



Impacto de los sesgos en software informático basado en Inteligencia Artificial

Impact of biases on AI-based computer software

Autores

* **Freddy Aníbal Jumbo Castillo** 

✉ fjumbo@utmachala.edu.ec

Johnny Paul Novillo Vicuña 

✉ jnovillo@utmachala.edu.ec

Camilly Yuliana Pacheco Ordoñez 

✉ cpacheco5@utmachala.edu.ec

Joselyn Katiuska Franco Avila 

✉ jfranco9@utmachala.edu.ec

Universidad Técnica de Machala,
Facultad de Ingeniería Civil, Carrera de
Tecnologías de la Información, Machala,
El Oro, Ecuador.

*Autor para correspondencia

Comó citar el artículo:

Jumbo Castillo, F.A., Novillo Vicuña, J.P., Pacheco Ordoñez, C.Y. & Franco Avila, J.K. 2025. Impacto de los sesgos en software informático basado en Inteligencia Artificial. *Informática y Sistemas*. 9(1), 41-51. <https://doi.org/10.33936/isrtic.v9i1.7406>

Enviado: 22/03/2025

Aceptado: 07/04/2025

Publicado: 07/04/2025

Resumen

El objetivo de la investigación es analizar los sesgos y su impacto en diversos campos relacionados con: La selección de personal, el reconocimiento facial, la predicción de reincidencia, el diagnóstico médico y la evaluación crediticia. Su presencia afecta a la efectividad y precisión de las aplicaciones basadas en Inteligencia Artificial (IA) favoreciendo las diferencias sociales. Para tal propósito se utilizó el método PRISMA, permitiendo la búsqueda de aportes teóricos, el filtrado de acuerdo con el cumplimiento de criterios de consulta y la selección de artículos científicos relevantes. Los hallazgos del trabajo revelan lo siguiente: Las aplicaciones de reconocimiento facial evidencian sesgos raciales debido al uso de datos desbalanceados; los sistemas de reclutamiento y predicción de reincidencia son susceptibles a errores en la identificación de género y raza, a raíz de la implementación de algoritmos complejos que inciden en su entendimiento; las aplicaciones de análisis médico arrojan diagnósticos clínicos fallidos a causa del uso de técnicas inapropiadas que afectan a ciertos grupos de personas. Lo mencionado anteriormente demuestra la compleja relación entre la data seleccionada, la codificación de los algoritmos y los aspectos éticos que deben regir la puesta en marcha de las aplicaciones basadas en IA. Las nuevas contribuciones científicas deben centrarse en la investigación de métricas de desempeño utilizando datos balanceados y desbalanceados, conjuntamente con las técnicas requeridas para corregir desigualdades presentes en los volúmenes de datos.

Palabras clave: sesgos; equidad algorítmica; sistemas informáticos; aprendizaje automático; inteligencia artificial.

Abstract

The objective of the research is to analyze biases and their impact in various fields related to personnel selection, facial recognition, recidivism prediction, medical diagnosis, and credit evaluation. Their presence affects the effectiveness and accuracy of AI-based applications, thus exacerbating social differences. To this end, the PRISMA method was used, allowing for the search of theoretical contributions, filtering according to consultation criteria, and selecting relevant scientific articles. The findings of the study reveal the following: facial recognition applications exhibit racial biases due to the use of unbalanced data; recruitment and recidivism prediction systems are susceptible to errors in gender and race identification as a result of implementing complex algorithms that influence their understanding; medical analysis applications deliver faulty clinical diagnoses due to the use of inappropriate techniques that adversely affect certain groups of people. The aforementioned points demonstrate the complex relationship between the selected data, algorithm coding, and the ethical aspects that should govern the implementation of AI-based applications. New scientific contributions should focus on researching performance metrics using both balanced and unbalanced data, along with the techniques required to correct inequalities present in data volumes.

Keywords: biases; algorithmic fairness; computer systems; machine learning; artificial intelligence.



1. Introducción

Las Tecnologías de la Información y la Comunicación (TIC), han evolucionado brindando nuevas herramientas, métodos y técnicas sofisticadas, las cuales han mejorado la vida de las personas a través del uso de aplicaciones basadas en Inteligencia Artificial (IA) orientadas a la salud, bancas, laboral, reconocimiento facial, entre otros. Sin embargo, persisten desafíos importantes relacionados con los sesgos que influyen negativamente en los resultados, lo cual afecta la toma de decisiones. Los sesgos reproducen desigualdades que disminuyen la confianza en los sistemas, lo cual afecta su funcionalidad y opaca los beneficios que la IA ofrece (DeCamp & Lindvall, 2023). En el aporte científico de (Tsamados et al., 2022), manifiestan que la ética está vigente en el desarrollo de algoritmos de IA, puesto que reconoce los argumentos moralistas que se deben manifestar en las herramientas de software.

La urgencia de abordar los sesgos presentes en los sistemas de IA, responde a la necesidad de promover entornos tecnológicos justos e inclusivos. Dichos sesgos no solo reproducen patrones de discriminación preexistentes, sino que también pueden intensificar las desigualdades sociales, especialmente en contextos que afectan a poblaciones vulnerables. Conscientes de este desafío, organismos internacionales como la Comisión Europea han establecido marcos regulatorios orientados a mitigar estos riesgos. En su documento *Ethics Guidelines for Trustworthy AI*, se subraya que la diversidad, la no discriminación y la equidad constituyen principios esenciales para el desarrollo de una IA confiable.

Estas orientaciones evidencian un compromiso ético con la creación de tecnologías que respeten los derechos fundamentales y favorezcan la equidad en el ámbito digital (Varona & Suárez, 2022).

Los sesgos son el reflejo de una variedad de factores, que incluyen el diseño, la selección y la administración de repositorios de datos (Seyyed-Kalantari et al., 2021). La IA favorece la gestión de procesos operativos en las organizaciones, buscando ser neutral en sus decisiones incorpora técnicas sofisticadas para mitigar las inequidades mediante el uso de datos balanceados (Bagga & Piper, 2020), esto permite que las desigualdades históricas no afectan a grupos vulnerables (Vela et al., 2022). La mitigación del impacto de los sesgos incide para el logro de sistemas confiables al servicio de las personas (Tang et al., 2023), lo que conlleva que para tal propósito deban considerarse aspectos técnicos y socioculturales (Simonetta et al., 2021).

Son varios los estudios relacionados con el tema de investigación, entre los cuales se destacan los aportes de

(Peng, 2023) y (Ferrara, 2024); el primero enfatiza en las diferencias de rendimiento de los modelos utilizados en aplicaciones de reconocimiento facial, lo cual se refleja en patrones consistentes de sesgo racial; el segundo trabajo señala que el uso de inadecuadas técnicas ocasiona que se afecten a grupos minoritarios, lo cual es fundamental corregir para el desarrollo de sistemas justos y equitativos.

Las contribuciones son valoradas desde la perspectiva de cómo abordar el desbalance de los datos y la estructura de algoritmos complejos, con el fin de obtener resultados equitativos.

Los trabajos desarrollados muestran un progreso en la detección de sesgos, la completa mitigación de los mismos sigue siendo algo distante, debido al aumento en los volúmenes de datos que son utilizados para el entrenamiento y evaluación, generando la necesidad de garantizar la equidad y minimizar su impacto.

Para tal propósito, se requiere un trabajo más meticuloso y una validación más sólida en contextos reales, lo cual plantea el desafío de encontrar un equilibrio entre las métricas de rendimiento utilizadas. Una correcta planificación incide en el desarrollo y resultados de proyectos con IA, por lo cual es importante abordar de forma concreta cada etapa establecida.

Ante la limitada literatura especializada, que establezca de forma rigurosa y respaldada, la relación entre los sesgos presentes en el software basado en IA y sus implicaciones tanto operativas como sociales en contextos reales de implementación, existen investigaciones que identifican sesgos o proponen soluciones de carácter técnico, siendo insuficientes los enfoques que examinan críticamente la forma en que dichos sesgos se incorporan, se ignoran o incluso se amplifican durante las etapas de diseño y despliegue de estos sistemas.

Por lo tanto, el estudio se distingue por ofrecer un diagnóstico integral que no solo identifi

El objetivo de esta investigación es analizar el impacto de los sesgos de IA, mediante la revisión de textos científicos que aporten en la fundamentación de su origen, el efecto en la toma de decisiones, así como las estrategias de mitigación.

El proceso investigativo se regirá bajo las directrices del método PRISMA, buscando sentar las bases para nuevos aportes que contribuyan al logro de sistemas justos y transparentes. Además, esta revisión contribuye a cerrar una brecha clave en el campo, al articular conocimientos técnicos y regulatorios con el fin de promover prácticas más responsables, transparentes y auditables en el desarrollo de herramientas inteligentes.

2. Materiales y Métodos

El desarrollo general la presente investigación tiene como base los siguientes métodos teóricos:

- **Método de Análisis-Síntesis:** Ayudó a identificar las fuentes de sesgo y a proponer soluciones separando y reorganizando el objeto de estudio.
- **Método Descriptivo:** Contribuyó en la recolección de información de los sesgos en sistemas basados en IA, lo cual permitió una redacción científica, estructurada y secuencial.

El estudio de los sesgos en sistemas informáticos se abordó teniendo como insumos principales los aportes científicos relacionados. Por lo tanto, la metodología se aplicó con base en las directrices PRISMA 2020 la cual se fundamenta en la Tabla 1.

Tabla 1. Selección de la metodología.

Fuente: Los autores.

Metodología usada	PRISMA 2020
Razón de elección	Las directrices de esta metodología permiten organizar la información de una forma estructurada y clasificar datos de manera precisa, siendo apropiado la obtención datos que se relacionen o ayuden a identificar el impacto de los sesgos en la IA.
Fases principales	Identificación, Cribado e Inclusión.

Gráficamente las fases definidas para la metodología, se pueden apreciar en la Figura 1.

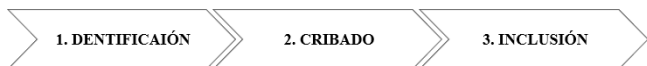


Figura 1. Fases de la metodología.

Fuente: Los autores.

En la Tabla 2, se describen las actividades definidas para cada etapa de la metodología:

Tabla 2. Actividades definidas para cada etapa de la metodología.

Fuente: Los autores.

Orden	Etapas	Actividades
1	Identificación	<ul style="list-style-type: none"> • Preguntas de investigación • Palabras claves • Selección de Bases de datos
2	Cribado	<ul style="list-style-type: none"> • Criterios de inclusión y exclusión • Selección de estudios
3	Inclusión	<ul style="list-style-type: none"> • Evaluación de relevancia • Extracción de datos

2.1. Identificación

2.1.1. Preguntas de investigación

Este estudio tiene como objetivo examinar los sesgos en sistemas de IA, enfocándose en sus causas, efectos y posibles soluciones. Dado que la IA se ha vuelto una herramienta indispensable en áreas como la toma de decisiones automatizada y la personalización de servicios; entender cómo surgen y se mantienen estos sesgos resulta importante. A través de investigaciones científicas, se busca identificar los principales tipos de sesgos en la IA y analizar su impacto en aspectos como la transparencia y la confiabilidad de sus algoritmos.

Tabla 3. Preguntas de investigación sobre sesgos en IA.

Fuente: Los autores.

Código de Pregunta	Pregunta de Investigación
PRG-1	¿Qué factores aportan en el surgimiento de sesgos en sistemas basados en IA?
PRG-2	¿Cómo afectan los sesgos a los usuarios que hacen uso de sistemas basados en IA?
PRG-3	¿Son efectivas las estrategias implementadas para la eliminación de sesgos?

Con base en este objetivo, se plantearon tres preguntas de investigación, las cuales se presentan en la Tabla 3:

Tabla 4. Términos clave empleados en la búsqueda de estudios.

Fuente: Los autores.

Términos	Contexto en el que se espera encontrarlo
Sesgos	Identificación de sesgos presentes en sistemas de IA.
Desafíos en algoritmos inteligentes	Análisis de los sesgos como desafíos en algoritmos y su incidencia en la toma de decisiones.
Filtros e inclusión	Explicación referente a las estrategias de filtrado de datos para mejorar la inclusión en sistemas de IA.
Controles y corrección de algoritmos basados	Descripciones sobre controles prácticos, para reducir la creación de sesgos algorítmicos en sistemas inteligentes y las medidas correctivas aplicadas.

2.1.2. Palabras clave y búsqueda bibliográfica

Para asegurar la relevancia de los estudios a considerar, se hizo uso de la selección de términos clave relacionados con la temática. A continuación, en la Tabla 4 se detallan los términos utilizados:

2.1.3. Selección de base de datos

La búsqueda de información se llevó a cabo en bases de datos científicas como:

- Scopus.
- Web of Science.



- Latindex.
- DOAJ.
- IEEE Xplore.
- SpringerLink.
- Google Scholar.
- Redalyc.

Tabla 5. Cadenas de búsqueda utilizadas en los motores de bases de datos.

Fuente: Los autores.

Idioma	Cadenas de Búsqueda
Inglés	Challenges OR barriers in intelligent algorithms AND their impact on decision-making AND ethical considerations OR optimization.
Inglés	Corrective measures to mitigate OR addressing algorithmic bias in artificial intelligence OR intelligent systems AND promote fairness OR transparency AND inclusion.
Español	Control de algoritmos O mitigación de sesgos Y brechas en sistemas inteligentes O automatizados.
Español	Algoritmos O sistemas inteligentes con desafíos O sesgos Y estrategias para su corrección O mitigación de fallas en decisiones.

Se utilizaron combinaciones de palabras clave incorporadas en cadenas de búsqueda, complementadas con operadores booleanos, con el objetivo de optimizar la precisión y relevancia de los estudios obtenidos en las bases de datos consultadas. A continuación, en la Tabla 5 se presentan las principales cadenas de búsqueda empleadas:

Además, los resultados fueron filtrados considerando que las fechas de publicación estén entre enero de 2020 y marzo de 2025, esperando contar con los puntos de vista más actualizados.

2.2. Cribado

2.2.1. Criterios de inclusión y exclusión

La filtración de los estudios importantes se dio en función de los argumentos de inclusión y exclusión, lo cual permitió focalizar el análisis en aportes teóricos de calidad relacionados con el objetivo del trabajo. A continuación, se listan los criterios establecidos:

- Criterios de inclusión
 - * Artículos científicos vigentes en revistas indexadas.
 - * Rango temporal de publicación desde el año 2020 al 2025.

Esto avala la actualidad de los descubrimientos científicos.

- * El contenido de las contribuciones científicas, se deben relacionar con los sesgos en sistemas informáticos basados en IA. Esto involucra aspectos tales como: Género, raza, cultura y factores socioeconómicos.

- * Artículos científicos completos que permitan un análisis minucioso.

• Criterios de exclusión

- * Artículos científicos que no hayan sido revisados por pares a doble ciego.

- * Publicaciones científicas que en su resumen bibliográfico no cuenten con un DOI (Identificador de Objeto Digital).

2.2.2. Selección de estudios

El proceso realizado tuvo como consecuencia un total de 58 aportes teóricos identificados en las bases de datos científicas. Sin embargo, se excluyeron 10 por condiciones de duplicidad y 5 por otros factores, quedando para su valoración 44 textos

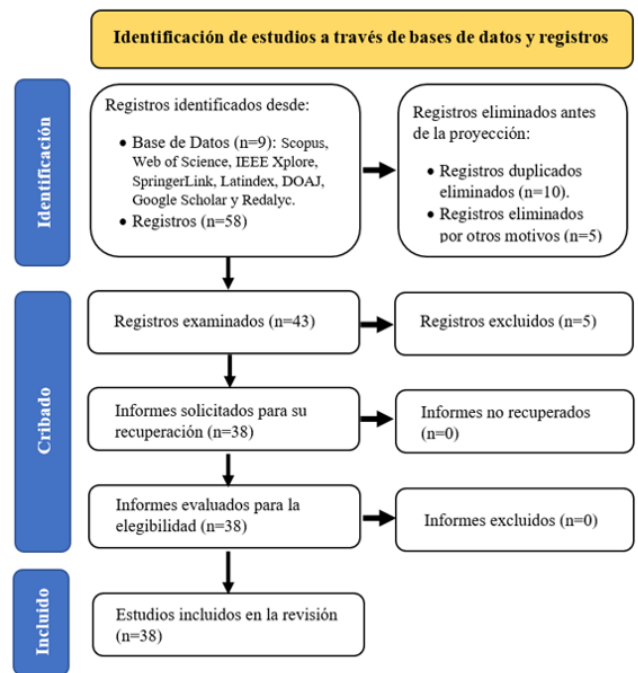


Figura 2. Diagrama PRISMA.

Fuente: Modificado de Santiago Arenas et al. (2023).

científicos. En la fase de cribado se descartaron 5, resultando en 38 estudios estimados para establecer su elegibilidad.

En la etapa final no se presidió de ningún trabajo, por lo cual el

desarrollo de la investigación se sustenta en las 38 contribuciones académicas. En la Figura 2, se representa gráficamente el proceso de acuerdo con la metodología PRISMA.

2.3. Inclusión

2.3.1. Evaluación de relevancia

La revisión sistemática de un trabajo científico se debe sustentar en contribuciones teóricas, que avalen resultados de calidad. Ante este desafío, se definieron los criterios de valoración considerando la relación que debe existir entre el tema planteado y los artículos cuyo contenido se corresponda con los sesgos en el software basado en IA. Por lo tanto, en la estructura de cada texto académico se analizaron los siguientes aspectos:

- Metodología utilizada.
- Alcance de los estudios con relación a la temática de investigación, específicamente en campos de aplicación tales como: selección de personal, reconocimiento facial, predicción de reincidencia, diagnóstico médico y evaluación crediticia.
- Rigor académico.

Tabla 6. Sistema de puntuación y resultados para la selección de estudios.

Fuente: Los autores.

Aspecto Evaluado	Detalle
Método de puntuación	Los cinco criterios considerados son: (1) definición clara del objetivo del estudio, (2) transparencia en la descripción metodológica, (3) profundidad en el tratamiento de los sesgos, (4) integración con marcos éticos o normativos pertinentes, y (5) aplicabilidad de los hallazgos en contextos reales de implementación. Cada criterio fue calificado con una escala ordinal de tres niveles: 0 (no cumple), 1 (cumple parcialmente) y 2 (cumple plenamente), estableciéndose un puntaje máximo de 10 puntos por estudio evaluado.
Rango de puntuación	Los estudios con una calificación superior a 7 puntos se seleccionaron considerándolos como relevantes para la investigación.
Total de estudios seleccionados	Se seleccionaron 38 artículos científicos, que estuvieron dentro del rango de puntuación.

Rango temporal de publicación en revistas científicas entre el año 2020 y 2025.

- La relación que debe existir entre el objetivo planteado en la investigación y los aportes teóricos identificados.

Se estableció un método para la valoración y selección de los estudios, el cual se detalla en la Tabla 6.

2.3.2. Extracción de datos

Permitió organizar y clasificar datos en distintas categorías, a través del uso de un protocolo estructurado con el fin de identificar información relevante de los estudios seleccionados con anterioridad. Los aspectos analizados para la extracción de

Tabla 7. Aspectos analizados a partir de los resultados de la extracción de datos.

Fuente: Los autores.

Aspecto Analizado	Detalle
Categorías de sesgos	Desigualdades relacionadas con el género, etnia, condición socioeconómica o estatus social, creencias y aspectos culturales del usuario.
Campos de aplicación	Reconocimiento facial, aplicaciones de reclutamiento, sistemas de predicción de reincidencia, análisis clínico y evaluación crediticia.
Medidas para mitigar sesgos	Enfoques prácticos, implementación de controles y métodos de corrección diseñados para prevenir y corregir estos problemas en sistemas o modelos basados en IA.

datos se estructuran en la Tabla 7.

3. Resultados y Discusión

3.1. Resultados

La investigación revela que los sistemas informáticos basados en IA, proporcionan resultados inequitativos de acuerdo con el contexto en que se desarrollan, esto debido a que, en función de los volúmenes de datos se seleccionan erróneamente los conjuntos para el entrenamiento y evaluación de los modelos. En la revisión de los estudios se destacan cinco tipos de sesgos, los

Tabla 8. Principales tipos de sesgos.

Fuente: Los autores.

Tipo de sesgo	Descripción
Sesgo de datos	El uso inapropiado de técnicas de muestreo y balanceo de datos, conlleva a que los algoritmos arrojen resultados injustos para ciertos grupos de la población.
Sesgo de selección	Ocurre cuando la muestra utilizada para un estudio, no es representativa de la población general.
Sesgo de confirmación	Los algoritmos incrementan los estereotipos, cuando los datos de entrenamiento se encuentran desbalanceados y esconden patrones no deseados.
Sesgo de Evaluación	La utilización incorrecta de métricas de rendimiento, afecta el desempeño y confiabilidad de un modelo ocultando la presencia de sesgos
Sesgo de implementación	La discordancia entre el ambiente de entrenamiento e implementación afecta significativamente la efectividad y eficacia del modelo.

mismo que se describen en la Tabla 8 (Akter et al., 2021).

Las investigaciones revelan un crecimiento significativo en el desarrollo de aplicaciones que incorporan algoritmos de IA, las cuales han traído consigo retos relacionados con la presencia de sesgos, cuya existencia se debe al uso de datos de entrenamiento que marginan a ciertos grupos, incidiendo para el logro de una representación justa y equitativa (de Lima et al., 2023). A



continuación, se describen diversos algoritmos utilizados en aplicaciones basadas en IA, enfatizando en sus fortalezas y debilidades frente a los sesgos:

- **Redes Neuronales Convolucionales (CNN):** Los algoritmos son frecuentemente utilizados para la clasificación de imágenes, reconocimiento facial, identificación de objetos, entre otros. A pesar de sus potencialidades son susceptibles a experimentar sesgos, si se utilizan métodos poco eficientes para la selección de datos de entrenamiento y evaluación, lo cual incide en búsqueda de una representatividad justa (Paredes Meneses, 2023).

- **Transformers:** Estos modelos son fuertes en tareas de Procesamiento de Lenguaje Natural (PLN), los cuales se desenvuelven procesando grandes cantidades de datos, de cuales en ocasiones se heredan las falencias relacionadas con la poca variabilidad de contenido favoreciendo la presencia de sesgos.

- **Árboles de Decisión y Random Forests:** Estos modelos requieren que los datos de entrenamiento se definan correctamente, para una implementación efectiva de los algoritmos (Ghosal & Hooker, 2020).

- **Algoritmos de Clustering:** Los modelos desfavorecen a ciertos grupos, si los datos no se sistematizan adecuadamente.

Los sistemas modernos involucran el uso de tecnologías emergentes, con el fin de optimizar las tareas de las aplicaciones informáticas relacionadas con la seguridad, justicia, talento humano, entre otros. La IA surge como un mecanismo válido para optimizar procesos críticos; sin embargo, están sujetos a experimentar sesgos debido a la poca variabilidad de los datos y la codificación de algoritmos poco transparentes, lo cual influye en la obtención de resultados inequitativos y discriminatorios. A continuación, se detallan los casos que han experimentado sesgos:

- **Sesgo racial en reconocimiento facial:** Los modelos de IA utilizados muestran un bajo rendimiento en la identificación de personas de piel oscura, debido a su escasa representatividad en los datos de entrenamiento, lo cual se sustenta en los resultados que arrojaron las aplicaciones de IBM afectando en un 34.7% a las mujeres afrodescendientes (Khalil et al., 2020). Los algoritmos implementados para el reconocimiento facial son frecuentemente utilizados por la policía, para la identificación de sospechosos o personas con antecedentes penales; sin embargo, no están exentos de sesgos lo cual se traduce en resultados erróneos. Ante esta situación (Sanabria Moyano et al., 2022), mencionan que esto genera dudas relacionadas con la certeza y aporte al desarrollo de los casos. La codificación de algoritmos incide en la toma de decisiones, por lo cual es necesario que se sustenten bajo sólidos criterios éticos, neutrales e inclusivos (Reyes Campos et al., 2023). De acuerdo con el criterio de (Ávila Bravo-Villasante, 2023), las aplicaciones implementadas como

medios de atención donde existe masiva concurrencia, pueden afectar la funcionalidad de los sistemas debido a la presencia de sesgos raciales (Beltramelli Gula et al., 2023). Sin embargo, existen otras herramientas que han experimentado algún tipo de

Tabla 9. Sesgos en herramientas de reconocimiento facial.

Fuente: Los autores.

Empresa/Aplicación	Descripción
Amazon Rekognition	Esta aplicación arrojó resultados que afectaban la precisión, perjudicando en el reconocimiento a personas de piel oscura o afrodescendientes.
Amazon Rekognition	En el reconocimiento de individuos de piel oscura la cantidad de falsos positivos fue muy alta.
IBM Watson Visual Recognition	El rendimiento de la aplicación arrojó resultados que afectaban significativamente las personas de sexo femenino afrodescendientes.

sesgo y se describen en la Tabla 9:

- **Sesgo de género en sistemas de reclutamiento:** Los sistemas inteligentes para la selección de personal, demuestran marcada polaridad hacia los candidatos hombres, discriminando a las candidatas mujeres, por lo cual es importante que las aplicaciones

no sólo se adapten a los avances tecnológicos, sino que además consideren las características de los usuarios que comúnmente se estructuran en las hojas de vida. Esto aporta transparencia a las tareas de revisión de currículums, mitigando el impacto de los sesgos y optimizando los tiempos de respuesta. El descubrimiento de postulantes en línea facilita el proceso de evaluación de actividades o tareas, mediante el uso de aplicaciones modernas de realidad aumentada, lo cual favorece la planificación establecida reduciendo los tiempos requeridos. La presencia de sesgos de género en los modelos de IA para el reclutamiento de talento humano incide en la inclusión laboral, por lo cual es pertinente considerar en el desarrollo de estrategias que ayuden a minimizar su impacto. Esto contribuye a una selección transparente donde candidatos masculinos y femeninos tengan las mismas opciones laborales (Pérez López, 2023). En la Tabla 10, se evidencian algunos casos de sesgos detectados en los sistemas informáticos de varias empresas:

- **Sesgo en sistemas de predicción de reincidencia:** Estas aplicaciones tienen el propósito de predecir cuando una persona reincide en el cometimiento de delitos, lo que las ha convertido en una necesidad emergente y de exigencia en los requerimientos científicos, tecnológicos y justicia penal, existiendo una marcada relación entre la ciencia de datos, las matemáticas y el derecho, las cuales de forma conjunta permiten el desarrollo de sistemas

Tabla 10. Sesgos de género identificados en organizaciones.

Fuente: Los autores.

Empresa/ Aplicación	Descripción
Amazon	El sistema favoreció a los hombres y perjudicó a las mujeres, debido a que el modelo fue entrenado con datos históricos, que reflejaban escasa variabilidad en su contenido (Chang, 2023).
LinkedIn	Los algoritmos de la aplicación InstaJob y People You May Know (PYMK), experimentaron sesgos de género que valoraban de mejor forma a los candidatos masculinos, en comparación con las candidatas femeninas reduciendo sus oportunidades en la plataforma en línea (Yu & Saint-Jacques, 2022).
InfoJobs	En esta aplicación los candidatos masculinos recibían avisos laborales y contratos estables con mejores sueldos y prestaciones, en comparación con las candidatas de sexo femenino (Martínez et al., 2021).

transparentes y efectivos. Esto conlleva a superar los retos que surgen entorno a su implementación, aprovechando los beneficios de la tecnología emergente y con una base jurídica sólida, lo cual se refleje en algoritmos correctamente estructurados y equitativos, que resuelvan de manera justa lo concerniente a infracciones (Bravo Bolado, 2023). El Level of Service Inventory-Revised (LSI-R), es un instrumento de evaluación ampliamente empleado en el sistema de justicia penal para medir el riesgo de reincidencia y determinar las necesidades criminógenas de las personas evaluadas. Este modelo considera una variedad de factores, incluidos antecedentes penales, historial educativo, laboral, relaciones familiares y estabilidad emocional, con el objetivo de generar puntuaciones que sirvan de apoyo a jueces y funcionarios. Estas calificaciones son utilizadas para tomar decisiones informadas sobre libertad condicional, planes de rehabilitación y otras intervenciones judiciales, buscando promover estrategias que reduzcan la reincidencia y fortalezcan la reintegración social (Zhang, 2016).

Sin embargo, varios estudios han señalado que la implementación del LSI-R presenta desafíos, especialmente en lo que respecta a la interpretación de sus puntuaciones en diferentes contextos culturales y sociales. Aunque es una herramienta transparente y bien documentada, su aplicación puede generar inquietudes cuando los datos utilizados para las evaluaciones no representan adecuadamente a ciertos grupos demográficos, lo que podría perpetuar desigualdades existentes. Estas limitaciones subrayan la necesidad de realizar auditorías regulares y ajustes contextuales para garantizar que las evaluaciones sean justas y equitativas para todos los individuos. Otra aplicación análoga es Prisoner Assessment Tool Targeting Estimated Risk and Needs (PATTERN), la cual valora el riesgo de reincidencia de los sujetos privados de libertad, teniendo como objetivo la reinserción de las personas a la sociedad (Hamilton et al., 2022). Estas herramientas utilizan modelos estadísticos de IA para analizar datos y efectuar predicciones.

• **Sesgo en sistemas de diagnóstico clínico:** En estas aplicaciones existen desafíos por resolver relativos a la equidad. La codificación de algoritmos ha permitido avances importantes en el ámbito

de la salud; sin embargo, estos modelos tienden a presentar inconsistencias en las respuestas, cuando los datos utilizados no se han preprocesado correctamente para una representación efectiva de la población. Esto conduce a valoraciones erróneas en ciertos grupos étnicos y socioeconómicos, enfatizando la necesidad de reconocer y manejar sesgos humanos y tecnológicos. La contribución científica de (Kudina & de Boer, 2021), explora la integración de sistemas de apoyo a la decisión basados en el aprendizaje automático, considerando que la inclusión de sesgos incide en la interpretación del diagnóstico de enfermedades del paciente, donde el juicio clínico juega un papel crucial. El no condicionamiento a la atención médica, incorpora modelos neutrales de IA que descartan la influencia de sesgos en los procesos de análisis.

Una de las herramientas desarrolladas en el campo de la medicina y basadas en IA es Watson, la cual mostró una alta tasa de correlación con equipos multidisciplinares sobre ciertos tipos de cáncer, especialmente cuando se aplican en diferentes contextos clínicos y geográficos. Sin embargo, otros estudios han señalado discrepancias en las recomendaciones que arroja la herramienta, las cuales pueden ser causados por un sistema entrenado principalmente con datos de instituciones específicas, que podrían no reflejar la diversidad de prácticas médicas a nivel mundial (Jie et al., 2021). De acuerdo con el criterio de (Gilbert et al., 2020), la aplicación Ada tuvo una precisión mayor en casos clínicos simples y con enfermedades comunes; no obstante, presentó un sesgo algorítmico relevante que limita su eficacia para detectar patologías raras o atípicas. Este sesgo se atribuye a la priorización de enfermedades prevalentes en sus datos de entrenamiento.

En el aporte teórico de (Golden et al., 2021), destacan que los National Standards for Culturally and Linguistically Appropriate Services (CLAS), son una herramienta determinante para lograr un trato consistente en la atención de salud, por tal motivo, para resolver estas discrepancias es esencial abordar directamente los requerimientos sociales de la salud. Con base en lo manifestado anteriormente, es imperiosa la necesidad de optimizar la data para el entrenamiento y evaluación de los modelos bajo protocolos de validación permanente. Es importante que los sistemas consideren, desde la capacitación en sesgos y diversidad hasta estrategias de estado en temas de salud para lograr una equidad y sostenibilidad.

La funcionalidad de estas aplicaciones se logra incorporando criterios de transparencia y equidad, que permitan la evaluación de los solicitantes de forma equitativa, considerando que la utilización de algoritmos de IA en la estimación del crédito ha transformado la forma en que se conceden los préstamos y se gestionan los riesgos financieros. Por lo tanto, es necesario que las instituciones financieras optimicen sus modelos, con el fin de evitar la implementación de sistemas sesgados que reducen las oportunidades de financiamiento para los grupos marginados. La presencia de registros históricos en la data, producen resultados sesgados que se sustentan en una incorrecta interpretación de variables tales como: raza, género y nivel socioeconómico (Fuster et al., 2021). Para hacer frente a estos desafíos, es necesario optimizar las técnicas de muestreo y las limitaciones de equidad,



permitiendo el desarrollo de sistemas de IA transparentes y justos.

Cabe señalar que en la contribución científica de (Chen et al., 2025), manifiestan que en el desarrollo de sistemas de evaluación de crédito debe buscarse un equilibrio entre la equidad y la exactitud de los modelos de evaluación del riesgo; sin embargo, eliminar completamente los sesgos en los algoritmos puede ser una tarea desafiante. La correcta ejecución de las tareas de diseño, codificación del algoritmo de IA y una auditoría permanente, contribuyen a la obtención de resultados equilibrados en las

Tabla 11. Sesgos de evaluación crediticia en organizaciones.

Fuente: Los autores.

Empresa/ Aplicación	Descripción
Apple Card	Un diseño no transparente puede ocultar sesgos asociados, lo que demuestra la necesidad de una supervisión regulatoria para la asignación de créditos justos. La valoración de perfiles otorgó diferentes límites de crédito según el género, afectando a las solicitantes mujeres (Li et al., 2023).
FICO Score	Las valoraciones crediticias frecuentemente reflejan las diferencias históricas de datos heredados, lo cual conduce al desafío de actualizar los conjuntos de datos utilizados en los sistemas para lograr resultados equitativos (Coraglia et al., 2024).
Kreditech	Está aplicación incorpora la IA junto con la Big Data para la valoración crediticia; sin embargo, no estuvo exenta de presencia de sesgos los cuales fueron resultado de escasos antecedentes bancarios, lo cual afectó a ciertos grupos de personas (Wang & Wang, 2020).

aplicaciones de evaluación de créditos financieros, lo cual beneficia a todos los segmentos de la sociedad. A continuación, se detallan en la Tabla 11 los sesgos detectados en aplicaciones de evaluación crediticia.

3.4. Discusión

La revisión de los estudios consultados revela un interés académico creciente por comprender la naturaleza y el alcance de los sesgos algorítmicos en los sistemas basados en IA. No obstante, persisten limitaciones estructurales que obstaculizan una evaluación integral del fenómeno. Una de las debilidades más recurrentes es la marcada heterogeneidad metodológica: la multiplicidad de métricas utilizadas para evaluar la equidad y la discriminación, tal como lo advierten (Simonetta et al., 2021) (Yu & Saint-Jacques, 2022), se dificulta establecer comparaciones sistemáticas entre modelos y ámbitos de aplicación. A pesar de los trabajos enfocados en sesgos visibles, como los relacionados

con género o raza, existe una tendencia a subestimar factores interseccionales como la clase social, la discapacidad o la edad. Esta omisión restringe la profundidad analítica desde una perspectiva de justicia social (Ávila Bravo-Villasante, 2023) (Bravo Bolado, 2023), y evidencia la necesidad de establecer marcos regulatorios que respondan de manera integral a las implicaciones éticas derivadas del uso de IA.

Otra limitación relevante identificada en la literatura es la escasa replicabilidad de los estudios empíricos. Trabajos como los de (Gilbert et al., 2020), centrados en aplicaciones clínicas de diagnóstico automatizado, o (Martínez et al., 2021), enfocados en la detección de sesgos de género en sistemas de alerta de empleo, presentan resultados sin ofrecer información detallada sobre la procedencia de los datos ni sobre las tasas específicas de error. Esta falta de transparencia metodológica impide evaluar la robustez y la generalización de sus hallazgos. La opacidad técnica se vuelve aún más problemática en contextos críticos como la justicia penal o el ámbito financiero, donde las decisiones automatizadas pueden tener implicaciones significativas para los derechos de las personas (Noiret et al., 2021)(Fuster et al., 2021). Los sistemas de reconocimiento facial han sido objeto de creciente investigación debido a los sesgos, que tienden a afectar de forma desproporcionada a determinados grupos demográficos, lo cual genera discrepancias que comprometen, no solo la fiabilidad técnica de los sistemas, sino que también generan preocupaciones éticas relevantes, particularmente cuando se aplican en contextos de alta sensibilidad social como la seguridad pública y la vigilancia(Khalil et al., 2020).

Por otro lado, la literatura también pone en evidencia controversias persistentes en torno a la tensión entre eficiencia técnica y justicia distributiva. Estudios como los de (Cary et al., 2023) y (DeCamp & Lindvall, 2023), documentan cómo ciertos algoritmos clínicos perpetúan inequidades raciales al diagnosticar erróneamente y sistemáticamente a pacientes pertenecientes a minorías étnicas, a pesar de haber sido diseñados con el propósito explícito de optimizar la atención médica. En la misma línea, (Chang, 2023) expone el caso del sistema de reclutamiento automatizado de Amazon, el cual discriminó de forma sistemática a mujeres, revelando que incluso organizaciones con capacidades tecnológicas avanzadas pueden desarrollar sistemas profundamente sesgados. Esto refuerza la idea de que el sesgo en los sistemas de IA no puede entenderse exclusivamente como un problema técnico, sino como un fenómeno sociotécnico que está intrínsecamente ligado a las dinámicas sociales y a los datos que alimentan dichos sistemas (Akter et al., 2021) (Ferrara, 2024).

El análisis puede fortalecerse mediante la inclusión de estudios que no solo respalden, sino también desafíen o complementen los hallazgos más consolidados en la literatura. Así, mientras algunos trabajos destacan el carácter discriminatorio de determinados modelos de puntuación crediticia (Fuster et al., 2021), otros como

(Coraglia et al., 2024) introducen herramientas como BRIO, diseñadas para mejorar la transparencia y reducir el sesgo desde las etapas iniciales del desarrollo algorítmico. En el campo de la salud, (Cary et al., 2023) advierten sobre el impacto adverso de los algoritmos clínicos en poblaciones vulnerables; sin embargo, (Vela et al., 2022) sostienen que, bajo enfoques de diseño responsable la IA puede constituirse en una herramienta eficaz para promover la equidad sanitaria. Estas posturas divergentes evidencian la necesidad de fomentar un debate académico

Tabla 12. Estudios analizados según el dominio.

Fuente: Los autores.

Dominio	Número de estudios
Selección de personal	4
Reconocimiento facial	6
Predicción de reincidencia	4
Diagnóstico médico	6
Evaluación crediticia	6

más plural, que contemple tanto las limitaciones como las oportunidades de la IA en contextos sociales complejos.

Clasificación de los estudios analizados por área de aplicación

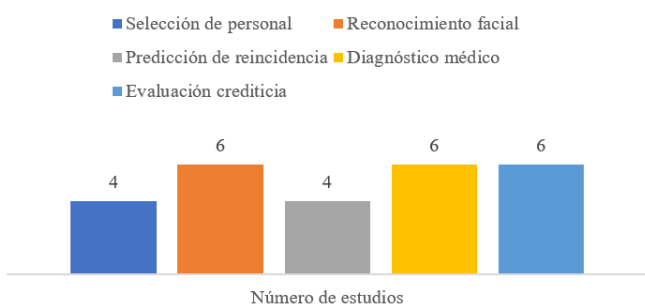


Figura 3. Clasificación de los estudios analizados por área de aplicación.

Fuente: Los autores.

Con el propósito de ofrecer una visión transversal del fenómeno del sesgo algorítmico en IA, se presenta en la Tabla 12, los estudios analizados según el dominio de aplicación.

A continuación, se representan gráficamente en la Figura 3, los resultados de la tabla anterior.

La gráfica revela que los dominios de reconocimiento facial, diagnóstico médico y evaluación crediticia concentran la mayor cantidad de estudios (seis en cada caso), lo cual sugiere una atención académica prioritaria hacia estos sectores, probablemente debido a su impacto directo en derechos fundamentales y en procesos de toma de decisiones sensibles.

En contraste, los ámbitos de selección de personal y predicción

de reincidencia muestran una menor cantidad de investigaciones (cuatro estudios en cada uno), lo que podría evidenciar la necesidad de una mayor profundización teórica y empírica en estas áreas, particularmente considerando los riesgos éticos y sociales que también conllevan sus aplicaciones algorítmicas.

4. Conclusiones

Los resultados de la investigación evidencian que los sesgos en el software basado en la IA, están latentes en varias aplicaciones que van desde el reconocimiento facial, diagnóstico médico, sistemas de evaluación del crédito, predicción de reincidencia y contratación de personal.

Esto debido a la codificación de algoritmos opacos que apoyan las desigualdades sociales actuales; por lo cual, es importante reconocer que para cada área de desarrollo se presentan desafíos únicos derivados de datos inestables y limitantes históricas.

Incluso se observa que los modelos modernos de estadística y aprendizaje automático, pueden experimentar sesgos inconscientes que afectan el desempeño de sistemas e impactan en el aumento de los resultados discriminatorios.

La aplicación de estrategias que busquen la justicia y transparencia, ayudan a minimizar el impacto del uso de datos desbalanceados y favorecen la programación exitosa de los requerimientos.

Se recomienda incorporar enfoques éticos y multidisciplinarios desde las fases iniciales del diseño de los sistemas basados en IA, así como establecer mecanismos de auditoría algorítmica que permitan evaluar su funcionamiento de forma transparente y responsable.

Asimismo, el uso de conjuntos de datos inclusivos y representativos resulta esencial para mitigar sesgos estructurales.

De cara a futuras investigaciones, sería especialmente valioso examinar casos en contextos históricamente subrepresentados en particular, África, América Latina, El Caribe, Asia y Oceanía, con el propósito de avanzar en el desarrollo de métricas de equidad que se ajusten a las diversas realidades socioculturales en las que se implementan estas tecnologías.

Entre las principales limitaciones de este estudio, se encuentra la posible parcialidad en la selección de fuentes bibliográficas, dado que se privilegió la revisión de publicaciones indexadas en bases de datos académicas de alta reputación.

Esta elección metodológica pudo haber excluido literatura relevante disponible en repositorios alternativos o en medios no tradicionales. Asimismo, el análisis se centró predominantemente en trabajos redactados en inglés y español, lo que introduce un sesgo lingüístico que podría haber restringido la diversidad de enfoques culturales y geográficos contemplados.

Estas consideraciones deben ser tenidas en cuenta al interpretar los resultados, y subrayan la necesidad de que futuras revisiones sistemáticas amplíen sus criterios de búsqueda y selección con el fin de incorporar una mayor heterogeneidad de voces, contextos y perspectivas.

Contribución de los autores

Freddy Aníbal Jumbo Castillo: Conceptualización, Metodología, Supervisión, Redacción – borrador original, Redacción – revisión y edición. **Johnny Paul Novillo Vicuña:** Investigación, Curación de datos, Análisis formal, Redacción – borrador original. **Camilly Yuliana Pacheco Ordoñez:** Validación, Análisis formal, Redacción – revisión y edición. **Joselyn Katuska Franco Ávila:** Validación, Interpretación de datos, Administración del proyecto, Redacción – revisión y edición.

Conflictos de interés

Los autores declaran no tener ningún conflicto de interés.

Referencias bibliográficas

- Ávila Bravo-Villasante, M. (2023). La agenda feminista ante la cuarta revolución industrial. Mujeres y algoritmización de la esfera pública. *Cuestiones de Género: De La Igualdad y La Diferencia*, 18, 132–155. <https://doi.org/10.18002/cg.i18.7573>
- Akter, S., McCarthy, G., Sajib, S., Michael, K., Dwivedi, Y. K., D’Ambra, J., & Shen, K. N. (2021). Algorithmic bias in data-driven innovation in the age of AI. *International Journal of Information Management*, 60. <https://doi.org/10.1016/j.ijinfomgt.2021.102387>
- Bagga, S., & Piper, A. (2020). Measuring the effects of bias in training data for literary classification. *LaTeCH-CLFL*.
- Beltramelli Gula, N., Ferro, C., Goñi Mazzitelli, M., Etcheverry, L., & Rocamora, M. (2023). Un concepto viajero. *Novos Rumos Sociológicos*, 10(18). <https://doi.org/10.15210/norus.v10i18.4847>
- Bravo Bolado, A. (2023). Justicia algorítmica: Un enfoque sociotécnico. *Estudios Penales y Criminológicos*, 1–42. <https://doi.org/10.15304/epc.44.8838>
- Cary, M. P., Zink, A., Wei, S., Olson, A., Yan, M., Senior, R., Bessias, S., Gadhomi, K., Jean-Pierre, G., Wang, D., Ledbetter, L. S., Economou-Zavlanos, N. J., Obermeyer, Z., & Pencina, M. J. (2023). Mitigating racial and ethnic bias and advancing health equity in clinical algorithms: A scoping review. *Health Affairs*, 42(10). <https://doi.org/10.1377/HLTHAFF.2023.00553>
- Chang, X. (2023). Gender bias in hiring: An analysis of the impact of Amazon’s recruiting algorithm. *Advances in Economics, Management and Political Sciences*, 23(1). <https://doi.org/10.54254/2754-1169/23/20230367>
- Chen, Z., Zhou, Y., Wang, Z., Liu, F., Leng, T., & Zhu, H. (2025). A bias evaluation solution for multiple sensitive attribute speech recognition. *Computer Speech & Language*, 93, 101787. <https://doi.org/10.1016/J.CSL.2025.101787>
- Coraglia, G., Genco, F. A., Piantadosi, P., Bagli, E., Giuffrida, P., Posillipo, D., & Primiero, G. (2024). Evaluating AI fairness in credit scoring with the BRIO tool. *arXiv*. <https://arxiv.org/abs/2406.03292>
- de Lima, R. M., Pisker, B., & Corrêa, V. S. (2023). Gender bias in artificial intelligence: A systematic review of the literature. *Journal of Telecommunications and the Digital Economy*, 11(2). <https://doi.org/10.18080/jtde.v11n2.690>
- DeCamp, M., & Lindvall, C. (2023). Mitigating bias in AI at the point of care. *Science*, 381(6654). <https://doi.org/10.1126/science.adh2713>
- Ferrara, E. (2024). Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *Sci*, 6(1). <https://doi.org/10.3390/sci6010003>
- Fuster, A., Goldsmith-Pinkham, P. S., Ramadorai, T., & Walther, A. (2021). Predictably unequal? The effects of machine learning on credit markets. *Journal of Finance*. <https://doi.org/10.2139/ssrn.3072038>
- Ghosal, I., & Hooker, G. (2020). Boosting random forests to reduce bias; One-step boosted forest and its variance estimate. *Journal of Computational and Graphical Statistics*, 30(2). <https://doi.org/10.1080/10618600.2020.1820345>
- Gilbert, S., Mehl, A., Baluch, A., Cawley, C., Challiner, J., Fraser, H., Millen, E., Montazeri, M., Multmeier, J., Pick, F., Richter, C., Türk, E., Upadhyay, S., Virani, V., Vona, N., Wicks, P., & Novorol, C. (2020). How accurate are digital symptom assessment apps for suggesting conditions and urgency advice? A clinical vignettes comparison to GPs. *BMJ Open*, 10(12), e040269. <https://doi.org/10.1136/bmjopen-2020-040269>
- Golden, S. H., Joseph, J. J., & Hill-Briggs, F. (2021). Casting a health equity lens on endocrinology and diabetes. *Journal of Clinical Endocrinology and Metabolism*, 106(4). <https://doi.org/10.1210/clinem/dgaa938>
- Hamilton, Z., Duwe, G., Kigerl, A., Gwinn, J., Langan, N., & Dollar, C. (2022). Tailoring to a mandate: The development and validation of the Prisoner Assessment Tool Targeting Estimated Risk and Needs (PATTERN). *Justice Quarterly*, 39(6). <https://doi.org/10.1080/07418825.2021.1906930>

- Jie, Z., Zhiying, Z., & Li, L. (2021). A meta-analysis of Watson for Oncology in clinical application. *Scientific Reports*, *11*(1), 5792. <https://doi.org/10.1038/s41598-021-84973-5>
- Kudina, O., & de Boer, B. (2021). Co-designing diagnosis: Towards a responsible integration of Machine Learning decision-support systems in medical diagnostics. *Journal of Evaluation in Clinical Practice*, *27*(3). <https://doi.org/10.1111/jep.13535>
- Li, R., Kingsley, S., Fan, C., Sinha, P., Wai, N., Lee, J., Shen, H., Eslami, M., & Hong, J. (2023). Participation and division of labor in user-driven algorithm audits: How do everyday users work together to surface algorithmic harms? Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems. <https://doi.org/10.1145/3544548.3582074>
- Martínez, N., Vinas, A., & Matute, H. (2021). Examining potential gender bias in automated-job alerts in the Spanish market. *PLoS ONE*, *16*(12), e0260409. <https://doi.org/10.1371/journal.pone.0260409>
- Noiret, S., Lumetzberger, J., & Kappel, M. (2021). Bias and fairness in computer vision applications of the criminal justice system. 2021 IEEE Symposium Series on Computational Intelligence (SSCI), 1–8. <https://doi.org/10.1109/SSCI50451.2021.9660177>
- Paredes Meneses, J. (2023). Aplicación informática para reconocimiento de la especie Camu Camu (*Myrciaria Dubia*) a través de redes neuronales convolucionales, en Iquitos Perú, durante el año 2017. Aplicación informática para reconocimiento de la especie Camu Camu (*Myrciaria Dubia*) a través de redes neuronales convolucionales, en Iquitos Perú, durante el año 2017. https://doi.org/10.37811/cli_w945
- Peng, Y. (2023). The role of ideological dimensions in shaping acceptance of facial recognition technology and reactions to algorithm bias. *Public Understanding of Science*, *32*(2). <https://doi.org/10.1177/09636625221113131>
- Pérez López, J. I. (2023). Inteligencia artificial y contratación laboral. *Revista De Estudios Jurídico Laborales Y De Seguridad Social (REJLSS)*, *7*, 186–205. <https://doi.org/10.24310/rejls7202317557>
- Reyes Campos, J. E. M., Castañeda Rodríguez, C. S., Alva Luján, L. D., & Mendoza de los Santos, A. C. (2023). Sistema de reconocimiento facial para el control de accesos mediante inteligencia artificial. *Innovación y Software*, *4*(1). <https://doi.org/10.48168/innosoft.s11.a78>
- Sanabria Moyano, J. E., Roa Avella, M. del P., & Lee Pérez, O. I. (2022). Tecnología de reconocimiento facial y sus riesgos en los derechos humanos. *Revista Criminalidad*, *64*(3), 61–78. <https://doi.org/10.47741/17943108.366>
- Santiago Arenas, A., Samboni, O., Villegas Trujillo, L. M., Zamora Córdoba, I., & Alfonso Morales, G. (2023). Tratamiento de la hipersensibilidad dentinaria primaria: Revisión exploratoria. *Salud Uninorte*, *39*(3), 1120–1152. <https://doi.org/10.14482/sun.39.03.741.258>
- Seyyed-Kalantari, L., Zhang, H., McDermott, M. B. A., Chen, I. Y., & Ghassemi, M. (2021). Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations. *Nature Medicine*, *27*(12). <https://doi.org/10.1038/s41591-021-01595-0>
- Simonetta, A., Trenta, A., Paoletti, M. C., & Vetrò, A. (2021). Metrics for identifying bias in datasets. *CEUR Workshop Proceedings*, *3118*.
- Tang, L., Li, J., & Fantus, S. (2023). Medical artificial intelligence ethics: A systematic review of empirical studies. *Digital Health*, *9*. <https://doi.org/10.1177/20552076231186064>
- Tsamados, A., Aggarwal, N., Cowls, J., Morley, J., Roberts, H., Taddeo, M., & Floridi, L. (2022). The ethics of algorithms: Key problems and solutions. *AI and Society*, *37*(1). <https://doi.org/10.1007/s00146-021-01154-8>
- Varona, D., & Suárez, J. L. (2022). Discrimination, bias, fairness, and trustworthy AI. *Applied Sciences*, *12*(12). <https://doi.org/10.3390/app12125826>
- Vela, M. B., Erondú, A. I., Smith, N. A., Peek, M. E., Woodruff, J. N., & Chin, M. H. (2022). Eliminating explicit and implicit biases in health care: Evidence and research needs. *Annual Review of Public Health*, *43*. <https://doi.org/10.1146/annurev-publhealth-052620-103528>
- Wang, S., & Wang, H. (2020). Big data for small and medium-sized enterprises (SME): A knowledge management model. *Journal of Knowledge Management*, *24*(4), 881–897. <https://doi.org/10.1108/JKM-02-2020-0081>
- Yu, Y., & Saint-Jacques, G. (2022). Choosing an algorithmic fairness metric for an online marketplace: Detecting and quantifying algorithmic bias on LinkedIn. arXiv. <https://arxiv.org/abs/2202.07300>
- Zhang, J. (2016). Testing the predictive validity of the LSI-R using a sample of young male offenders on probation in Guangzhou, China. *International Journal of Offender Therapy and Comparative Criminology*, *60*(4). <https://doi.org/10.1177/0306624X14557471>

