

Limitaciones del uso de la inteligencia artificial: revisión crítica de su aplicación en el campo de la demografía, a partir del caso de la tuberculosis en la provincia de Córdoba (1895-1914).

Virginia Alba Adauto y Victor Maximiliano Adauto.

Cita:

Virginia Alba Adauto y Victor Maximiliano Adauto (2025). *Limitaciones del uso de la inteligencia artificial: revisión crítica de su aplicación en el campo de la demografía, a partir del caso de la tuberculosis en la provincia de Córdoba (1895-1914)*. XVIII Jornadas Argentinas de Estudios de Población - V Congreso Internacional de Población del Cono Sur. Asociación de Estudios de Población de la Argentina, Córdoba.

Dirección estable: <https://www.aacademica.org/xviii.jornadas.aepa/55>

ARK: <https://n2t.net/ark:/13683/exQq/ENn>



Esta obra está bajo una licencia de Creative Commons.
Para ver una copia de esta licencia, visite
<https://creativecommons.org/licenses/by-nc/4.0/deed.es>.

Acta Académica es un proyecto académico sin fines de lucro enmarcado en la iniciativa de acceso abierto. Acta Académica fue creado para facilitar a investigadores de todo el mundo el compartir su producción académica. Para crear un perfil gratuitamente o acceder a otros trabajos visite:
<https://www.aacademica.org>.



Limitaciones del uso de la inteligencia artificial: revisión crítica de su aplicación en el campo de la demografía, a partir del caso de la tuberculosis en la provincia de Córdoba (1895–1914)

Adauto, Virginia Alba

Facultad de Filosofía y Humanidades

Universidad Nacional de Córdoba

Virginia.adauto@mi.unc.edu.ar

Adauto, Victor Maximiliano

Facultad de Ciencias Químicas

Universidad Nacional de Córdoba

Victor.adauto@mi.unc.edu.ar

Resumen

Este estudio propone un análisis crítico de la aplicación de técnicas de Inteligencia Artificial (IA) en el análisis de fenómenos demográficos históricos.

Tomando como caso de estudio la provincia de Córdoba (Argentina) entre 1895 y 1914, momento en el que hay una epidemia de tuberculosis, que pese a la alta letalidad (especialmente entre mujeres jóvenes de sectores populares), no se refleja la presencia de la enfermedad en los censos nacionales de la época; los registros muestran un crecimiento poblacional.

A través de este estudio de caso, se buscará problematizar la presunta neutralidad de los sistemas de IA cuando se aplican a conjuntos de datos históricos que presentan: omisiones sistemáticas, limitaciones en su recolección original, y sesgos ideológicos intrínsecos a su producción. Cuestionamos hasta qué punto las herramientas computacionales pueden compensar las distorsiones presentes en los registros demográficos históricos, y qué nuevos sesgos pueden las IA's introducir en el proceso de análisis de datos.

Introducción

El presente trabajo consiste en la puesta en escena de un ensayo del proceso de análisis de datos antiguos mediante Inteligencia Artificial. Teniendo en cuenta la alta accesibilidad y concurrencia que tienen las inteligencias artificiales hoy en día, quisimos examinar la eficacia y posibles usos de la IA como herramienta prometedora para posiblemente ser usada en el campo de la investigación.

Para el abordaje del estudio de los posibles usos y limitaciones, es que se optó por trabajar con la puesta en escena de un análisis de las fuentes en relación a Córdoba entre 1895 y 1914, período de una epidemia de tuberculosis, fenómeno disruptivo ampliamente documentado en la literatura histórica y otros trabajos académicos, lo que permite contrastar los resultados algorítmicos con evidencia ya establecida.

El estudio se centra específicamente en el proceso de transformación de documentos escaneados como fichas censales en bases de datos estructuradas que permitan su explotación sistemática en investigación. Esta problemática adquiere relevancia al observar que a pesar de tratarse de registros de origen público, la ausencia de versiones digitalizadas accesibles y procesables computacionalmente limita severamente su análisis cuantitativo y comparativo. Para abordar este vacío, se implementa una metodología comparativa aplicada a distintos modelos de IA de código abierto, con el fin de extraer, clasificar e interpretar datos demográficos de fuentes primarias previamente digitalizadas en formato imagen, carentes de reconocimiento óptico de caracteres (OCR) o transcripción.

Se analizan críticamente cuestiones como la capacidad de estos sistemas para identificar sesgos inherentes a las fuentes (omisiones sistemáticas, construcciones ideológicas) o, por el contrario, su potencial para reproducirlos y amplificarlos. Asimismo, se reflexiona sobre la difícil demarcación entre el análisis automatizado y la interpretación humana, así como sobre la presunta neutralidad de los algoritmos al operar con datos social e históricamente situados. Esto, teniendo en cuenta que la inteligencia artificial constituye

un caso paradigmático: diversos modelos se encuentran en constante rediseño y evolución, lo que abre un abanico amplio de posibilidades en lo que respecta al análisis de datos.

Tal como plantea la literatura reciente, la irrupción de la inteligencia artificial generativa en la comunidad científica ha sido descrita como un fenómeno con un carácter ambivalente: por un lado, representa una oportunidad inédita para agilizar procesos de investigación, gestionar grandes volúmenes de información y ampliar la capacidad analítica de los investigadores; por otro, genera un conjunto de desafíos relacionados con la confianza profesional, la dependencia tecnológica y la ausencia de consensos éticos sobre su implementación (Wiley, 2024).

Consideramos que al no ser sencillo manejar las ciencias computacionales a niveles complejos, es relevante la constitución de bases de datos estructuradas que faciliten la investigación demográfica y en las ciencias humanas.

Marco conceptual y antecedentes

En relación al uso, manejo y procesamiento de datos demográficos provenientes de registros estatales y extraestatales de carácter histórico, diversos trabajos previos han señalado una serie de dificultades inherentes a este tipo de fuentes.

Al indagar, es posible observar que estos tipos de fuentes presentan desafíos metodológicos significativos. Estos problemas se originan, en primer lugar, en el momento mismo de su registro, manifestándose mediante problemáticas como la ausencia sistemática de datos, errores de digitación y la consignación imprecisa o mal interpretada de la información. Además de estas limitaciones inherentes a la fuente, surgen problemas de distinta naturaleza durante el análisis longitudinal; destacan la variabilidad onomástica (cambios en nombres y apellidos), que obstaculiza el seguimiento del ciclo vital de un individuo a través de diferentes registros, y la movilidad geográfica de la población

(eventos migratorios e inmigratorios), la cual complica la contextualización y vinculación coherente de los datos en el tiempo y el espacio.

Ante esto, varios investigadores exploraron formas de esquivar esos problemas con la utilización de fórmulas estadísticas, etc para tratar de mejorar la información y la precisión.

Entre estos, un precedente relevante es el estudio realizado por Marks-Anglin et al. (2024) en este trabajo se abordó el desafío de realizar análisis de supervivencia a partir de un cruce de datos entre el censo nacional de Estados Unidos de 1930 y registros vitales de defunción extraestatales. La principal limitación metodológica de este estudio, inherente al uso de fuentes históricas, es la ausencia de identificadores únicos, lo que impide una vinculación certera de los registros. Esta imperfecta vinculación da lugar a un problema de datos faltantes no aleatorios y a potenciales errores de clasificación por la inclusión de falsos emparejamientos (falsos positivos) o la exclusión de emparejamientos válidos (falsos negativos).

La pertinencia de citar esta investigación radica en su enfoque explícito en el desarrollo y comparación de métodos estadísticos diseñados para corregir los sesgos introducidos por una vinculación imperfecta. Los autores evaluaron diversas técnicas de imputación y ajuste, concluyendo que los métodos más óptimos fueron aquellos que aprovecharon la información covariante de individuos con perfiles sociodemográficos similares para imputar el estado de supervivencia.

Este enfoque permitió generar estimaciones de hazard ratios más precisas y menos sesgadas, demostrando su utilidad en una aplicación de epidemiología histórica para evaluar el impacto de la exposición laboral al asbesto. Si bien este estudio provee un marco estadístico sólido para el manejo de datos vinculados de forma imperfecta, no explora la potencial aplicación de técnicas de inteligencia artificial o aprendizaje automático para optimizar la fase inicial de vinculación de registros (record linkage) o para mejorar los procesos de imputación. Por lo tanto, su trabajo establece una base

metodológica crucial sobre la cual es posible construir, integrando avanzadas técnicas computacionales para abordar problemas análogos.

Otro antecedente fundamental lo constituye el estudio de Helgertz et al. (2021), cuyo propósito fue crear y validar una estrategia innovadora para conectar registros individuales entre los censos de 1900 y 1910, con aplicabilidad potencial a un rango más amplio de periodos (1850-1940). Esta investigación se inserta en el marco del proyecto IPUMS Multigenerational Longitudinal Panel, cuyo objetivo es construir un panel longitudinal multigeneracional que permita seguir trayectorias de vida y dinámicas familiares desde mediados del siglo XIX hasta mediados del XX.

El desafío central identificado por los autores es la ausencia de identificadores únicos en los censos históricos, lo que dificulta rastrear a las mismas personas a lo largo del tiempo y genera sesgos de representatividad en los vínculos obtenidos. Los métodos tradicionales de record linkage, como el enfoque probabilístico clásico de Fellegi–Sunter, suelen presentar limitaciones como, por ejemplo, la tendencia a favorecer perfiles demográficos específicos (varones, población blanca e individuos con nombres poco comunes) y a dejar subrepresentados a migrantes, mujeres y minorías.

La innovación metodológica del trabajo consiste en un enfoque híbrido que combina distintas herramientas técnicas: Uso de variables relativamente estables en el tiempo (edad aproximada, lugar de nacimiento, ocupación, estructura del hogar), que permiten asignar probabilidades de coincidencia más consistentes; algoritmos de similitud de texto (e.g., métricas de Levenshtein), diseñados para manejar variaciones ortográficas y errores de enumeración en nombres y apellidos; modelo multigeneracional, que incorpora las relaciones familiares (padres, cónyuges, hijos) como evidencia contextual para reforzar la validación de emparejamientos; enfoque probabilístico/bayesiano, que asigna puntajes de enlace y descarta coincidencias dudosas, equilibrando la precisión con la cobertura; entre otros.

Los resultados muestran que esta estrategia incrementa notablemente la tasa de vinculaciones válidas y reduce sesgos de selección respecto de métodos previos, lo que permite la construcción de un recurso más sólido para el análisis de movilidad intergeneracional, dinámicas familiares y epidemiología histórica. El panel resultante ofrece la posibilidad de rastrear familias a lo largo de varias generaciones, posibilitando así, un análisis tanto de movilidad intergeneracional como dinámica familiar, ampliando de manera significativa el potencial de la investigación en historia social y demografía cuantitativa.

A pesar de que los autores (Helgertz, et al, 2021) reconocen limitaciones inherentes, como los trade-offs entre precisión y exhaustividad, la dependencia de la calidad de los datos originales y la posibilidad residual de errores de clasificación (falsos positivos y negativos); consideramos valiosa la experiencia del manejo de algoritmos computacionales avanzados y de un marco probabilístico sofisticado. Sin embargo, a pesar de la combinación técnicas de procesamiento de datos, un algoritmo de aprendizaje automático supervisado, para la calificación probabilística, y un modelo contextual multigeneracional; en su marco no hayamos incorporación de técnicas más modernas de IA como el aprendizaje automático, lo que delimita su aporte y señala un área de potencial desarrollo futuro.

En estos trabajos se observó las maneras de optimizar enfoques estadísticos y probabilísticos tradicionales, sin embargo no se aplicaron técnicas modernas de IA o machine learning, es decir, el sistema no “aprende” de ejemplos previos para mejorar progresivamente su desempeño, sino que aplica reglas estadísticas predefinidas con mayor refinamiento.

El uso del matching learning y el deep learning de la IA, es algo que consideramos relevante tener en cuenta porque la investigación en ciencias humanas puede beneficiarse de las múltiples posibilidades de utilización que ofrece la IA. Tal como dice Diaz Sanchez (2024), para gente no versada en computación, la existencia de una herramienta como la IA que puede ser usada como un recurso para acelerar y hacer más eficiente la

investigación, teniendo la capacidad de digitalizar, examinar e interpretar grandes volúmenes de datos con precisión y velocidad, es una cualidad a tener en cuenta.

El proceso de aprendizaje de los modelos de inteligencia artificial, particularmente aquellos basados en redes neuronales profundas, requiere disponer de grandes volúmenes de datos de calidad. La efectividad del aprendizaje profundo se fundamenta en la capacidad del sistema de procesar información en distintos niveles de manera simultánea, lo que hace imprescindible no solo la cantidad, sino también la confiabilidad y coherencia de los datos utilizados. De acuerdo con lo señalado por Whang (2023), si bien existen algoritmos y técnicas que permiten mitigar parcialmente el impacto de datos imperfectos, los problemas derivados de la insuficiencia, el ruido, el sesgo o la contaminación de los datos representan un desafío metodológico sustancial.

En este sentido, aunque se dispone de estrategias de validación, limpieza e integración, no todas resultan plenamente compatibles con los modelos de aprendizaje profundo. Aun en escenarios en los que los datos no puedan ser depurados completamente, es posible gestionarlos mediante técnicas de entrenamiento robusto que permiten que el modelo mantenga un desempeño aceptable frente a imperfecciones.

Habiendo señalado la relevancia que adquieren las bases de datos, resulta pertinente hacer referencia a los censos poblacionales como fuente primaria de información.

En el caso de los censos nacionales de 1895 y 1914, los registros disponibles señalan la presencia en el país de enfermedades como fiebre tifoidea, tos seca y bocio, entre otras. No obstante, en la provincia de Córdoba únicamente se consigna la “tos seca” en el censo de 1914, sin que en ningún momento se mencione la tuberculosis. Este silencio resulta significativo, ya que en otros documentos provinciales de fines del siglo XIX y comienzos del siglo XX sí aparecen referencias a dicha enfermedad, lo que evidencia las limitaciones de los censos como fuentes exclusivas para el estudio demográfico.

En este sentido, la organización de los censos y las categorías empleadas para clasificar la información no eran neutrales, sino que respondían a decisiones políticas, ideológicas y técnicas propias de cada época. Tal como sostiene Sacco (s.f.), los censos constituyen una tecnología estatal que produce representaciones sociales atravesadas por dichas elecciones y limitaciones. En consecuencia, las modalidades de relevamiento incidían de manera directa en la representación de los distintos sectores de la población, dejando en muchos casos a grupos enteros invisibilizados. Particularmente, las poblaciones en condiciones de mayor vulnerabilidad, como mujeres jóvenes de bajos recursos, trabajadores con acceso limitado a los servicios de salud, habitantes de conventillos y barrios precarios solían ser registradas de manera incompleta o incluso quedar excluidas.

Ello se debía tanto a dificultades metodológicas para acceder a esos espacios como a prácticas propias del empadronamiento, entre ellas la de anotar únicamente a quienes se encontraban presentes en el hogar al momento de la visita censal o la exclusión sistemática de personas institucionalizadas.

En este marco, la tuberculosis, una enfermedad que afectaba desproporcionadamente a los sectores más desfavorecidos, aparece como un ejemplo claro de los sesgos presentes en los censos históricos. La subrepresentación de los grupos más pobres y hacinados contribuyó a invisibilizar la verdadera magnitud de la enfermedad en Córdoba y en otras regiones del país, limitando la posibilidad de reconstruir de manera precisa su impacto social y sanitario.

Más allá de los censos, incluso las fuentes médicas y estadísticas oficiales presentan limitaciones importantes. Entre 1906 y 1918, tanto los certificados de defunción como los registros civiles muestran un patrón de omisiones y errores sistemáticos en la consignación de la causa de muerte por tuberculosis.

Finalmente, como postulan Carbonetti y Beltramone (2021), no puede desestimarse que, en algunos casos, la omisión de la tuberculosis en los certificados de defunción respondiera a una intencionalidad más profunda. Debido a la fuerte estigmatización social

que implicaba, algunos médicos habrían optado deliberadamente por disfrazar la causa de muerte con diagnósticos alternativos menos comprometedores para los familiares. Así, los silencios y las imprecisiones en torno a la tuberculosis no fueron solo el resultado de limitaciones técnicas, sino también de prácticas sociales y culturales que buscaban amortiguar el peso simbólico de una enfermedad asociada con la pobreza, la marginalidad y la exclusión.

Objetivos

El objetivo general del texto es evaluar críticamente las posibilidades y limitaciones del uso de la inteligencia artificial en el análisis de datos históricos, tomando como caso de estudio los censos argentinos de 1895 y 1914 en la provincia de Córdoba.

Los objetivos específicos son:

1. Analizar la capacidad de las inteligencias artificiales para digitalizar y estructurar documentos históricos censales: Se pretende evaluar en qué medida las herramientas de reconocimiento óptico de caracteres (OCR), junto con técnicas más avanzadas de procesamiento de lenguaje natural y visión computacional, permiten transformar fichas censales manuscritas o mecanografiadas en bases de datos estructuradas y confiables para la investigación histórica y demográfica. Este objetivo no se limita únicamente a la fase de transcripción, sino que también contempla la normalización de los datos, la resolución de problemas derivados de la legibilidad (tachaduras, caligrafía deficiente, deterioro físico de los documentos), y la creación de un formato accesible que pueda ser reutilizado en distintos entornos de investigación. La relevancia de este objetivo radica en que, pese a la importancia de los censos de 1895 y 1914 como fuentes primarias, no existen versiones digitalizadas de acceso libre y sistemático, lo cual dificulta el análisis cuantitativo y comparativo.
2. Examinar la potencialidad de los algoritmos para identificar patrones demográficos ocultos y variables disruptivas no registradas explícitamente en las fuentes históricas: Se apuntó a explorar cómo distintos modelos de inteligencia

artificial, en particular los de aprendizaje automático y aprendizaje profundo, pueden aportar nuevas perspectivas al análisis de los censos, revelando dinámicas poblacionales que no emergen de una lectura lineal o descriptiva de los registros. Buscamos determinar, si los algoritmos permiten detectar correlaciones entre variables aparentemente inconexas, rastrear la incidencia de fenómenos coyunturales como epidemias o sequías, o estimar la magnitud de procesos sociales más amplios como la migración interna e internacional. Tuvimos un abordaje comparativo en dos direcciones: por un lado, contrastar diferentes modelos de IA entre sí (supervisados, no supervisados, generativos, etc.), y por otro, cotejar los resultados obtenidos con los métodos tradicionales de análisis histórico-demográfico. De esta forma, no solo se pondrá a prueba la capacidad de la IA para generar conocimiento nuevo, sino también su valor como complemento (no sustituto) de las metodologías humanísticas y sociales.

3. Problematicar los límites epistemológicos del uso de la inteligencia artificial en el análisis de fuentes históricas: Más allá de la eficacia técnica de las herramientas, se considera imprescindible reflexionar sobre las implicaciones epistemológicas de su aplicación en ciencias sociales e históricas. Cuestionamos la frontera entre el procesamiento automatizado y la interpretación humana, entendiendo que los datos censales no son neutrales ni completos, sino que están atravesados por sesgos, omisiones y construcciones ideológicas propias de su contexto de producción. En este marco, se analizará la capacidad o incapacidad de los algoritmos para reconocer estos sesgos y, en caso de no hacerlo, el riesgo de que los reproduzcan o incluso los amplifiquen. También abordamos la cuestión de la supuesta “neutralidad” de la inteligencia artificial, cuestionando la idea de que los modelos computacionales ofrecen lecturas objetivas del pasado. Postulamos la necesidad de una relación dialógica entre máquina e historiador, en la cual el procesamiento automatizado aporte nuevas posibilidades sin desplazar la mirada crítica y contextualizada propia de la investigación histórica.

Metodología y fuentes

Este trabajo se fundamenta en el análisis exhaustivo de las fichas censales originales del Segundo (1895) y Tercer Censo Nacional (1914) de Argentina, disponibles únicamente en formato de imágenes digitales a través de repositorios en línea. El problema observado es que no se encontró alguna base pública que contenga en su interior los datos provenientes de estos documentos.

Es esta observación la que nos llevó a reflexionar acerca de la posibilidad de experimentar realizar un proceso de digitalización de datos antiguos mediante herramientas de inteligencia artificial de código abierto, y observar las capacidades y limitaciones que podría llevar la elaboración de este trabajo con el uso de Inteligencia Artificial

Para poder realizar lo propuesto, se procedió a dividir el trabajo en dos fases claves: En la primera, procesamos los datos brutos extraídos por la IA sin proporcionarle ningún contexto histórico, geográfico o temporal, lo que esperamos que haya permitido evaluar su capacidad para identificar patrones demográficos intrínsecos y anomalías estadísticas sin influencia de sesgos interpretativos previos.

Para evaluar cuál es la Inteligencia Artificial de código abierto y que sea poseedora de un rendimiento superior en lo respectivo a la precisión en la lectura de documentos históricos con variaciones caligráficas y deterioro propio de su antigüedad; se realizó una evaluación con fines comparativos en el cual el elemento principal a observar fue el rendimiento en el reconocimiento óptico de caracteres (OCR).

Para la examinación, el procedimiento elegido consistió en la escritura de un mismo prompt a diferentes inteligencias artificiales, siendo estas: *chat gpt-5*, *geminis ia* y *deepseek* con el objetivo de que realicen un análisis de datos y realicen un cuadro comparando los resultados (habitantes por departamento) incluyendo el cambio absoluto y porcentual de habitantes. Cabe mencionar que junto al prompt se les subió a las IAs a

evaluar, dos documentos pertenecientes al censo de 1895 y 1914 respectivamente, pero previamente recortados para que aparezcan únicamente las páginas que presentaban datos de la provincia de Córdoba.

A partir de los resultados alcanzados en la primera instancia, se resolvió avanzar hacia una segunda fase en la cual se incorporaron metadatos específicos relativos a la procedencia cordobesa de los registros, a su carácter censal y a la delimitación temporal precisa que abarcan (1895-1914). Con este procedimiento se quiso analizar mediante contrastación, de qué manera la inclusión de información contextual en el prompt puede incidir en la respuesta producida, generando además inferencias complementarias elaboradas por el propio sistema. Conviene destacar, asimismo, que en este último prompt se solicitó no sólo la interpretación de los datos sino también la explicitación de su significado. De igual modo, se indicó que, al tratarse de registros censales correspondientes a dichos años, se valorará la eventual existencia de acontecimientos históricos que pudieran haber influido en las cifras allí consignadas.

Resultados

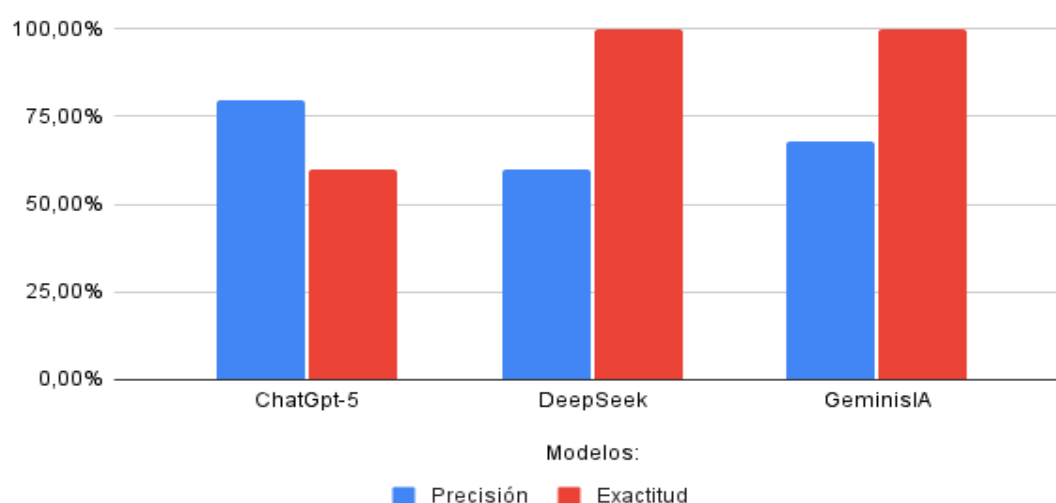
En la primera fase de la digitalización de datos, los resultados al comparar IAs hicieron que seleccionaremos específicamente al sistema *Gemini 2.5* como el más apto en la lectura de los documentos que se le fueron presentados, esto debido a que fue la Inteligencia Artificial que dio mejores resultados en el reconocimiento óptico de caracteres.

En el cotejo realizado, *Géminis IA* presentó un desempeño destacado, mientras que *DeepSeek* logró efectuar una lectura pero de manera deficiente con una cantidad considerable de errores. Un aspecto relevante identificado en este proceso es la necesidad de elaborar prompts altamente específicos y descriptivos, dado que el modelo tiende a desviarse del objetivo textual cuando las instrucciones no se formulan con suficiente precisión. Por su parte, *ChatGPT-4* no consiguió procesar los documentos; y, en la etapa de verificación, se constató que, si bien *Géminis* presentó algunos errores de lectura que

derivaron en variaciones respecto del resultado esperado, la magnitud de dichas imprecisiones fue menor en comparación con las observadas en los otros modelos.

La siguiente tabla gráfica muestra los resultados obtenidos de cada modelo de inteligencia artificial, en ella se encuentran cuantificadas las propiedades de “Precisión” y “Exactitud” a la lectura de datos. Cabe aclarar que un 100% de Exactitud radica en que presentaron un cuadro correctamente con todos los departamentos y describiendo un ascenso o descenso de población cuantitativo y porcentual de cada uno; y la “Precisión” se refiere a que tan correctos son los datos numéricos de la población

Precisión y Exactitud de los Datos



Cómo es posible observar, *DeepSeek* y *GeminisIA* presentan una exactitud del 100% debido a que cumplen correctamente todo lo especificado en el prompt pero cometen errores al momento de transcribir los datos poblacionales teniendo una precisión de 60% y 68% respectivamente.

Es pertinente señalar que, durante el desarrollo del presente trabajo, se hizo pública la actualización de la IA *ChatGPT* a la versión “.5”. Si bien esta versión continuó siendo incapaz de reconocer los documentos analizados, se decidió (dada su amplia difusión como una de las IAs más utilizadas) evaluar la herramienta denominada deep research, con el propósito de contrastar si la información recabada de manera autónoma por el

modelo podía aproximarse a los datos contenidos en el documento original. Los resultados mostraron que el sistema logró identificar 19 de los 25 departamentos de la provincia de Córdoba, con un nivel de exactitud relativamente alto, aunque no exento de errores, por eso presenta una exactitud de 60% pero una precisión de 80%, ya que al sacar los datos de internet coinciden mayormente con los del documento.

Cabe mencionar también que las tres inteligencias artificiales, en el momento de lectura de las fuentes tuvieron como mayor dificultad la identificación de los números “9,6,3 y 8”, el error consistió en la confusión entre 9 y 6 con 0 y también entre el 3 y 8. Este tipo de error se estima que es producto del tipo de fuente usada en la imprenta con la cual se escribieron los documentos.

Pasando a lo producido en la segunda fase, se ofreció en el prompt contexto temporal y espacial, y se pidió interpretación de los datos con una explicación y justificación de la respuesta dada. Los resultados fueron los que aparecen en la siguiente tabla:

Aspecto	Geminis IA	Deep Seek
Crecimiento poblacional	Subraya el crecimiento exponencial del 109% en 19 años. Lo atribuye al crecimiento vegetativo y sobre todo a la inmigración.	Señala un aumento del 66,8% entre 1869-1895 y casi duplicación entre 1895-1914. Resalta la combinación de inmigración y crecimiento natural.
Inmigración extranjera	Remarca que en 1914 los extranjeros eran el 20,5% de la población. Enfatiza el fenómeno como parte de la gran inmigración europea. Identifica italianos y españoles como principales grupos.	Destaca que los extranjeros eran el 10,1% en 1895 y 20,5% en 1914. Identifica nacionalidades (italianos, españoles, franceses, alemanes, austríacos). Analiza distribución regional (departamentos con alta inmigración agrícola).
Urbanización	Enfatiza el crecimiento acelerado de la ciudad de	Describe una urbanización incipiente: crecimiento de

	Córdoba, con un aumento del 146%, mayor que la tasa provincial. Interpreta a Córdoba capital como polo de atracción para inmigrantes.	Córdoba capital, Río Cuarto, Villa María, Bell Ville. Relaciona este aspecto con el ferrocarril, comercio y administración.
Transformación productiva y social	Destaca la integración de inmigrantes en actividades agrícolas en el interior y en servicios y oficios urbanos en la capital.	Señala colonización agrícola en el este y sur (Marcos Juárez, Unión, Juárez Celman, San Justo) con alta inmigración. Contrasta con departamentos del oeste y norte (Minas, Pocho, Sobremonte) más tradicionales.
Educación y alfabetización	No hace mención	Precisa que en 1895 solo 36,4% sabía leer y escribir, baja asistencia escolar (20-30%), con desigualdades regionales.
Fenómenos de larga duración	Focaliza en la gran inmigración europea como motor de cambio demográfico, social y cultural.	Identifica inmigración masiva, expansión agrícola, urbanización y nuevas élites, desigualdades regionales persistentes, y transición demográfica incipiente.
Enfoque interpretativo	Más sintético y focalizado: interpreta todo el proceso como expresión de la gran ola inmigratoria y su impacto en Córdoba.	Más analítico y multifactorial: combina variables demográficas, económicas, sociales, educativas y regionales.

Cómo es posible dar cuenta, un modelo plantea un enfoque mayor en la inmigración, identificándolo como un gran evento; mientras que el otro se puede decir que realizó un

análisis más cualitativo que no solo abarca la inmigración sino que también habla de la expansión agrícola y desigualdades regionales.

Sin embargo, ambos ignoraron el hecho de que en Córdoba a finales del siglo XIX y comienzos del siglo XX, existieron varias enfermedades infecciosas, algunas con un porcentaje de letalidad más alto que otras. Entre ellas se encuentran la tuberculosis, enfermedades respiratorias agudas como bronquitis, neumonías, viruela, sarampión, difteria y tos convulsa, gripe, cólera, sífilis, fiebre tifoidea, entre otras.

Como última cuestión a resaltar es notorio que, aunque le hayamos indicado un marco temporal preciso (1895/1914), los modelos de inteligencia artificial decidieron abarcar un espectro temporal más amplio para realizar su análisis demográfico.

Discusión

El presente estudio se propuso avanzar en la comprensión de problemáticas como: ¿qué ocurre cuando estos sistemas procesan datos que no reflejan de manera fidedigna la realidad histórica? ¿Hasta qué punto son capaces de incorporar variables contextuales externas o identificar ausencias significativas en los corpus? ¿Existe el riesgo de que lleguen a conclusiones erróneas al operar con información incompleta o intrínsecamente sesgada?

La elección de trabajar con modelos de código abierto responde a la necesidad de abordar estas cuestiones de forma accesible para la comunidad académica, reconociendo que no siempre se cuenta con la capacitación especializada que requieren herramientas de código cerrado. El objetivo fue evaluar la capacidad de interpretación y procesamiento de estos modelos en un escenario real de análisis documental. Los resultados demuestran que, efectivamente, existen limitaciones sustanciales en su capacidad de análisis contextual. Por ejemplo, los modelos omitieron por completo la dimensión epidemiológica del período estudiado, lo que plantea una interrogante más amplia: ¿qué otros factores contextuales críticos pueden estar siendo sistemáticamente ignorados?

Asimismo, se observó que, a pesar de tener acceso potencial a vastos repositorios de información en línea, los sistemas no incorporaron datos históricos-geográficos relevantes para el análisis, como el cambio en la división departamental de la provincia de Córdoba entre 1895 y la actualidad. Esta falta de contextualización afecta directamente la precisión y profundidad de sus outputs.

Cabe destacar que, al momento de la experimentación, ninguna de las tres plataformas evaluadas poseía la capacidad para una interpretación completamente precisa de texto generado por reconocimiento óptico de caracteres (OCR). No obstante, dado el ritmo acelerado de desarrollo en esta área (con actualizaciones y mejoras que se suceden de manera prácticamente semestral), los resultados aquí expuestos deben considerarse dentro de un marco temporal específico.

Surge, por tanto, la pregunta sobre la vigencia futura de estos hallazgos y en qué plazo una mejora en los algoritmos permitirá una lectura precisa que podría eventualmente superar las limitaciones aquí identificadas. Esta volatilidad tecnológica subraya la necesidad de abordar los estudios de IA con una perspectiva crítica y temporalmente situada, reconociendo que las conclusiones obtenidas están sujetas a una rápida evolución.

Conclusiones

El recorrido realizado permite hacer notorio que la aplicación de la inteligencia artificial en el análisis de fuentes demográficas históricas constituye una línea de investigación aún por explorar, aunque todavía está lleno de tensiones y desafíos.

Las experiencias de trabajo con los censos de Córdoba de 1895 y 1914 demostró que estas herramientas son útiles para transformar documentos en bases de datos manejables, pero al mismo tiempo dejan en evidencia limitaciones técnicas importantes, en particular en la lectura de caracteres y en la exactitud numérica. Esto lleva a reflexionar acerca de que la digitalización asistida por IA no como un proceso cerrado, sino como una instancia

que requiere de la presencia de gente profesional que esté capacitada para realizar tanto revisiones como correcciones sistemáticas.

Asimismo, la comparación entre distintos modelos permitió observar que, aunque se obtienen resultados de diversa calidad según la herramienta utilizada, persiste un obstáculo común: la dificultad para detectar las ausencias o silencios de las fuentes. La omisión de la tuberculosis como fenómeno epidemiológico central de la Córdoba de comienzos del siglo XX constituye un ejemplo paradigmático de cómo los algoritmos replican los vacíos y sesgos de origen, sin capacidad real de problematizarlos.

Otro aspecto significativo fue la influencia del contexto en los resultados. Cuando se aportaron metadatos históricos y temporales, los modelos lograron producir interpretaciones más detalladas, aunque todavía insuficientes para alcanzar un análisis crítico de calidad. Esto confirma que la IA no puede reemplazar la mirada académica: su potencial aparece cuando se integra de manera reflexiva y en diálogo con el conocimiento experto

La neutralidad atribuida a los algoritmos queda, de este modo, cuestionada ya que su desempeño depende tanto de la calidad del corpus como de la precisión de las instrucciones con que se los guíe.

Si bien el riesgo de reproducir los sesgos de las fuentes es innegable, el perfeccionamiento de estas tecnologías abre la posibilidad de acelerar la digitalización, identificar tendencias ocultas y explorar nuevos caminos de investigación.

Más que una herramienta neutral o un sustituto de la interpretación, la IA se perfila como un insumo complementario que, bajo un uso crítico y contextualizado, puede enriquecer la producción de conocimiento histórico y demográfico.

Referencias Bibliográficas

- Carbonetti, A., & Beltramone, J. (2021). Historiar una enfermedad. Fuentes y estrategias para abordar el estudio de la tuberculosis en la ciudad de Córdoba entre principios y mediados del siglo XX. *Universidade Federal do Piauí. Contraponto*, 1(1), 27–48.
- Chanipal-Heras, D., & Diaz-Sanchez, C. (2024). A review of AI applications in Human Sciences research. *Digital Applications in Archaeology and Cultural Heritage*, 32.
- Helgertz, J., Price, J., Wellington, J., Thompson, K. J., Ruggles, S., & Fitch, C. A. (2021). A new strategy for linking U.S. historical censuses: A case study for the IPUMS multigenerational longitudinal panel. *Historical Methods: A Journal of Quantitative and Interdisciplinary History*, 55(1), 12–29.
- Marks-Anglin, A. K., Barg, F. K., Ross, M., et al. (2024). Survival analysis under imperfect record linkage using historic census data. *BMC Medical Research Methodology*, 24, 67.
- Novick, S. (1995). Políticas de población en la Argentina (1870-1989). Una visión del Estado. *Estudios Demográficos y Urbanos*, 10(2).
- Otero, H. (2006). Estadística y nación. Una historia conceptual del pensamiento censal de la Argentina moderna, 1869-1914. Buenos Aires: Editorial Prometeo, Colección Historia Argentina.
- Pu, Z., Shi, C. L., Jeon, C. O., Fu, J., Liu, S. J., Lan, C., Yao, Y., Liu, Y. X., & Jia, B. (2024). ChatGPT and generative AI are revolutionizing the scientific community: A Janus-faced conundrum. *iMeta*, 3(2), Article e178. <https://doi.org/10.1002/imt2.178>
- República Argentina. Dirección de Estadísticas e Investigaciones Económicas de Mendoza. (s.f.). Segundo Censo Nacional de Población, 1895: Fichas censales.

<https://deie.mendoza.gov.ar/#!/censos-nacionales-de-poblacion/1895-segundo-censo-nacional-18>

República Argentina. Dirección de Estadísticas e Investigaciones Económicas de Mendoza. (s.f.). Tercer Censo Nacional de Población, 1914: Fichas censales. <https://deie.mendoza.gov.ar/#!/censos-nacionales-de-poblacion/1914-tercer-censo-nacional-38>

Sacco, N. (s.f.). Los censos de población modernos. En Libro de cocina para el análisis de las clases sociales en Argentina. Universidad Nacional de General Sarmiento.

Whang, S. E., Roh, Y., Song, H., & Lee, J.-G. (2023). Data collection and quality challenges in deep learning: A data-centric AI perspective. *The VLDB Journal*, 32(4), 791–813. <https://doi.org/10.1007/s00778-022-00775-9>