de competição” (homem x tecnologia), discutindo o resguardo do “humano” diante da crescente hegemonia tecnológica. Pouca atenção é concedida aos efeitos sobre os não humanos e mesmo sobre os humanos historicamente excluídos do desenvolvimento em epígrafe.

Antes, todavia, de voltar-se mais decisivamente para essa miopia, Coeckelbergh faz uma espécie de prolepse, antecipando possíveis críticas à sua reivindicação de atenção à pauta da inteligência artificial.

Vale citar um trecho mais extenso:

Olhando para a agenda referente ao desenvolvimento sustentável (os chamados Objetivo do Desenvolvimento Sustentável) das Nações Unidas de 2015, e sua visão geral das questões globais relativas ao que o então secretário-geral da ONU Ban Ki Moon chamou de “pessoas e planeta”, vemos muitas questões globais que exigem atenção ética e política: desigualdades crescentes dentro dos países e entre eles, guerra e extremismo violento, pobreza e desnutrição, falta de acesso à água potável, falta de instituições democráticas e eficazes, envelhecimento da população, doenças infecciosas e epidêmicas, riscos relacionados à energia nuclear, falta de oportunidade para crianças e jovens, desigualdade de gênero e várias formas de discriminação e exclusão, crises humanitárias e todo o tipo de violações dos direitos humanos, problemas relacionados a migração e os refugiados, além de mudanças climáticas e problemas ambientais [...] como desastres naturais mais frequentes e intensos, e formas de degradação ambiental como a seca e perda de biodiversidade. Diante desses enormes problemas, a IA deveria ser a nossa prioridade máxima? Será que a IA nos distrai de assuntos mais importantes? (p. 170-171)

A defesa da pauta tecnológica feita pelo autor se agarra primeiro à tópica do possível “agravamento dos problemas existentes”, sociais e ambientais. O desenvolvimento imprudente da inteligência artificial ameaça elevar ainda mais o consumo de energia *per capita*, aprofundar o fosso entre ricos e pobres, além de aumentar o “controle” sobre a vida no planeta, mais e mais em mãos de inteligências artificiais. O autor considera também a possível ajuda da inteligência artificial na mitigação da atual ameaça climática, mas adverte que essa ajuda pode levar a tecno-solucionismos e, por aí, a novas formas de autoritarismo. Tudo isso justificaria manter no topo da lista de prioridades éticas a inteligência artificial.

Não deixa de ser curiosa, nesse ponto, a falta de atenção do autor ao problema retórico trazido pela nova economia digital, problema da transformação radical no processo de formação de opinião e, por aí, da definição do que tem ou não tem importância, capitaneado pelos algoritmos e redes sociais. A própria emergência de negacionismos científicos e políticos de todas as espécies não pode ser negligenciada na sua relação como os novos caminhos digitais de formação de subjetividades.

O último alerta trazido pelo autor é o perigo de seduções tecnicistas do tipo “colonização do espaço” e congêneres, sobretudo pela cifra de alienação que gera em relação aos problemas terrenos concretos e presentes, tendo como epítome o ideário de “sobrevivência dos mais ricos”, com referência a autores como Hannah Arendt¹⁷ e Douglas Rushkoff¹⁸.

A finalização do capítulo e do livro aponta para o necessário advento de uma espécie híbrida de sabedoria. Nas palavras do autor: “Tal sabedoria pode ser abastecida por processos cognitivos abstratos e análise de dados, mas também tem como base experiências corporificadas, relacionais e situacionais no mundo, na lida com outras pessoas, com a materialidade e com o nosso ambiente natural” (p. 184).

IV. Considerações finais

O livro contém um elenco bastante maior de autores e obras citadas, além de boas leituras sugeridas, sendo seu exame direto fortemente recomendado como introdução aos problemas éticos trazidos pelo presente desenvolvimento das inteligências artificiais. Como também já ponderado na introdução, a aceleração característica destes tempos recentes faz com que os mais de três anos

¹⁷ ARENDT, Hannah. **A condição humana** [1958]. Rio de Janeiro, Forense, 2007.

¹⁸ RUSHKOFF, Douglas. Survival of the Richest, **Medium**, 5, jul. 2018.

transcorridos desde a publicação original ocasione que o leitor possa se ressentir de algumas lacunas referentes a questões recentíssimas, das quais os desdobramentos das *inteligências artificiais gerativas* constituem o grupo mais importante. Um exemplo importante são as questões éticas ligadas ao uso dessas ferramentas na pesquisa científica e na produção de trabalhos acadêmicos geral. Mas essa constatação é antes um índice do fenômeno tecnológico em análise que uma real lacuna, ajudando a compreender o efetivo nível de disrupção inaugurado pela abertura a público do ChatGPT pela Open AI. Ou seja, o texto permanece atual da sua inusitada “datação”.

Edgar Lyra

ORCID: <https://orcid.org/0000-0003-3664-3537>

PUC-Rio, Programa de Pós-graduação em Filosofia, Rio de Janeiro (RJ), Brasil

Doutor em Filosofia pela PUC-Rio

E-mail: edlyra@puc-rio.br

Recebido em: 15 de setembro de 2024.

Aprovado em: 20 de setembro de 2024.

Este artigo é publicado em acesso aberto (Open Access) sob a licença Creative Commons Attribution Non-Commercial (CC-BY-NC 4.0), que permite que outros remixem, adaptem e criem a partir do seu trabalho para fins não comerciais, e embora os novos trabalhos tenham de lhe atribuir o devido crédito e não possam ser usados para fins comerciais, os usuários não têm de licenciar esses trabalhos derivados sob os mesmos termos.