

Sesgos y discriminaciones sociales de los algoritmos en Inteligencia Artificial: una revisión documental

Social biases and discriminations of algorithms in Artificial Intelligence: a literature review

Rodrigo Ramirez Autrán*

Artículo recibido: 27-09-23

Artículo aprobado: 05-11-23

Palabras clave:

Inteligencia Artificial, agentes artificiales, sesgos y brechas, discriminación social, Social Construction of Technology (SCOT).

Keywords:

Artificial Intelligence, artificial agents, biases and gaps, social discrimination, Social Construction of Technology (SCOT).

Cómo citar este artículo

Ramirez Autrán, R. (2023). Sesgos y discriminaciones sociales de los algoritmos en Inteligencia Artificial: una revisión documental. *Entretextos*, 15(39), 1-17. <https://doi.org/10.59057/iberoleon.20075316.202339664>.

Resumen

El objetivo de este trabajo es analizar el concepto y aplicación de la Inteligencia Artificial (IA) a la luz de los posibles sesgos, discriminaciones y desigualdades al momento de ser implementada en la sociedad. Mediante una investigación en repositorios especializados, se revisa una serie de externalidades negativas de los sistemas de IA, las cuales se podrían manifestar en aquellos grupos e individuos que carecen de la posibilidad de acceso, no cuentan con los medios para evaluar la calidad y veracidad de los datos masivos o que simplemente son segregados y/o discriminados por los algoritmos debido a su color de piel, etnia, género u otro aspecto. Este escenario fomentaría con mayor recurrencia dichas externalidades e incluso podría crear nuevas brechas digitales, sociales, políticas y económicas. Una de las preguntas clave que guiaron la investigación fue: ¿qué consecuencias éticas, morales y sociales podría implicar la presente generación en Inteligencia Artificial?

* Instituto de Investigaciones Sociales de la Universidad Nacional Autónoma de México / Centro de Estudios Antropológicos de la Facultad de Ciencias Políticas y Sociales de la Universidad Nacional Autónoma de México.

La principal contribución de este artículo es la de proveer una revisión documental y ejemplificada con casos particulares a nivel internacional sobre el impacto social de dicha tecnología, particularmente desde una perspectiva de las Ciencias Sociales.

Abstract

The objective of the following work is to analyze the concept and application of Artificial Intelligence (AI) in light of the possible biases, discriminations and inequalities when it is implemented in society. Through research in specialized repositories, a series of negative externalities of AI systems are addressed, which could be manifesting in those groups and individuals who: lack the possibility of access, among those who do not have the means to evaluate the quality and veracity of massive data, or that are simply segregated and/or discriminated against by algorithms due to their skin color, ethnicity, gender or other. Such a scenario would encourage these externalities with greater recurrence and could even create new digital, social, political and economic divides. One of the key questions that guided the research was: what ethical, moral and social consequences do the present generation in Artificial Intelligence could imply? The main contribution of this article is to provide a document review exemplified with particular cases at an international level on the social impact of said technology, particularly from a Social Sciences perspective.

Introducción

Como en otros momentos de la historia, uno de los motores del cambio social son las transformaciones tecnológicas en los modos de producción, en el mundo del trabajo y en la socialización, las cuales a su vez pueden traer como consecuencia importantes transformaciones socioeconómicas, también conocidas como revoluciones industriales (Rifkin, 2015), que han tenido como base el descubrimiento de nuevas formas de producción, de energía, de nuevos medios de comunicación y nuevos arreglos sociales, así como de organización.

En los últimos años, las discusiones enmarcadas en la denominada Cuarta Revolución Industrial y una de sus tecnologías, la Inteligencia Artificial (IA), se han extendido más allá de lo tecnológico, lo empresarial y las comunidades académicas hacia el ámbito público, en general (UNESCO, 2023; Wright, 2018; Floridi *et al.*, 2018). Un síntoma de ello es la creciente información en los medios masivos de comunicación y las redes sociales digitales, los cuales continuamente están mostrando que la IA desempeñará un papel cada vez más importante en nuestra vida diaria, y ejercerá una creciente influencia en las economías nacionales y mundiales (Polak, 2021).

Estamos en un momento de avance tecnológico, donde cada semana somos bombardeados por noticias e investigaciones nuevas sobre las maravillas que traerá consigo la IA en nuestra vida, junto con el crecimiento exponencial de nuevas empresas de base tecnológica que la utilizan como herramienta determinante en muchos de sus procesos, al mismo tiempo que poco o nada sabemos los usuarios y consumidores sobre el impacto sociocultural, ético y moral que traerá consigo la actual vorágine tecnocientífica.

Somos testigos del uso cada vez mayor de una variedad de algoritmos basados en IA, también llamados agentes artificiales (Osoba y Welsler, 2017), que son utilizados para resolver problemas de la vida cotidiana, mientras que se busca la mejora en la productividad y la eficiencia, al mismo tiempo que “las sumas astronómicas que hoy invierten los países desarrollados en IA ilustran claramente el enorme potencial que incuba” (Madrigal, 2020, p. 107).

Algunos datos importantes sobre la fuerza económica y el impacto sociocultural de la IA son:

Los asistentes de voz impulsados por IA alcanzarán los 8 mil millones para 2023. Para 2025, se espera que el mercado global de IA sea de casi \$60 mil millones. La mayor cantidad de habilidades de inventario de Alexa están disponibles en los EE. UU.: alrededor de 66 000 habilidades. El PIB mundial crecerá en 15.7 billones de dólares para 2030 gracias a la IA. La IA puede aumentar la productividad empresarial en un 40%. El número de nuevas empresas de IA creció 14 veces en las últimas dos décadas. La inversión en nuevas empresas de IA creció 6 veces desde 2000. Ya el 77% de los dispositivos que usamos cuentan con una forma de IA. (Techjury, 2021)

Hoy sabemos que la IA, más allá de nuestro uso inmediato y cotidiano, tiene un impacto arrollador en muchos ámbitos de la vida: espacios laborales, educativos, médicos, industriales, financieros, deportivos y culturales, en los que pareciera ser que se están modificando diversas formas de interacción, aprendizaje, procesos de automatización, toma de decisiones, entre muchas otras. Al mismo tiempo, herramientas tecnológicas como la IA muestran cierta precisión arbitraria, con lo que se modifican aspectos de nuestra vida, debido a que se mueven en el terreno de las decisiones y los escenarios (Rezaev y Tregubova, 2018), lo cual podría tener un fuerte impacto en esferas sociales de tipo ético y moral.

Por ejemplo, un sector que a nivel global se ha beneficiado de la IA es el de los servicios financieros. Los algoritmos diseñados a partir de la IA han auxiliado a este sector en procesos tales como la calificación crediticia y la evaluación de riesgos; sin embargo, su uso exacerbado y exponencial podría repercutir negativamente en una buena parte de la población si dichos procesos de calificación y evaluación crediticia no toman en

cuenta los valores, así como las necesidades del conjunto de la sociedad. Eventos tan desastrosos como el llamado *Flash Crash* del año 2010, una de las mayores quiebras bursátiles en la historia suscitada por agentes inteligentes, evidencian estos riesgos en los sectores financieros.

Más allá del ejemplo anterior, algunos autores se han preguntado: “si permitimos que la IA tome decisiones que cambien la vida, ¿cuáles son los costos, de oportunidad, compensaciones e implicaciones para los sistemas sociales, económicos, técnicos, legales y ambientales?” (Elliott *et al.*, 2021, p. 179). Los costos e implicaciones resaltados por Elliot *et al.*, inevitablemente están relacionados con un tipo de delegación de poder y decisión a los sistemas basados en algoritmos, cuyos funcionamientos son opacos (Bucher, 2018) e inobservables para buena parte de sus usuarios, denominándose modelos de caja negra (Pinch y Bijker, 1989), concepción de la tecnología como un artefacto del que se conocen sus condiciones de entrada (*input*) y de salida (*output*), pero no su funcionamiento interno.

De acuerdo con Sued (2022): “los algoritmos tienen una doble materialidad: digital en su constitución y cultural en su funcionalidad” (p. 45). De esta forma, la tecnología puede ser vista como un reflejo de las relaciones sociales en un contexto histórico particular y, por ello, no puede entenderse como neutra. La contingencia de observar estos sistemas inteligentes como opacos, cerrados y no neutros, recae en la necesidad imperante de una comprensión de su diseño, crucial a la hora de analizar las formas sociales de transferencia y apropiación de la tecnología, entendiendo a los algoritmos como objetos sociotécnicos (Pinch, 2009).

Al tomar en cuenta lo anterior, se propone estudiar el impacto social de la IA con base en la teoría de la construcción social de la tecnología (Social Construction of Technology [SCOT, por sus siglas en inglés]), la cual se suscribe dentro del campo de los estudios de Ciencia y Tecnología (CTS). Los constructivistas sociales sostienen que la tecnología no determina la acción humana, sino que la acción humana da forma a la tecnología. Ésta podría ayudar, entre otras cosas, a ver que estos sistemas no vienen acompañados solamente de promesas de bonanza económica para los países, las empresas y las organizaciones, y que su propio diseño en muchas ocasiones entraña serios riesgos en diversos ámbitos, que cada vez demandarán mayor atención y análisis por parte de los gobiernos, organizaciones y las personas, ya que “la tecnología y la sociedad se construyen mutuamente, y tecnología, sociedad y materialidad se encuentran en interacción continua” (Pinch, 2009, p. 45).

Un ejemplo sobre la importancia del diseño y la construcción social de la tecnología es el siguiente: cotidianamente en nuestra vida utilizamos, cada vez con mayor frecuencia, diversos agentes de IA generativa, como *Dall-e* o *Stable Diffusion* (rama de la IA que se

enfoca en la generación de contenido original a partir de datos existentes, y utiliza algoritmos y redes neuronales avanzadas para aprender de textos e imágenes, para luego generar contenido nuevo y único [Granieri, 2023]), así como los famosos *chatbots* (programas informáticos basados en IA y procesamiento natural del lenguaje [PNL] para comprender las preguntas realizadas por parte de sus usuarios y, con ello, automatizar las respuestas, simulando una conversación cotidiana con un ser humano [IBM, s.f.]). Éstos se encuentran en una etapa de desarrollo y perfeccionamiento, por lo cual es común encontrarnos con su mal funcionamiento o con las llamadas alucinaciones, encontradas en la literatura como confabulaciones o delirios, que son calificativos que revelarían parte del comportamiento sesgado de la IA (Ji *et al.*, 2022), que no son otra cosa que respuestas seguras (cómodas) de una IA, que no parecen estar justificadas por sus datos de entrenamiento, posiblemente con un sesgo de programación o falta de información desde la forma de ser alimentada por los datos.

Por su parte, en su proceso de pre-entrenamiento, la también IA generativa *ChatGPT* no se rige estrictamente por principios éticos y, de hecho, no distingue entre las cosas “buenas” y las “malas”, entre la verdad y la falsedad, ya que su procesamiento tiene como base los datos algorítmicos de lógica simbólica y proposicional, que son convertidos en números obtenidos y procesados de información proveniente de internet y otras bases de datos (UNESCO, 2023); ello constituye lo que se puede designar como un *sesgo cognitivo digital* de las IA generativas. El uso cotidiano y exponencial de dichas tecnologías a nivel global, supondría un nuevo salto tecnológico y social a gran escala, del cual aún desconocemos sus impactos.

Descripciones como las anteriores muestran que el entusiasmo inicial por los descubrimientos e innovaciones tecnocientíficas no siempre contempla las posibles consecuencias no deseadas de su aplicación práctica. Por ejemplo, en el proceso de transferencia de la nueva tecnología, se encuentra una serie de obstáculos y consecuencias no previstas, que incluyen desde una fuerte restricción en el campo de su aplicación, pasando por un uso diferenciado por parte de los consumidores, hasta aspectos de desigualdad y exclusiones tecnológicas (Shi, 2023).

Rakowski, Polak y Kowalikova (2021) afirmaron que dichas consecuencias no previstas pueden compensarse por efectos positivos en áreas particulares, como la económica o la política, siempre y cuando no entren en conflicto directo con algún tipo de normativas legislativas o de orden social. Los mismos autores también señalaron que “cuanto más rápido se entiendan los desarrollos tecnológicos y lo más importante, el papel social que juegan, más cuidadosamente nos prepararemos para sus impactos en la vida de un individuo o del sistema social” (p. 197).

Con base en lo anterior, podemos afirmar que la existencia de fuertes diferencias entre nuevas tecnologías y sus posibilidades están influenciadas por factores económicos, ideológicos, culturales, legales y organizacionales, entre aquellos que crean la nueva tecnología y los que la utilizan. Profundizar en esta serie de influencias ayudaría a comprender tanto las externalidades positivas como las negativas de la tecnología.

Como sistemas cerrados, opacos o de caja negra, los de IA están diseñados por grupos de trabajo y por personas con sus propias visiones del mundo, prejuicios, valoraciones de los hechos y sesgos adquiridos a lo largo de su experiencia de vida. Éstos se filtran en el diseño y en la definición de criterios de evaluación para los modelos algorítmicos (Murphy, 2012), con lo que se puede decir que, si esos grupos de trabajo no son lo suficientemente “diversos” e “inclusivos” como para reflejar una amplia variedad de visiones, muy probablemente no lleguen siquiera a darse cuenta de la existencia de los sesgos y, por tanto, a corregirlos (Ferrante, 2021, p. 35). No obstante, es importante señalar que para la UNESCO, los humanos somos moralmente responsables de la IA, cualquiera que se crea, implementa y/o usa. En investigaciones recientes se constata que ante alguna controversia de tipo legal los *softwares*, así como cualquier otro sistema inteligente, no pueden tener un carácter de sujeto legal.

En investigaciones como las antes presentadas, se apela a no perder de vista que la mayoría de las veces los programadores desconocen fuertemente los alcances de los resultados de sus innovaciones. Como desarrolladores de la tecnología, las personas u organizaciones que crean estos sistemas son responsables de que los algoritmos se alimenten de los conjuntos de datos, pero “no solo los conjuntos de datos pueden estar sesgados, también el algoritmo en sí” (Henz, 2021, p. 6).

Fue en los inicios del siglo XXI cuando se comenzó a decidir que los puestos de trabajo podrían ser sustituidos por computadoras, que las implicaciones éticas comenzaron a visibilizarse; es decir, al asumir a la IA como un campo de la informática centrado en el desarrollo de sistemas informáticos que funcionan como humanos, y al comprobar que las capacidades de ésta podrían sobrepasar las humanas, la situación se tornó adversa.

Efectos como el desempleo, el requerimiento de mayor capacitación o formación en nuevas profesiones, como en el uso de las Tecnologías de Información y Comunicación (TIC), así como las ciencias de la computación e información automatizada, opina García-Vigil (2021), fueron algunos de los motivos que crearon un “conflicto en la convivencia” de la inteligencia humana con la IA.

En términos institucionales, a nivel global, la Organización para la Cooperación y el Desarrollo Económicos (OCDE) ha definido a la IA como “un sistema computacional que puede, para un determinado conjunto de objetivos definidos por humanos,

hacer *predicciones y recomendaciones* o tomar decisiones que influyen en entornos reales o virtuales. Los sistemas basados en algoritmos de IA están diseñados para operar con distintos niveles de autonomía” (Cabrol *et al.*, 2020, p. 10). Como lo veremos a lo largo del artículo, por medio de ejemplos específicos, las condiciones tecnocientíficas y socioculturales en las que se generan dichas predicciones y recomendaciones autónomas serán puntos centrales de los actuales debates académicos en torno a los efectos y externalidades negativas de la IA.

El aporte de la presente investigación radica principalmente en la necesidad de un mayor entendimiento, tanto en el ámbito humano como en los avances de la IA, lo cual es fundamental para la comprensión profunda de su integración en la vida social, en procesos como el trabajo o la impartición de justicia y los beneficios surgidos de la colaboración entre la mente humana y la artificial. Además, permitiría reforzar la posición central del ser humano, encaminar los progresos hacia su mejora integral y fortalecer los elementos que contribuyen con su esencia: solidaridad, búsqueda del bienestar colectivo, acceso al conocimiento y creatividad (Barrios, Pérez y Guerra, 2020).

Asimismo, se busca contribuir en los debates que analizan las implicaciones de la opacidad y cerrazón de los agentes inteligentes, estudiando a las tecnologías como no intrínsecamente neutrales, donde el diseño “recrea” un tipo de ideología, aquella de quienes las desarrollan (Bowles, 2018). Lo anterior supone que en el diseño de una tecnología deberían considerarse disposiciones éticas, pues a partir de éste incluso se puede cambiar la forma en que vemos el mundo, cómo podemos actuar dentro de él (Bowles, 2018) y, de hecho, se podrían reestructurar las creencias sobre cómo debemos convivir con objetos y entornos que las personas usarán y habitarán. Consecuentemente, el objetivo de este trabajo es analizar el concepto y aplicación de la IA a la luz de los posibles sesgos, discriminaciones y desigualdades al momento de ser implementada en la sociedad, específicamente dando seguimiento a tres temas recurrentes en la literatura especializada contemporánea: los planteamientos éticos de la IA, los debates de género en torno a la feminización de las IA y los desafíos en torno a la impartición de justicia basada en algoritmos inteligentes.

Metodología

Para conformar el trabajo, se realizó una búsqueda de artículos científicos en repositorios especializados como: EBSCO, JSTOR, Scopus y Google Scholar, utilizando palabras clave, como: *Artificial Intelligence*, *Social gap* + Inteligencia Artificial, *Artificial Intelligence + bias*, *Artificial Intelligence + exclusion + bias*. El periodo de análisis y elección de artículos fue de enero a abril del 2023, privilegiando materiales con no más de cinco años de antigüedad. En esa búsqueda se encontraron debates en relación con los avances de la IA y su impacto en áreas como la justicia penal, la salud, el género, el trabajo

y la vigilancia digital. Se eligieron, de forma preliminar, 60 textos que cumplieran con los requerimientos considerados previamente (se privilegiaron textos provenientes de Ciencias Sociales, Humanidades o Administrativas, de corte preferentemente cualitativo o no estadístico, así como textos que no fueran relatorías o documentos elaborados por empresas), donde se analizaron los títulos, las palabras clave, así como las consideraciones finales de cada uno. Por último, fueron seleccionados 32 artículos científicos para ser revisados a profundidad en la muestra final.

Discusión

Inteligencia Artificial (IA)

Especialistas han afirmado que, al hablar de los sesgos de la IA, es importante poner atención en los datos, los modelos y las personas para construir con ello, eventualmente, una IA más justa (Hagerty y Rubinov, 2019). Esto debido a que la implementación de los modelos y sistemas inteligentes puede acarrear múltiples efectos negativos y riesgos éticos asociados a la desinformación, discriminación de grupos o de individuos, vulneraciones a la privacidad y seguridad o al uso indebido de los datos recolectados, entre muchos otros.

Un ejemplo de lo anterior es lo que se expone en el documental “Prejuicio cifrado” (“*Coded Bias*”), dirigido por la cineasta Shalini Kantayya y estrenado en 2020, donde se narra cómo la Dra. Buolamwini tomó conciencia del sesgo racial existente en los algoritmos de reconocimiento facial, al tiempo que analiza sus consecuencias. La Dra. Buolamwini es una mujer afrodescendiente, especialista en informática, activista y fundadora de la Liga por la Justicia Algorítmica (*Algorithmic Justice League*), que hace algunos años descubrió que varios sistemas comerciales de reconocimiento facial diseñados por Amazon, IBM y Microsoft funcionaban mejor con el rostro de sus amigos blancos que con el suyo. A este tipo de tendencias se le ha denominado *sesgo algorítmico*, es decir: “son sistemas cuyas predicciones benefician sistemáticamente a un grupo de individuos frente a otro, resultando así injustas o desiguales” (Ferrante, 2021, p. 29).

El concepto de sesgo algorítmico surgió de la observación de tecnologías de reconocimiento facial y vigilancia predictiva (Govia, 2020), las cuales están reproduciendo desigualdades existentes, expandiéndose en gran medida hacia nuevas construcciones de tipo coloniales, que continúan colocando a las comunidades negras, indígenas y racializadas bajo la vigilancia dirigida por los Estados. De acuerdo con Singer (2019), los sesgos algorítmicos se encuentran en el terreno del debate de lo ético, el cual históricamente se conformó como un campo de pensamiento en ocasiones ajeno a los tecnólogos e ingenieros.

Se ha dado tanta relevancia a todo el tema, que recientemente se ha desarrollado un subcampo de estudio denominado *Machine Ethics* (ME), cuya preocupación son los agentes

inteligentes, las máquinas y su relación con las modificaciones éticas, imperantes en la realidad social. Sobre estas circunstancias, hemos encontrado que los sistemas de IA en ocasiones “promueven una objetividad sin responsabilidad pública, y sin el sesgo social incrustado en ellos” (Manasi *et al.*, 2022, p. 6); pero esto se complica por el hecho de que, aunque las máquinas inteligentes pueden identificar estrategias óptimas, es posible “que no puedan decir si una elección fue correcta o incorrecta” (Manasi *et al.*, p. 8). Estas reflexiones reintroducen las cuestiones de responsabilidad cuando se enfrentan a resultados dañinos de IA (Govia, 2020) y, al mismo tiempo, podrían promover nuevas acciones para nuevos conjuntos de precedentes legales, sociotécnicos, derechos y debates sobre la posición de las tecnologías por parte de aquellos considerados responsables del bien público.

Hagerty y Rubinov (2019) han puntualizado que las tecnologías impulsadas por IA poseen un patrón que tiende a profundizar las divisiones sociales y a incrementar la desigualdad social, particularmente entre los grupos históricamente desfavorecidos, marginados y vulnerables. Este patrón existe a escala mundial y sugiere que los países de ingresos bajos y medianos pueden ser más vulnerables a los impactos sociales negativos de la IA, por lo cual es menos probable que se beneficien de las ganancias concomitantes.

En los mismos debates éticos se presenta una fuerte vulnerabilidad de sistemas de IA debido a técnicas como el envenenamiento de datos (Barrios *et al.*, 2020). Este es un método de explotación en el que un actor social, como un programador, puede manipular los datos de entrenamiento de los algoritmos para alterar las decisiones de un sistema. Dicho proceso de envenenamiento puede traer consigo graves consecuencias que van de la mano con la violación a la privacidad de las personas y sus datos, o el riesgo de acceder, manipular y utilizar los datos personales en contra de los usuarios.

Todo lo anterior ha levantado preocupaciones entre organismos internacionales que buscan regular dichas tecnologías. Éticamente, lo relacionado con la propiedad y utilidad de los datos personales, resguardados y utilizados por empresas y agencias gubernamentales, podría profundizar los problemas de discriminación social y justicia (Lupton, 2016). Igualmente, este tipo de cuestiones ponen en evidencia la posibilidad de usar la IA para influir y limitar la libertad de las personas (Kalluri, 2020). Así, como se ha reiterado en varias ocasiones, el verdadero desafío ya no es la innovación tecnológica y digital o su acceso, sino la gobernanza de dichas tecnologías y su impacto en la vida social (Floridi *et al.*, 2018).

Finalmente, sobre este punto identificamos el trabajo de Arguelles y Amaro (2023), quienes evalúan una serie de preocupaciones éticas que surgen en torno a los *chatbots* de asesoramiento médico y de vacunación, implementados por el gobierno mexicano en el contexto de la COVID-19. En sus hallazgos, los autores describen cómo algunas de las

preocupaciones éticas entre los usuarios de dichas tecnologías son la transparencia, la rendición de cuentas y la privacidad, las cuales según su perspectiva “generan una falta de confianza por parte de la ciudadanía hacia los *chatbots*, que se ha traducido en un bajo nivel de uso” (p. 108).

Los debates de género en la feminización de la IA

Sutko (2019) afirmó que las relaciones de género se materializan también en la tecnología. De esta manera, las perspectivas que se enfocan en analizar las IA y su entendimiento con las relaciones de género podrían ofrecer un ejemplo de cómo estas divisiones se “naturalizan y reproducen” (p. 569) a través de la tecnología.

Se ha encontrado que asistentes virtuales como Siri, Cortana y Alexa (todas con nombres de mujeres o diosas), simbolizan un proceso de feminización de las IA, donde las divisiones del trabajo por género se normalizan por medio de la asociación de “feminidad con trabajo simbólico y comunicativo” (Sutko, 2019, p. 569), en el que estos dispositivos y asistentes virtuales están diseñados con personalidades sumisas (Manasi *et al.*, 2022). Perspectivas como la anterior ayudan a entender la existencia de una “domesticación de la IA”, que resulta de una “asociación con la feminidad como dócil, receptiva y cariñosa”.

Investigaciones a nivel internacional (UNESCO, 2019) demuestran, sin ambigüedad, que los sesgos de género que persisten en los conjuntos de datos, algoritmos y dispositivos de capacitación de la IA, tienen el potencial de propagar y reforzar estereotipos de género perjudiciales. Estos sesgos se manifestarían durante el desarrollo del algoritmo, el entrenamiento de los conjuntos de datos o mediante la toma de decisiones generada por la IA (Manasi *et al.*, 2022), y podrían llegar a estigmatizar aun más a las mujeres, con el peligro de quedar relegadas en varios ámbitos de la vida económica, política y social, y así retrasar el progreso en materia de igualdad de género.

Investigaciones contemporáneas han resaltado el aumento de una “robotización de la sociedad y un antropomorfismo de lo tecnológico” (Dobrovestnova, Hannibal y Reinboth, 2022, p. 2), aunado a una humanización de la IA (Isola, 2020). No obstante, ciertos sectores y trabajos que se perciben como pertenecientes al “terreno de las mujeres”, como hotelería, turismo, comercio minorista, atención médica y educación, son más propensos a enfatizar posibles sesgos de género; por ejemplo, en el año 2015, Japón inauguró su primer hotel con recepcionistas únicamente de robots mujeres (Zimmerman, 2018). En consecuencia, entenderemos a los robots humanoides como “artefactos culturales”, imbuidos de “características antropomórficas”, los cuales tienden a reflejar las preferencias y suposiciones de sus creadores (Lamola, 2021, p. 120).

Similar a lo observado en las IA de asistencia virtual, a los robots mujeres se les asigna una variedad de tareas sociales (Dobrosovestnova *et al.*, 2022), como saludar a los clientes, proporcionar información o mantener un diálogo, tareas que han sido caracterizadas como un tipo de trabajo afectivo y relacionando directamente con las características femenino-maternales. Por su parte, Manasi *et al.* (2022) profundizan con mayor detalle sobre estos trabajos afectivos: en éstos, casi siempre relacionados con el sector de los servicios, las emociones que puedan llegar a mostrar los robots humanoides son percibidas como recursos que se capitalizan, lo que a su vez muestra una similitud con la forma en que se trata a la fuerza laboral feminizada de dicho sector.

Para ahondar en lo dicho anteriormente, se encontró el caso paradigmático de Sophia, el robot humanoide desarrollado por Hanson Robotics. En su fabricación se hizo un especial hincapié en hacerla lucir “excepcionalmente atractiva” (Manasi *et al.*, 2022, p. 9) y evocar un sentimiento de “mecánico-erotismo”, basado particularmente en la investigación existente en el campo de la interacción humano-robot o Human-Robot Interaction (HRI), el cual ha estudiado el impacto de los robots y sus reacciones emulando lo humano.

De esta forma y al aceptar esta robotización de la sociedad, resaltada por Dobrosovestnova *et al.* (2022), la aparición de tecnologías transversales, ubicuas y con tanta penetración como la IA y la robótica, podría reforzar la brecha digital actual e introducir nuevas formas de exclusión; por ejemplo, el *PricewaterhouseCoopers* (PwC), una de las firmas de consultoría, planeamiento y asesoramiento legal y jurídico más importantes del mundo, afirmó que de los 15,7 billones de dólares en riqueza que la IA generará a nivel mundial para el 2030, el 70% corresponderá exclusivamente a China y a los Estados Unidos (Rao, s.f.). De hecho, el gobierno japonés estableció un plan para que, en 2025, todos los hogares adopten un “estilo de vida robótico” que implique una vida segura, cómoda y conveniente con la ayuda de máquinas complementarias (Lobel, 2022).

Entonces, una nueva brecha, la “brecha robótica” (Toboso y Aparicio, 2019, p. 174), cuyo origen proviene del checo *robot* (trabajo forzado) y *rabota* (servidumbre), y fue popularizada por el narrador y dramaturgo checo Karel Čapek (1890-1938) en su obra *R.U.R (Rossumovi Univerzální Roboti [Los Robots Universales de Rossum])*, escrita en 1920 y en la cual apareció por primera vez, podría vislumbrarse en las próximas décadas, trazada por la exclusión del acceso y uso de los dispositivos robóticos para personas y grupos que no tienen oportunidades para adquirirlos, en contraste con los especialistas, programadores y aquellos más integrados digitalmente. Consecuentemente, podrían surgir “nuevos parias ubicados en la periferia de los grandes datos”, con el riesgo de “ser ignorados y excluidos de las decisiones basadas en los datos acopiados por la IA y de la utilidad de la información generada” (Barrios *et al.*, 2020, p. 98), lo que acarrearía falta de justicia social, fuertes desigualdades y exclusiones de acceso y uso, que podrían potencializar las ya existentes.

Aspectos cuestionables de la IA en el ámbito profesional

La IA se desarrolla rápidamente, lo que supone importantes desafíos. Aspectos sensibles y controversiales son los jurídicos y éticos, que se encuentran en discusión entre los organismos internacionales y los gobiernos a nivel mundial (Madrigal, 2020). La construcción de normas y leyes referentes al desarrollo y uso de la IA está lejos de alcanzar el desarrollo tecnológico. Además, como lo señala Salazar (2020), ciertos países van más adelante que otros. Lo anterior se convirtió en un debate entre entidades globales y gobiernos locales, especialmente después de la implementación de ciertas herramientas inteligentes en la impartición de justicia. Esto no es algo novedoso, ya que existen agentes artificiales en el sistema de justicia penal, particularmente en regiones anglosajonas desde hace varios años.

No obstante, actualmente el sistema de justicia penal de los EE. UU. recurre cada vez más a herramientas algorítmicas que ayudan a aliviar la carga de administrar un sistema tan grande. Por ello, cualquier sesgo algorítmico sistemático en estas herramientas tendría un alto riesgo y eventualmente podría transformarse en una desventaja acumulativa. Algunos especialistas han observado el uso de algoritmos en la fase de sentencia y libertad condicional; por ejemplo, Angwin *et al.* (2016) encontraron historias que muestran cómo los sistemas inteligentes tergiversan los datos, aumentando los riesgos de reincidencia de los convictos al ser valorados negativamente por el sistema de evaluación de riesgos penales, denominado Correctional Offender Management Profiling for Alternative Sanctions (COMPAS). Este tema puede revisarse en la página web del Departamento de Correccional del Estado de Wisconsin (State of Wisconsin Department of Corrections, s.f.).

Bornstein (2017), por su parte, reafirma los postulados anteriores y muestra en su investigación importantes indicios de un tipo de racismo estructural inmerso en los algoritmos inteligentes. Investigaciones como éstas rastrean una serie de evidencias que “insinúan un sesgo racial sistemático en la estimación del riesgo; los convictos negros estaban siendo calificados como más peligrosos que los no negros, incluso cuando estos últimos tenían delitos más graves” (Osoba y Welser, 2017, p. 9).

En junio del 2020 se publicó un texto que causó polémica entre los estudiosos de la IA aplicada en el terreno de la impartición de justicia: “Predictive policing algorithms are racist. They need to be dismantled”, un artículo de divulgación y análisis del *MIT Technology Review* en el cual se afirma que “existe una falta de transparencia y los datos de capacitación sesgados, muestran que estas herramientas no son adecuadas para su propósito. Si no podemos arreglarlas, deberíamos deshacernos de ellas” (Heaven, 2020).

El autor señala que el problema radica en los datos de los que se alimentan los algoritmos. Los algoritmos predictivos son fácilmente sesgados por las tasas de arresto. Según su

análisis, con base en las cifras del Departamento de Justicia de EE. UU., se tiene más del doble de probabilidades de ser arrestado si el sospechoso es de tez negra que si su tez es blanca. En el mismo texto, se plantea que la principal justificación del uso de algoritmos para ciertos procesos policiales, judiciales y de justicia en Estados Unidos, es la creencia generalizada de que éstos son más objetivos que los humanos.

A partir de lo anterior, se puede señalar que estamos frente a “nuevas” formas de impartición de justicia, soportadas en el procesamiento de datos y en el que los procesos de decisión descansan en las herramientas de la IA, donde “los prejuicios humanos se han integrado porque los modelos de aprendizaje automático se entrenan con datos policiales sesgados”. Con ello, lejos de evitar el racismo, simplemente pueden “ser mejores para ocultarlo”, por lo que “sin los ajustes culturales necesarios, se reproducen y profundizan las discriminaciones estructurales” (Muradas, 2021, p. 105). Muchos críticos de estos sistemas ven estas herramientas como una forma de lavado de tecnología (*tech-washing*), donde una apariencia de objetividad cubre los mecanismos que perpetúan las desigualdades en la sociedad. El concepto describe afirmaciones falsas sobre las capacidades reales de la tecnología, que incluye desde transformaciones digitales exageradas hasta la adopción de un lenguaje de marketing engañoso (Chartered Banker Institute, 2023).

Uno de los grandes desafíos en los debates sobre los efectos sociales de estos sistemas será la creación de procesos de transparencia y de rendición de cuentas por parte de los programadores y desarrolladores de los algoritmos inteligentes (European Parliament Procedure, 2021; Shi, 2023). De esta forma, los usuarios y consumidores conseguirían mejores herramientas para su eventual protección. Crear procesos de transparencia y de rendición de cuentas será un paso positivo para abrir las cajas negras de la IA, y atacar directamente lo que se ha denominado *ceguera de taller*, término que se usa cuando algo nos resulta tan normal y cotidiano, que hace perder de vista las oportunidades y riesgos siempre presentes.

Consideraciones finales

Al inicio del trabajo se hizo la pregunta: ¿qué consecuencias éticas, morales y sociales implica la presente generación en Inteligencia Artificial? De esta forma, a lo largo del texto se ha dado voz a casos particulares que apelan a la necesidad de fortalecer y profundizar el conocimiento de *lo humano*, su experiencia y sus capacidades, para afrontar los nuevos avances y, a la vez, preguntarse: ¿cuál es el proyecto humano para la era digital? (Floridi *et al.*, 2018).

Dicho cuestionamiento permite enriquecer la comprensión profunda de la naturaleza humana, para con ello poder desarrollar los conocimientos necesarios sobre los proce-

sos, alcances y potencialidades de los sistemas de IA. Estos conocimientos podrían facultar tanto a sus desarrolladores como a los usuarios, con el fin de enfrentar desafíos y promover soluciones relacionadas con los límites en la gestión y configuración de vastos conjuntos de datos, así como en el funcionamiento de las nuevas máquinas dotadas de sistemas en IA (Barrios *et al.*, 2020; Kalluri, 2020). El imperativo anterior instaría a los desarrolladores de dichas tecnologías a “ser capaces de discernir el impacto positivo o negativo que están infligiendo a la sociedad” (Floridi *et al.*, 2018, p. 700), y a ser plenamente conscientes de las responsabilidades inherentes de los resultados que generan.

En este sentido, se hace imperativo comprender, tanto las “oportunidades que los avances de la IA ofrecen para mejorar la calidad de vida humana” como consolidar la capacidad de “mantener el control sobre estos avances y sus implicaciones” (Floridi *et al.*, 2018, p. 693). En un escenario de incertidumbre ética, sumado a un avance acelerado de las tecnologías digitales, reflexiones académicas como las expuestas en el texto contribuyen con el campo de conocimiento identificado y propician la discusión sobre la relación estrecha entre algoritmos-dominación, en la cual se podría imponer una suerte de autoridad algorítmica, donde las decisiones son tomadas por los diseñadores de software que desempeñan un papel preponderante en la configuración, en buena medida, de las oportunidades de vida de personas (Barrios *et al.*, 2020), debido a que los usuarios, en la mayoría de los casos, desconocen cómo utilizar y analizar los datos digitales, además de que se tiene casi nulo control sobre ellos.

Para finalizar, preguntas que pueden ser planteadas hacia el futuro, a manera de cierre, serían: ¿desde dónde se construye/constituye la ética en la tecnología o cuáles son los términos sociales y culturales en la construcción de tecnologías como la IA? ¿Quiénes son los responsables: los gobiernos, las empresas, los usuarios o los especialistas, por avanzar en las agendas regulatorias sobre la ética social de la IA? ¿Cómo pueden participar activamente los grupos más vulnerables en la reducción de tipo de exclusiones y sesgos tecnológicos?

Referencias

- Angwin, J., Larson, J., Mattu, S. y Kirchner, L. (2016, 23 de mayo). Machine Bias. There is software that is used across the county to predict future criminals. And it is biased against blacks. *ProPublica*. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Arguelles, E. y Amaro, M. (2023, 10 de noviembre). Preocupaciones éticas en el uso de inteligencia artificial, transparencia y derecho de acceso a la información. El caso de los chatbots en el gobierno de México, en el contexto de la COVID-19. *Estudios en Derecho a la Información*, 1(15), 85-111. <https://doi.org/10.22201/ijj.25940082e.2023.15.17472>.
- Barrios, H., Pérez, V. y Guerra, Y. (2020, diciembre). Subjetividades e inteligencia artificial: desafíos para “lo humano”. *Veritas*, 47, 81-107. <https://www.scielo.cl/pdf/veritas/n47/0718-9273-veritas-47-81.pdf>.

- Bornstein, A. M. (2017, 5 de diciembre). Are algorithms building the new infrastructure of racism? *Nautilus*. <https://nautil.us/are-algorithms-building-the-new-infrastructure-of-racism-236911/>.
- Bowles, C. (2018). *Future Ethics*. Now Next Press.
- Bucher, T. (2018). *If... Then: Algorithmic power and politics*. Oxford University Press.
- Cabrol, M., González, N., Pombo, C. y Sánchez, R. (2020, enero). *Adopción ética y responsable de la inteligencia artificial en América Latina y el Caribe*. Banco Interamericano de Desarrollo. <http://dx.doi.org/10.18235/0002169>.
- Chartered Banker Institute. (2023, 3 de febrero). *Dissecting Techwashing*. https://www.charteredbanker.com/resource_listing//dissecting-techwashing.html.
- Dobrosovstnova, A., Hannibal, G. y Reinboth, T. (2022, junio). Service robots for affective labor: a sociology of labor perspective. *AI & Society*, 37(2), 487–499. <https://doi.org/10.1007/s00146-021-01208-x>.
- Elliott, K., Price, R., Shaw, P., Spiliotopoulos, T., Ng, M., Coopamootoo, K. y Van Moorsel, A. (2021, junio). Towards an equitable digital society: Artificial Intelligence (AI) and Corporate Digital Responsibility (CDR). *Society*, 58, 179–188. <https://doi.org/10.1007/s12115-021-00594-8>.
- European Parliament Procedure. (2021). *Artificial intelligence: questions of interpretation and application of international law in so far as the EU is affected in the areas of civil and military uses and of state authority outside the scope of criminal justice*. <https://oeil.secure.europarl.europa.eu/oeil/popups/printficheglobal.pdf?id=710011&l=>.
- Ferrante, E. (2021, julio-agosto). Inteligencia artificial y sesgos algorítmicos. ¿Por qué deberían importarnos? *Nueva Sociedad*, 294, 27-37. <https://nuso.org/articulo/inteligencia-artificial-y-sesgos-algoritmicos/>.
- Floridi, L., Cows, J., Beltrametti, M., Chatila, R., Chazerand, P., Dignum, V., Luetge, C., Madelin, R., Pagallo, U., Rossi, F., Schafer, B., Valcke, P. y Vayena, E. (2018, diciembre). AI4People—An ethical framework for a good AI society: opportunities, risks, principles, and recommendations. *Minds and Machines*, 28(4), 689–707. <https://doi.org/10.1007/s11023-018-9482-5>.
- García-Vigil, J. (2021, abril). Reflexiones en torno a la ética, la inteligencia humana y la inteligencia artificial. *Gaceta Médica de México*, 157, 311-314. <https://www.scielo.org.mx/pdf/gmm/v157n3/0016-3813-gmm-157-3-311.pdf>.
- Govia, L. (2020, 30 de julio). Coproduction, ethics and Artificial Intelligence: a perspective from cultural anthropology. *Journal of Digital Social Research*, 2(3), 42-64. <https://doi.org/10.33621/jdsr.v2i3.53>.
- Granieri, M. (2023, 5 de marzo). ¿Qué es la Inteligencia Artificial Generativa? *OBS Business School*. <https://www.obsbusiness.school/blog/que-es-la-inteligencia-artificial-generativa>.
- Hagerty, A. y Rubinov, I. (2019, 18 de julio). Global AI Ethics: A review of the social impacts and ethical implications of Artificial Intelligence. *Computer and Society*, 1-27. <https://doi.org/10.48550/arXiv.1907.07892>.
- Heaven, W. (2020, 17 de julio). Predictive policing algorithms are racist. They need to be dismantled. *MIT Technology Review*. <https://www.technologyreview.com/2020/07/17/1005396/predictive-policing-algorithms-racist-dismantled-machine-learning-bias-criminal-justice/>.
- Henz, P. (2021, 22 de septiembre). Ethical and legal responsibility for Artificial Intelligence. *Discover Artificial Intelligence*, 1(2), 1-10. <https://link.springer.com/article/10.1007/s44163-021-00002-4>.
- IBM. (s.f.). ¿Qué es un chatbot? <https://www.ibm.com/mx-es/topics/chatbots>.

- Isola, N. J. (2020, 11 de septiembre). Humanizar la Inteligencia Artificial. *Expansión*. <https://expansion.mx/opinion/2020/09/10/humanizar-la-inteligencia-artificial>.
- Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Jin, Y., Madotto, A. y Fung, P. (2022, 8 de febrero). Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), 1-47. <https://doi.org/10.48550/arXiv.2202.03629>.
- Kalluri, P. (2020, 7 de julio). Don't ask if Artificial Intelligence is good or fair, ask how it shifts power. *Nature*, 583(169). <https://www.nature.com/articles/d41586-020-02003-2>.
- Lamola, M. (2021, 11 de abril). The future of artificial intelligence, posthumanism and the inflection of Pixley Isaka Seme's African humanism, *AI & SOCIETY*, 37, 131-141. <https://doi.org/10.1007/s00146-021-01191-3>.
- Lobel, O. (2022, 22 de octubre). In Japan, humanoid robots could soon become part of the family. *Big Think*. <https://bigthink.com/the-future/the-equality-machine/>.
- Lupton, D. (2016, 15 de abril). The diverse domains of quantified selves: self-tracking modes and dataveillance. *Economy and Society*, 45(1), 101-122. <http://dx.doi.org/10.1080/03085147.2016.1143726>.
- Madrigal, A. (2020). América Latina busca su propia ruta hacia la Inteligencia Artificial (IA). En W. Weck y L. Salazar (eds.), *Inteligencia Artificial en Latinoamérica* (pp. 105-131). Fundación Konrad Adenauer. <https://dialogopolitico.org/wp-content/uploads/2023/04/Inteligencia-Artificial-en-Latinoamerica.pdf>.
- Manasi, A., Panchanadeswaran, S., Sours, E. y Ju, S. (2022, 8 de noviembre). Mirroring the bias: gender and Artificial Intelligence. *Gender, Technology and Development*, 3(26), 295-305. <https://doi.org/10.1080/09718524.2022.2128254>.
- Muradas, D. (2021, julio-agosto). Inteligencia Artificial: el derecho y el revés. *Nueva Sociedad*, 294. <https://nuso.org/articulo/inteligencia-artificial-el-derecho-y-el-reves/>.
- Murphy, K. (2012). *Machine Learning: A Probabilistic Perspective*. MIT Press.
- Osoba, O. y Welsler, W. (2017, enero). An Intelligence in our image: the risks of bias and errors in Artificial Intelligence. *RAND Corporation*. https://www.rand.org/pubs/research_reports/RR1744.html.
- Pinch, T. (2009). The Social Construction of Technology (SCOT): The Old, the New and the NonHuman. En P. Vannini (ed.), *Material Culture and Technology in Everyday Life: Ethnographic Approaches* (pp. 45-58). Peter Lang.
- Pinch, T. y Bijker, W. (1989). The social construction of facts and artifacts: or how the Sociology of Science and Sociology of Technology might benefit each other. En W. Bijker, T. Hughes y T. Pinch. (eds.), *The Social Construction of Technological Systems*. MIT Press.
- Polak, P. (2021, 14 de junio). Welcome to the Digital Era—the Impact of AI on Business and Society. *Society*, 58, 177-178. <https://link.springer.com/article/10.1007/s12115-021-00588-6>.
- Rakowski, R., Polak, P. y Kowalikova, P. (2021, 25 de mayo). Ethical aspects of the impact of AI: the status of humans in the Era of Artificial Intelligence. *Society*, 58, 196-203. <https://www.springerprofessional.de/en/ethical-aspects-of-the-impact-of-ai-the-status-of-humans-in-the-/19200602>.
- Rao, A. (s.f.). Sizing the prize. PwC's Global Artificial Intelligence Study: Exploiting the AI Revolution. What's the real value of AI for your business and how can you capitalise? *PWC*. <https://www.pwc.com/gx/en/issues/data-and-analytics/publications/artificial-intelligence-study.html>.
- Rezaev, A. y Tregubova, N. (2018, noviembre). Are sociologists ready for "artificial sociality"? Current issues and future prospects for studying Artificial Intelligence in the Social Sciences. *The Monitoring of Public Opinion Economic & Social Changes*, 147(5), 91-108. <https://doi.org/10.14515/monitoring.2018.5.10>.

- Rifkin, J. (2015). *La sociedad de costo marginal cero. El internet de las cosas, el procomún colaborativo y el eclipse del capitalismo*. Paidós.
- Salazar, L. (2020). Inteligencia Artificial: una oportunidad mundial. En W. Weck y L. Salazar (eds.), *Inteligencia Artificial en Latinoamérica* (pp. 11-28). Fundación Konrad Adenauer. <https://dialogopolitico.org/wp-content/uploads/2023/04/Inteligencia-Artificial-en-Latinoamerica.pdf>.
- Shi, Z. (2023, 25 de julio). *Learning and planning towards AI for social good* [Tesis de Doctorado, Carnegie Mellon University, Software and Societal Systems Department School of Computer Science]. Repositorio. https://kilthub.cmu.edu/articles/thesis/Learning_and_Planning_Towards_AI_for_Social_Good/23681178.
- Singer, N. (2019, 3 de septiembre). When apps get your medical data, your privacy may go with it. *The New York Times*. <https://www.nytimes.com/2019/09/03/technology/smartphone-medical-records.html>.
- State of Wisconsin Department of Corrections. (s.f.). COMPAS. <https://doc.wi.gov/Pages/AboutDOC/COMPAS.aspx>.
- Sued, G. (2022, 31 de agosto). Culturas algorítmicas: conceptos y métodos para su estudio social. *Revista Mexicana de Ciencias Políticas y Sociales*, 246, 43-73. <https://doi.org/10.22201/fcpys.2448492xe.2022.246.78422>.
- Sutko, D. (2019, 26 de septiembre). Theorizing femininity in Artificial Intelligence: a framework for undoing technology's gender troubles. *Cultural Studies*, 34(4), 567-592. <https://doi.org/10.1080/09502386.2019.1671469>.
- Techjury. (2021, mayo). *101 Artificial Intelligence Statistics*. <https://techjury.net/blog/ai-statistics/#gref>.
- Toboso, M. y Aparicio, M. (2019, 31 de mayo). Entornos de funcionamientos robotizados. ¿Es posible una robótica inclusiva? *Dilemata*, 30, 171-185. <https://www.dilemata.net/revista/index.php/dilemata/issue/view/31>.
- UNESCO. (2019). *Preliminary study on the Ethics of Artificial Intelligence*. <https://unesdoc.unesco.org/ark:/48223/pf0000367823.locale=es>.
- UNESCO. (2023). *ChatGPT e inteligencia artificial en la educación superior: guía de inicio rápido*. https://unesdoc.unesco.org/ark:/48223/pf0000385146_spa.locale=es.
- Wright, N. (2018, 10 de julio). How Artificial Intelligence will reshape the global order. *Foreign Affairs*, 5(22). <https://www.semanticscholar.org/paper/How-Artificial-Intelligence-Will-Reshape-the-Global-Wright/2c6057a67191fc65a66f6e-485166da79c1c06d4a>.
- Zimmerman, M. (2018). *Teaching AI: Exploring new frontiers for learning*. International Society for Technology in Education.

Agradecimientos

El presente artículo es resultado de la investigación *Inteligencia Artificial para el desarrollo y bienestar social: seguimiento, reconstrucción y evaluación antropológica de iniciativas mexicanas*, llevado a cabo en el Instituto de Investigaciones Sociales de la UNAM. Un agradecimiento por el apoyo y la asesoría a la Dra. María Josefa Santos Corral y por el financiamiento al Consejo Nacional de Humanidades, Ciencias y Tecnologías en su convocatoria de Estancias Posdoctorales en México 2022(3), con número 3752006.
