

تأثیر هوش مصنوعی بر تحول هنجارهای اخلاقی: رویکردی مبتنی بر

روان شناسی اخلاق

علی اکبرپور^۱

^۱دانشجوی کارشناسی ارشد روان شناسی تربیتی، دانشگاه تهران، Ali.akbarpoor@ut.ac.ir

چکیده

این پژوهش به بررسی نقش هوش مصنوعی در دگرگونی هنجارهای اخلاقی انسان از منظر روان شناسی اخلاق می پردازد. هدف مطالعه، تبیین چگونگی تأثیرگذاری سامانه های هوش مصنوعی، به ویژه الگوریتم های توصیه گر و پلتفرم های تعاملی، بر شکل گیری و بازتعریف ارزش هایی چون عدالت و همدلی است. در این پژوهش، با رویکردی نظری و مبتنی بر مرور انتقادی ادبیات، مفاهیم اصلی روان شناسی اخلاق و کارکردهای شناختی - عاطفی در قضاوت های اخلاقی تحلیل شده و با ویژگی های فناوری های نوین هوش مصنوعی مقایسه می شوند. نتایج تحلیل نشان می دهد که هدایت توجه و اولویت بندی محتوا توسط الگوریتم ها می تواند الگوهای قضاوت اخلاقی افراد را تغییر داده و در برخی موارد به کاهش حساسیت نسبت به پیامدهای انسانی تصمیم ها بینجامد. در عین حال، طراحی سامانه های هوش مصنوعی بر پایه اصول روان شناختی و توجه به سازوکارهای همدلی و انصاف، این ظرفیت را دارد که به تقویت هنجارهای اخلاقی همسو با ارزش های انسانی کمک کند. بر این اساس، ضرورت ادغام ملاحظات روان شناختی در فرایند توسعه و استقرار فناوری های هوش مصنوعی مورد تأکید قرار می گیرد.

کلید واژه: ارزش های انسانی، روان شناسی اخلاق، فناوری هوش مصنوعی، هنجارهای اخلاقی

۱- مقدمه:

در دهه های اخیر، هوش مصنوعی به یکی از تعیین کننده ترین عوامل تحول در زندگی فردی و اجتماعی انسان بدل شده است. از الگوریتم های توصیه گر در شبکه های اجتماعی و پلتفرم های محتوایی گرفته تا دستیارهای هوشمند و سیستم های تصمیم یار، بخش قابل توجهی از تجربه های روزمره ما در بستری شکل می گیرد که به طور نامرئی توسط سامانه های هوش مصنوعی سازمان دهی می شود. این تحول صرفاً جنبه فنی یا کارکردی ندارد، بلکه به تدریج نحوه ادراک ما از خود، دیگران و جهان اجتماعی را نیز دگرگون می کند. در چنین زمینه ای، یکی از پرسش های اساسی آن است که هوش مصنوعی چگونه می تواند بر هنجارهای اخلاقی انسان و شیوه های داوری درباره «درست» و «نادرست» اثر بگذارد.

اخلاق هوش مصنوعی عمدتاً در قالب مباحثی چون شفافیت الگوریتمی، مسئولیت پذیری، حریم خصوصی و تبعیض الگوریتمی پیگیری شده است؛ اما کمتر به این موضوع پرداخته شده که خود مواجهه مداوم با سامانه های هوش مصنوعی چگونه می تواند ساختار قضاوت های اخلاقی انسان و حساسیت او نسبت به ارزش هایی مانند عدالت و همدلی را تغییر دهد. به بیان دیگر، مسئله فقط این نیست که «سیستم اخلاقی باشد یا نه»، بلکه این است که «تعامل با سیستم، اخلاق ما را چگونه بازسازی می کند». این پرسش در حوزه روان شناسی اخلاق، که به سازوکارهای شناختی و عاطفی قضاوت های اخلاقی می پردازد، اهمیت ویژه ای دارد.

نظریه های کلاسیک و معاصر در روان شناسی اخلاق - از الگوهای رشد اخلاقی تا رویکردهای عاطفی-شناختی - نشان می دهند که قضاوت های اخلاقی برآیند تعامل میان فرآیندهای شناختی (مانند استدلال، توجه و تفسیر موقعیت) و فرآیندهای عاطفی (مانند همدلی، احساس گناه یا انزجار اخلاقی) هستند. در فضای دیجیتال متکی بر هوش مصنوعی، هر دو بُعد شناختی و عاطفی در معرض شکل دهی و هدایت الگوریتمی قرار می گیرند. الگوریتم های توصیه گر با اولویت بندی انواع خاصی از محتوا، می توانند بر آنچه می بینیم، به چه چیز حساس می شویم و چه نوع روایت هایی از عدالت و رنج انسانی در معرض دید ما قرار می گیرد، اثر بگذارند. این فرایند می تواند به مرور، مرزهای «طبیعی» و «پذیرفته شده» در قضاوت های اخلاقی را جابه جا کند.

از سوی دیگر، هوش مصنوعی صرفاً تهدیدی برای ارزش های اخلاقی به شمار نمی آید، بلکه می تواند به عنوان ابزاری برای تقویت هنجارهای اخلاقی انسانی نیز به کار گرفته شود؛ به شرط آن که طراحی آن بر پایه فهم دقیق تری از سازوکارهای روان شناختی اخلاق انجام گیرد. برای مثال، می توان سامانه هایی طراحی کرد که به طور هدفمند تعاملات همدلانه را تقویت کنند، تنوع دیدگاه ها را افزایش دهند و از تشدید حباب های اطلاعاتی و قطبی سازی اخلاقی جلوگیری نمایند.

با وجود این ظرفیت ها، شکاف محسوسی میان ادبیات اخلاق هوش مصنوعی و یافته های روان شناسی اخلاق وجود دارد. بسیاری از بحث های رایج در حوزه اخلاق فناوری، بیش از آن که بر فرآیندهای روان شناختی تحول هنجارهای اخلاقی تمرکز کنند، معطوف به مقررات، اصول کلی یا پیامدهای اجتماعی هستند. این مقاله می کوشد با اتخاذ رویکردی مبتنی بر روان شناسی اخلاق، به این پرسش بپردازد که تعامل مستمر با سامانه های هوش مصنوعی چگونه می تواند به بازتعریف هنجارهای اخلاقی مرتبط با عدالت و همدلی بینجامد.

در این چارچوب، ابتدا به تبیین مفاهیم بنیادین در روان شناسی اخلاق و جایگاه عواطف و شناخت در قضاوت اخلاقی پرداخته می شود. سپس نقش سازوکارهای الگوریتمی در جهت دهی به ادراک ها و داورهای اخلاقی تحلیل می گردد و در ادامه، پیامدهای این تحلیل برای طراحی اخلاق محور سامانه های هوش مصنوعی مورد بحث قرار می گیرد. در نهایت، با جمع بندی یافته های نظری، بر ضرورت ادغام ملاحظات روان شناختی در توسعه و سیاست گذاری مرتبط با هوش مصنوعی برای حفظ و تقویت ارزش های انسانی تأکید می شود.

۲-۲- روش تحقیق

این پژوهش از نوع مطالعه نظری - تحلیلی است و با رویکردی مبتنی بر مرور نظام مند مفاهیم و ادبیات علمی انجام شده است. هدف، تبیین سازوکارهای روان شناختی مؤثر در تحول هنجارهای اخلاقی در تعامل با سامانه های هوش مصنوعی است؛ بنابراین روش تحقیق بر تحلیل مفهومی و تطبیقی استوار است و شامل مراحل زیر بوده است:

۲-۱- جمع آوری ادبیات تخصصی

منابع معتبر حوزه روان شناسی اخلاق، نظریه های عاطفی-شناختی، مطالعات اخلاق فناوری و آثار پژوهشی مرتبط با تأثیر الگوریتم ها بر رفتار انسان بررسی شده اند. برای انتخاب منابع، از پایگاه هایی مانند *Scopus*، *Google Scholar* و *IEEE Xplore* استفاده شده است. تحلیل نظری چارچوب های روان شناختی نظریه های کلیدی از جمله رشد اخلاقی کالبرگ، نظریه بنیادهای اخلاقی هیدت، مدل های همدلی، و الگوهای شناختی قضاوت اخلاقی استخراج و مبانی مفهومی آنها دسته بندی شد. مقایسه تطبیقی میان داده های روان شناختی و سازوکارهای هوش مصنوعی نحوه عملکرد الگوریتم های توصیه گر، سیستم های تعاملی و سازوکارهای شخصی سازی محتوا تحلیل شد و با یافته های روان شناختی مرتبط با توجه، پردازش عاطفی، سوگیری های شناختی و قضاوت اخلاقی تطبیق داده شد. استخراج الگوی اثرگذاری هوش مصنوعی بر هنجارهای اخلاقی بر اساس تطبیق مفاهیم، یک چارچوب نظری پیشنهادی برای نشان دادن مسیرهای احتمالی تأثیرگذاری هوش مصنوعی بر ارزش هایی مثل عدالت و همدلی تدوین شد. تحلیل پیامدها و پیشنهادات اخلاقی نتایج تحلیلی در قالب پیامدهای نظری و کاربردی برای طراحی سامانه های اخلاق محور هوش مصنوعی ارائه شده است. به دلیل ماهیت نظری پژوهش، ابزارهای میدانی و آماری به کار گرفته نشده و تمرکز اصلی بر تحلیل استدلالی، انسجام مفهومی و استنتاج نظری بوده است.

۳-۳- مبانی نظری و پیشینه پژوهش

۳-۱- روان شناسی اخلاق و الگوهای کلاسیک قضاوت اخلاقی

روان شناسی اخلاق به مطالعه سازوکارهای شناختی و عاطفی دخیل در تصمیم گیری اخلاقی می پردازد. یکی از نظریه های بنیادین در این حوزه، نظریه رشد اخلاقی کالبرگ^۱ است که تحول قضاوت اخلاقی را در سه سطح پیش قراردادی، قراردادی و پس قراردادی توضیح می دهد. این مدل نشان می دهد که ظرفیت انسان برای داور اخلاقی در تعامل با تجارب اجتماعی، آموزش و تکامل شناختی شکل می گیرد. در کنار این رویکرد شناختی، دیدگاه های نوین تر بر تعامل عواطف و شناخت تأکید دارند. جان اتان هیدت^۲ با ارائه نظریه بنیادهای اخلاقی^۳ بیان می کند که قضاوت های اخلاقی محصول واکنش های سریع عاطفی مانند همدلی، خشم و انزجار هستند و استدلال اخلاقی معمولاً در مرحله ای ثانویه وارد می شود. در این چارچوب، ارزش هایی مانند عدالت و مراقبت مستقیماً تحت تأثیر محیط اجتماعی و محرک های بیرونی قرار دارند.

^۱Kohlberg, 1984

همچنین جاشوا گرین^۴ با مدل «دو فرایندی» خود نشان می دهد که تصمیم های اخلاقی در تعامل دو سیستم شکل می گیرند:

(۱) سیستم سریع، هیجانی و شهودی

(۲) سیستم کند، تحلیلی و مبتنی بر استدلال

این مدل به طور مستقیم قابل استفاده برای تحلیل نحوه تأثیر الگوریتم ها بر قضاوت های انسانی است؛ زیرا هوش مصنوعی اغلب سیستم سریع شهودی را تحریک می کند و فرصت پردازش تحلیلی را کاهش می دهد.

۲-۳- تغییر هنجارهای اخلاقی و نقش محیط اجتماعی

هنجارهای اخلاقی مجموعه باورهای مشترکی هستند که رفتار پذیرفته شده را در جامعه تعریف می کنند. پژوهش های الایر تورنر^۵ و شور و وینبرگر^۶ نشان می دهد که هنجارهای اخلاقی در اثر سازوکارهای فرهنگی، رسانه ای و ارتباطی تغییر می کنند. با گسترش رسانه های دیجیتال، زمینه های جدیدی برای شکل گیری ارزش ها پدید آمده است؛ به گونه ای که کانال های اطلاعاتی نقش مهمی در گسترش سوگیری های شناختی، کاهش تنوع دیدگاه ها و تغییر حساسیت اخلاقی ایفا می کنند.

۳-۳- نقش هوش مصنوعی و الگوریتم ها در شکل دمی ارزش های اخلاقی

مطالعات معاصر در حوزه اخلاق هوش مصنوعی نشان می دهد که سامانه های هوشمند تنها ابزارهای فنی نیستند، بلکه بخشی از محیط اجتماعی-شناختی انسان محسوب می شوند.

لوتجیانو فلوریدی^۷ به عنوان یکی از برجسته ترین نظریه پردازان «اخلاق اطلاعات»، معتقد است که فناوری های هوشمند «محیط اخلاقی» را بازسازی می کنند و بر نحوه درک انسان از مسئولیت، عدالت و هنجارهای اجتماعی اثر می گذارند.

همچنین کیت کراوفورد^۸ و همکاران (در حوزه سوگیری الگوریتمی) نشان داده اند که الگوریتم های توصیه گر با ایجاد «حباب های اطلاعاتی» و شخصی سازی شدید محتوا می توانند حساسیت اخلاقی افراد را تغییر دهند. این سامانه ها با هدایت توجه به انواع خاصی از محتوا، شیوه مواجهه افراد با رنج انسانی، عدالت اجتماعی و رفتار اخلاقی را دست کاری می کنند.

در حوزه طراحی سیستم های اخلاق محور، وندل والچ^۹ و پیتیر آلن^{۱۰} چارچوب هایی مطرح کرده اند که بر لزوم ادغام اصول اخلاقی و روان شناختی در طراحی سامانه ها تأکید دارد. این پژوهش ها نشان می دهند که هوش مصنوعی می تواند هم هنجارهای اخلاقی را تضعیف کند و هم در شرایط طراحی شده و نظارت شده، آن ها را تقویت نماید.

۴-۳- جمع بندی پیشینه علمی

ادبیات علمی نشان می دهد که:

کالبرگ، هیدت و گرین سازوکارهای روان شناختی اخلاق را توضیح داده اند.

فلوریدی، کراوفورد و ونکاتاسوبرامانیان نقش فناوری را در تغییر هنجارهای اخلاقی بررسی کرده اند.

اما ترکیب این دو حوزه و تحلیل روان شناختی اثر الگوریتم ها بر هنجارهای اخلاقی هنوز محدود است.

^۴Haidt, 2001

^۵Moral Foundations Theory

^۶Greene, 2013

^۵Turner, 1991

^۶Shweder & Weinberger, 1992

^۷Luciano Floridi

^۸Kate Crawford, 2021

^۹Wendell Wallach

^{۱۰}Peter Allen

از این رو، پژوهش حاضر تلاش می کند با ادغام این دو سطح نظری، سازوکاری را نشان دهد که طی آن هوش مصنوعی می تواند به بازتعریف ارزش هایی مانند عدالت و همدلی منجر شود.

۴- بحث و تحلیل

تحلیل تحول هنجارهای اخلاقی در عصر هوش مصنوعی مستلزم ترکیب دو حوزه نظری است:

(۱) روان شناسی اخلاق (کالبرگ، هیدت، گرین)

(۲) اخلاق فناوری و هوش مصنوعی (فلوریدی، کراوفورد، ونکاتاسوبرامانیا، والاچ)

از این منظر، تعامل انسان با سامانه های هوشمند تنها یک کنش تکنیکی نیست، بلکه فرایندی است که بر ساختارهای شناختی و عاطفی قضاوت اخلاقی اثر می گذارد.

۴-۱- تأثیر الگوریتم ها بر سازوکارهای شناختی اخلاق

مطابق نظریه رشد اخلاقی کالبرگ، قضاوت اخلاقی انسان در سطوحی از استدلال مبتنی بر اصول شکل می گیرد. اما سامانه های توصیه گر با ایجاد محیط های شناختی محدود و شخصی سازی شده، ظرفیت افراد برای پردازش تحلیلی و استدلال اخلاقی را کاهش می دهند.

به تعبیر جاشوا گرین (۲۰۱۳)، تصمیم های اخلاقی نیازمند تعادل میان سیستم هیجانی سریع و سیستم کند شناختی هستند؛ اما محیط الگوریتمی رسانه ها با تسریع جریان محرک های هیجانی و کوتاه کردن فرصت پردازش، باعث تقویت سیستم سریع و تضعیف سیستم کند می شود. این وضعیت می تواند به قضاوت های سطحی تر، واکنش های احساسی شدید و کاهش عمق استدلال اخلاقی منجر شود.

۴-۲- اثر هوش مصنوعی بر حساسیت های عاطفی اخلاقی

بر اساس نظریه بنیادهای اخلاقی هیدت، ارزش هایی همچون همدلی، مراقبت و عدالت ریشه های عاطفی دارند. الگوریتم های مبتنی بر اولویت بندی محتوای هیجانی می توانند الگوهای عاطفی را بازتعریف کنند.

کیت کراوفورد (۲۰۲۱) نشان می دهد که پلتفرم ها تمایل دارند محتواهای برانگیزاننده ُ خشم یا انزجار را بیشتر نمایش دهند؛ زیرا این واکنش ها تعامل بیشتری تولید می کنند.

این امر به تدریج حساسیت نسبت به رنج انسانی را کاهش می دهد و واکنش های همدلانه (که زیربنای اخلاق مراقبت اند) را تضعیف می کند. همچنین تحلیل های سورس ونکاتاسوبرامانیا در حوزه سوگیری الگوریتمی نشان می دهد که فناوری های توصیه گر با تقویت کلیشه های اخلاقی یا اجتماعی ممکن است «چارچوب های عاطفی» افراد را در جهت گیری های خاص ثابت نگه دارند. این فرایند می تواند در بلندمدت به کاهش انعطاف اخلاقی و کاهش توانایی انسان در درک دیدگاه های متفاوت منجر شود.

۴-۳- بازتعریف هنجارهای عدالت در محیط های هوشمند

عدالت یکی از بنیادی ترین ارزش های اخلاقی در نظریه های کالبرگ و هیدت است. در محیط های دیجیتال، نوع مواجهه افراد با مسائل عدالت محور تحت تأثیر نحوه توزیع اطلاعات و تجربه های شخصی سازی شده قرار می گیرد.

لوسانو فلوریدی اشاره می کند که سامانه های هوش مصنوعی «معماری اخلاقی محیط» را تغییر می دهند؛ یعنی تعیین می کنند که افراد چه چیزی را ببینند و چه چیزی خارج از میدان توجه آن ها قرار گیرد. این بازطراحی محیطی به بازطراحی هنجارهای عدالت نیز می انجامد.

۴-۴- ظرفیت های هوش مصنوعی برای تقویت هنجارهای اخلاقی

در کنار پیامدهای چالش برانگیز، پژوهشگران اخلاق فناوری مانند وندل والاچ و پیتر آلن بر پتانسیل هوش مصنوعی برای تقویت ارزش های اخلاقی تأکید کرده اند. به شرط طراحی صحیح، سامانه های هوشمند می توانند:

تنوع دیدگاه های اخلاقی را افزایش دهند،

محتوای همدلانه را تقویت کنند،

سوگیری ها را کاهش دهند،

کاربران را در تصمیم گیری های اخلاقی یاری دهند.

برای نمونه، طراحی سامانه هایی که مداخلات مبتنی بر همدلی ارائه می کنند یا محتواهای اخلاقی را با دقت بیشتری برجسته می کنند، می تواند هنجارهای عدالت و مراقبت را تقویت نماید.

۵-۴- جمع بندی تحلیل

تحلیل حاضر نشان می دهد که هوش مصنوعی دو سطح از سازوکارهای اخلاقی انسان را درگیر می کند: سطح شناختی: جهت دهی به توجه، کاهش استدلال اخلاقی، تقویت سیستم سریع تصمیم گیری سطح عاطفی: تغییر حساسیت های همدلانه، تقویت واکنش های هیجانی، دست کاری قالب های عاطفی نتیجه این فرایند، بازتعریف تدریجی هنجارهای اخلاقی است؛ به این معنا که عدالت، همدلی و مسئولیت پذیری در بستری شکل می گیرند که توسط الگوریتم ها طراحی و هدایت می شود.

۵- نتیجه گیری

یافته های این پژوهش نشان می دهد که هوش مصنوعی نه فقط به عنوان یک ابزار فناوری، بلکه به عنوان یک بافت (context) روان شناختی جدید عمل می کند که در آن، فرایندهای رشد اخلاقی، قضاوت اخلاقی و هویت اخلاقی افراد شکل می گیرند و دگرگون می شوند. از منظر روان شناسی اخلاق، تعامل مداوم با سامانه های هوشمند می تواند سطوح مختلف تجربه اخلاقی را تحت تأثیر قرار دهد: در سطح شناختی، الگوریتم های توصیه گر و محیط های دیجیتال شخصی سازی شده، الگوهای توجه، تفسیر موقعیت های اجتماعی و دسترسی به اطلاعات را تغییر می دهند. این تغییرات می تواند مسیر طی شده از قضاوت های پیش قراردادی تا پس قراردادی را، آن گونه که کالبرگ توصیف می کند، دچار دگرگونی کند؛ زیرا فرد بخش قابل توجهی از مواجهه های اخلاقی خود را در فضایی تجربه می کند که توسط منطق الگوریتمی فیلتر شده است. در نتیجه، فرصت مواجهه با تعارض های اخلاقی پیچیده و تمرین استدلال اصول محور کاهش یافته و قضاوت های اخلاقی بیش از پیش به الگوهای تثبیت شده و گروه محور وابسته می شود. در سطح عاطفی-شهودی، نظریه بنیادهای اخلاقی هیدت به ما یادآور می شود که ارزش هایی مثل مراقبت و انصاف ریشه در واکنش های عاطفی و شهودی دارند. محیط های مبتنی بر هوش مصنوعی، از طریق برجسته سازی مکرر برخی محرک های هیجانی (مانند خشم، ترس یا تحقیر) و کم رنگ سازی موقعیت های برانگیزاننده همدلی، به تدریج «پروفایل عاطفی اخلاقی» فرد را بازتنظیم می کنند. این فرایند می تواند به کاهش حساسیت همدلانه، عادی شدن رنج دیگران و تقویت قضاوت های قطبی شده بینجامد؛ پدیده ای که برای سلامت روان اجتماعی و انسجام جمعی مخاطره آمیز است. در سطح هویت اخلاقی و خودآگاهی، زندگی در بافتی که واکنش ها، ترجیحات و قضاوت های فرد مدام توسط سامانه های هوشمند بازتاب و تقویت می شوند، می تواند تجربه فرد از «من اخلاقی» را تغییر دهد. فرد ممکن است کمتر در معرض بازخوردهای چالش برانگیز قرار گیرد، کمتر نیاز به بازنگری خودانتقادی در باورهای اخلاقی اش احساس کند و در نتیجه، فرایندهای رشد هویت اخلاقی کند یا منحرف شود. برای روان شناسان تربیتی، مشاوران و درمانگران، این نکته پیامدهای مستقیم دارد؛ چراکه مداخلات اخلاقی و آموزشی دیگر نمی توانند مستقل از بافت هوش مصنوعی طراحی شوند. در عین حال، نتایج این مطالعه نشان می دهد که همین سازوکارها می توانند به صورت معکوس به خدمت تقویت اخلاق انسانی نیز درآیند. اگر طراحان سامانه های هوش مصنوعی با تکیه بر یافته های روان شناسی اخلاق عمل کنند، می توان محیط هایی ساخت که:

- مواجهه با دیدگاه های متنوع اخلاقی را تسهیل کند،
- موقعیت های برانگیزاننده همدلی را برجسته سازد،
- فرصت تمرین استدلال اخلاقی را افزایش دهد،
- و رشد هویت اخلاقی خودآگاه و مسئول را حمایت کند.

بنابراین، این پژوهش بر ضرورت شکل گیری گفت و گویی فعال بین روان شناسی و علوم داده تأکید می کند. برای حفظ و تقویت سلامت اخلاقی فرد و جامعه در عصر هوش مصنوعی، کافی نیست که صرفاً درباره «اخلاقی بودن» یا «نبودن» فناوری صحبت کنیم؛ بلکه باید به طور مشخص بررسی کنیم که این فناوری ها چگونه بر پردازش های شناختی، الگوهای هیجانی، هویت اخلاقی و روابط میان فردی اثر می گذارند و بر اساس این شناخت، هم در سطح طراحی سامانه ها و هم در سطح مداخلات آموزشی و درمانی، رویکردهای جدیدی توسعه دهیم.

مراجع سبک (APA)

Greene, J. D. (2013). *Moral tribes: Emotion, reason, and the gap between us and them*. Penguin Press.

Haidt, J. (2001). *The emotional dog and its rational tail: A social intuitionist approach to moral judgment*. *Psychological Review*, 108(4), 814-834.

Kohlberg, L. (1984). *The psychology of moral development: The nature and validity of moral stages*. Harper & Row.

Crawford, K. (2021). *The atlas of AI: Power, politics, and the planetary costs of artificial intelligence*. Yale University Press.

Floridi, L. (2013). *The ethics of information*. Oxford University Press.

Venkatasubramanian, S. (2019). *The social consequences of algorithmic bias*. *Annual Review of Statistics and Its Application*, 6, 29–43.

Wallach, W., & Allen, C. (2009). *Moral machines: Teaching robots right from wrong*. Oxford University Press.