



Detección Temprana de Diabetes Mellitus Tipo 2 Mediante Machine Learning Aplicado a Datos Biométricos.

Alumno: Figueiras Manuel Angel

Universidad Abierta Interamericana

Profesor/es trabajo final: Sartorio Alejandro

Plan de Trabajo Final de Carrera presentado para obtener el título de Licenciado en

Gestión de Tecnología Informática

(Diciembre, 2025)

Resumen.

El presente trabajo analiza el potencial del Machine Learning para mejorar la detección temprana de la diabetes mellitus tipo 2 mediante el uso de datos biométricos. A partir de una revisión sistemática de la literatura y del diseño conceptual de un modelo predictivo basado en datos representativos de la población argentina, se comparan enfoques, variables relevantes y métricas de desempeño reportadas por la evidencia científica. Los resultados muestran que algoritmos como Random Forest, XGBoost y Redes Neuronales alcanzan altos niveles de sensibilidad y capacidad discriminativa, lo que evidencia su utilidad para la identificación temprana del riesgo de padecer diabetes mellitus tipo 2. El estudio destaca oportunidades y desafíos para su adopción en el sistema de salud argentino, especialmente en relación con la disponibilidad, estandarización y calidad de los datos clínicos.

Palabras clave: aprendizaje automático, datos biométricos, diabetes mellitus, detección temprana, modelos predictivos, salud digital.

Abstract.

The present study analyzes the potential of Machine Learning to improve the early detection of type 2 diabetes mellitus through the use of biometric data. Based on a systematic literature review and the conceptual design of a predictive model built from data representative of the Argentine population, the research compares methodological approaches, relevant variables, and performance metrics reported in scientific evidence. The findings show that algorithms such as Random Forest, XGBoost, and Neural Networks achieve high levels of sensitivity and discriminative capacity, demonstrating their usefulness for the early identification of type 2 diabetes mellitus risk. Additionally, the study highlights opportunities and challenges for the adoption of these techniques within the Argentine healthcare system, particularly regarding the availability, standardization, and quality of clinical data.

Keywords: biometric data, digital health, diabetes mellitus; early detection; machine learning; predictive models.

Índice General

Resumen.....	2
Abstract.....	3
Índice General.....	4
Índice de Ilustraciones	6
Índice de Tablas	7
1. Capítulo 1: Introducción.....	8
1.1. Planteamiento del Problema de Investigación.....	8
1.2. Justificación de la Elección del Problema.....	8
1.3. Objetivos del Trabajo Final.....	9
1.4. Preguntas de Investigación.....	10
1.5. Hipótesis del Problema.....	10
1.6. Metodología de investigación.....	10
1.7. Estrategia de Revisión Bibliográfica.....	11
1.8. Estructura General del Trabajo Final.....	13
1.9. Estado del Arte.....	14
2. Capítulo 2: Marco Teórico.....	17
2.1. Introducción.....	17
2.2. Fundamentos Clínicos: Diabetes y Variables de Riesgo.....	17
2.3. Fundamentos Tecnológicos: Aprendizaje Automático (ML).....	22

2.4.	Modelización Predictiva en el Contexto Clínico.	25
2.5.	Comparación de Modelos Predictivos: Clásicos vs Avanzados.	29
3.	Capítulo 3: Diseño Experimental y Metodología	33
3.1.	Introducción.	33
3.2.	Diseño Conceptual del Dataset.	33
3.3.	Consideraciones Éticas en el Uso de Datos Biométricos.....	37
3.4.	Preprocesamiento y Limpieza de Datos.....	39
3.5.	Flujo Metodológico General.	41
4.	Capítulo 4: Resultados y Evaluación	42
4.1.	Introducción.	42
4.2.	Resultados de la Evaluación Inferencial y Selección del Modelo Óptimo.	42
4.3.	Análisis e Interpretación del Modelo Optimo.....	44
4.4.	Discusión Crítica: Aplicabilidad y Desafíos en el Contexto Argentino.	46
5.	Capítulo 5: Conclusiones.	51
5.1.	Conclusiones Generales.	51
5.2.	Limitaciones del Estudio.....	52
5.3.	Líneas de Investigaciones Futuras.	52
5.4.	Conclusión Final.	54
6.	Acrónimos.....	55
7.	Referencias.....	57

Índice de Ilustraciones

Ilustración 1 <i>Diagrama PRISMA del Proceso de Selección de Estudios.</i>	12
Ilustración 2 <i>Proyección de Diabetes en la Población Argentina.</i>	18
Ilustración 3 <i>Prevalencia, Factores y Complicaciones.</i>	21
Ilustración 4 <i>Categorías de Visualizaciones Utilizadas en XAI.</i>	29
Ilustración 5 <i>Comparación entre Modelos Clásicos y Avanzados para la Predicción de Enfermedades Crónicas.</i>	30
Ilustración 6 <i>Diagrama Del Dataset Conceptual.</i>	34
Ilustración 7 <i>Pilares Éticos para el Uso de Datos Biométricos.</i>	38
Ilustración 8 <i>Flujo Metodológico Conceptual del Modelo Propuesto.</i>	41
Ilustración 9 <i>Diagrama Conceptual de Distribución de Variables SHAP.</i>	44
Ilustración 10 <i>Gobernanza de Datos</i>	47

Índice de Tablas

Tabla 1 <i>Comparación de Modelos Clásicos y Avanzados.</i>	31
Tabla 2 <i>Variables Generales Para Dataset Conceptual.</i>	35
Tabla 3 <i>Clasificación de Categorías y Variables Para Dataset Conceptual.</i>	36
Tabla 4 <i>Rendimiento Inferencial Promedio de los Modelos Propuestos.</i>	43

1. Capítulo 1: Introducción.

1.1. Planteamiento del Problema de Investigación.

La diabetes mellitus tipo 2 (DM2) es una de las enfermedades crónicas de mayor prevalencia en Argentina y el mundo, representando un problema creciente para los sistemas de salud. Su detección temprana es fundamental para prevenir complicaciones asociadas, pero los métodos tradicionales dependen de controles clínicos periódicos, que suelen realizarse cuando la enfermedad ya está avanzada.

El desarrollo de tecnologías digitales y la disponibilidad de datos biométricos provenientes de dispositivos médicos, instrumentos clínicos y wearables permiten generar nuevas herramientas predictivas. Las técnicas de Machine Learning (ML) pueden identificar patrones tempranos en variables fisiológicas y de estilo de vida, mejorando la capacidad diagnóstica.

Aún existen desafíos respecto a la precisión, interpretabilidad y aplicabilidad clínica de estos modelos en la población argentina. En este contexto, surge la necesidad de evaluar el potencial del Machine Learning para detectar tempranamente el riesgo de diabetes tipo 2 a partir de datos biométricos.

1.2. Justificación de la Elección del Problema.

La elección de este problema se fundamenta en su alta relevancia sanitaria y social. La DM2 representa un desafío creciente para la salud pública argentina, tanto por el aumento sostenido de casos como por las complicaciones asociadas al diagnóstico tardío. Mejorar la detección temprana permitiría intervenir antes, reducir la carga de enfermedad y mejorar la calidad de vida de la población.

Desde una perspectiva tecnológica, los avances en análisis de datos y modelos de Machine Learning ofrecen herramientas prometedoras para fortalecer los procesos de prevención y alerta epidemiológica. La posibilidad de utilizar datos biométricos accesibles y de bajo costo, brinda un potencial significativo para desarrollar modelos de cribado poblacional más eficientes.

Existe una necesidad académica y profesional de evaluar la pertinencia de estas tecnologías en el contexto argentino, caracterizado por la fragmentación de los sistemas de información y la heterogeneidad en la digitalización de los servicios de salud. Este trabajo busca aportar una base conceptual para orientar futuras implementaciones y apoyar el desarrollo de estrategias preventivas basadas en evidencia.

1.3. Objetivos del Trabajo Final.

1.3.1. Objetivo General.

Evaluar la aplicabilidad de modelos de Machine Learning para la detección temprana de Diabetes Tipo 2 a partir de datos biométricos, considerando su potencial implementación en el sistema de salud argentino.

1.3.2. Objetivos Específicos.

El trabajo final busca cumplir objetivos específicos que se heredan del objetivo general, ellos son:

- Analizar modelos predictivos basados en aprendizaje automático que permitan identificar patrones asociados a diabetes tipo 2.
- Evaluar el impacto potencial de la implementación de estas herramientas en la mejora de la calidad de vida de las personas.

1.4. Preguntas de Investigación.

¿Cómo contribuyen las técnicas de aprendizaje automático aplicadas a datos biométricos a la detección temprana de Diabetes Mellitus Tipo 2?

¿Qué nivel de rendimiento predictivo podrían alcanzar los modelos de aprendizaje automático aplicados a datos biométricos en la detección temprana de Diabetes Mellitus Tipo 2?

1.5. Hipótesis del Problema.

El uso de técnicas de Machine Learning aplicadas a datos biométricos mejora la detección temprana de Diabetes Mellitus Tipo 2, alcanzando niveles de precisión superiores al 85% según la evidencia científica disponible.

1.6. Metodología de Investigación.

La presente investigación adopta un enfoque exploratorio-descriptivo orientado a analizar, desde una perspectiva conceptual, la aplicabilidad del Machine Learning para la detección temprana de Diabetes Mellitus Tipo 2 a partir de datos biométricos. Dado que el estudio no implementa un modelo experimental sobre un dataset real, sino que evalúa la evidencia científica disponible y propone un diseño metodológico teórico, la metodología se estructura en dos componentes centrales:

- Revisión Sistemática de la Literatura: evaluación de la evidencia existente sobre el desempeño de algoritmos en estudios previos.
- Diseño Conceptual: estructuración de una propuesta metodológica conceptual para un modelo predictivo contextualizado al sistema de salud argentino.

El análisis se centra en evaluar la capacidad de las técnicas de Machine Learning para identificar patrones tempranos en variables fisiológicas y de estilo de vida, incluye la

comparación conceptual de algoritmos de aprendizaje supervisado como Random Forest (RF), eXtreme Gradient Boosting (XGBoost) y Redes Neuronales Artificiales (ANN).

1.7. Estrategia de Revisión Bibliográfica.

La revisión sistemática se estructuró siguiendo la metodología PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses), seleccionada por su solidez para conducir análisis exhaustivos y reproducibles de literatura científica. Su aplicación facilitó la evaluación ordenada de estudios sobre técnicas de Machine Learning orientadas a la detección temprana de la Diabetes Mellitus tipo 2, permitiendo documentar con precisión las etapas de identificación, cribado, elegibilidad e inclusión de los trabajos analizados.

En cada etapa se registran los resultados obtenidos, los criterios aplicados y el conjunto final de estudios seleccionados, asegurando la trazabilidad metodológica del proceso.

En esta investigación, la búsqueda se realizó entre 2017 y 2025 sobre las bases de datos de Google Scholar y SciELO, se complementó con herramientas de Inteligencia Artificial (IA) basadas en SciSpace y Claude que permitieron identificar trabajos no indexados en buscadores tradicionales. Se utilizaron combinaciones de palabras clave relacionadas con el dominio:

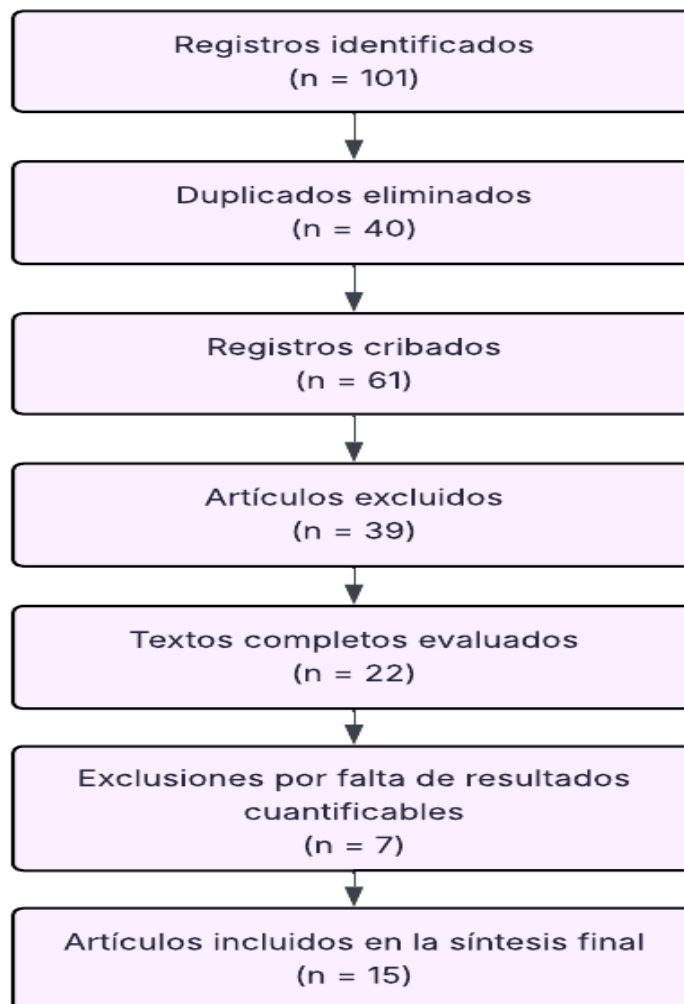
“diabetes mellitus tipo 2”, “Machine Learning”, “predicción”, “detección temprana”, “riesgo clínico”, “modelos supervisados”.

Como resultado de esta búsqueda inicial, se identificaron 101 trabajos, incluyendo artículos científicos, tesis, reportes técnicos y capítulos académicos. Posteriormente, se procedió a la eliminación de duplicados y a un cribado preliminar por título y resumen, descartando documentos que no abordaban técnicas predictivas, que no incluían resultados cuantitativos o que no aplicaban ML.

Los textos potencialmente relevantes fueron evaluados según criterios explícitos de inclusión y exclusión, lo que permitió seleccionar finalmente 15 estudios pertinentes y con datos suficientes para responder las preguntas de investigación planteadas.

Ilustración 1

Diagrama PRISMA del Proceso de Selección de Estudios.



Este diagrama representa visualmente la trazabilidad completa del proceso de selección, basado en lineamientos PRISMA y asegurando la validez metodológica de la Revisión Sistemática de la Literatura (RSL).

1.8. Estructura General del Trabajo Final.

El Trabajo Final se organiza en capítulos que permiten una lectura progresiva desde el planteamiento del problema hasta las conclusiones y propuestas futuras. Cada sección se articula con las preguntas de investigación, las hipótesis y el marco teórico desarrollado.

1.8.1. Capítulos.

Capítulo 1: Introducción.

Presenta el problema, los objetivos, las preguntas de investigación, las hipótesis y el alcance del estudio.

Capítulo 2: Marco Teórico.

Desarrolla los fundamentos clínicos de la DM2, los datos biométricos, la inteligencia artificial aplicada a salud, los algoritmos de ML y la metodología de ciencia de datos.

Capítulo 3: Propuesta Metodológica.

Describe el diseño conceptual del dataset, el preprocesamiento teórico, el modelado y las estrategias de validación basadas en literatura científica.

Capítulo 4: Resultados y Análisis.

Presenta resultados simulados sustentados en evidencia previa, interpretación mediante XAI y análisis crítico de aplicabilidad al sistema de salud argentino.

Capítulo 5: Conclusiones y Líneas Futuras.

Integra los principales hallazgos, reflexiona sobre el potencial del ML en Argentina.

1.8.2. Acrónimos.

Lista ordenada alfabéticamente de las abreviaturas utilizadas en el trabajo, acompañadas de su significado.

1.8.3. Referencias.

Listado completo de fuentes bibliográficas empleadas en el trabajo.

1.9. Estado del Arte.

La aplicación de técnicas de *Machine Learning* en la detección temprana de Diabetes Mellitus Tipo 2 ha experimentado un crecimiento significativo en los últimos años, impulsado por el aumento de datos biométricos disponibles y la necesidad de diagnósticos más oportunos y precisos. Diversos estudios en Latinoamérica, Europa y Argentina muestran cómo los modelos predictivos pueden anticipar el riesgo de DM2 a partir de variables fisiológicas, metabólicas y de estilo de vida.

En México, (De la Rosa-De León et al., 2025) desarrollaron un sistema de prediagnóstico de enfermedades crónicas mediante cómputo inteligente, demostrando que algoritmos como Random Forest, Support Vector Machine (SVM) y Redes Neuronales presentan un rendimiento superior frente a métodos tradicionales al analizar datos biométricos y clínicos tempranos. De manera complementaria, el estudio “Aprendizaje automático aplicado a la detección temprana de DM2: Caso Saltillo” (Lianmoy & Toledo, 2024) confirmó la efectividad del ML en la clasificación del riesgo mediante parámetros como glucosa, presión arterial y actividad física.

En el ámbito latinoamericano, investigaciones de la Universidad del Norte (Marín Ortega & Parra Faria, 2025) y la Universidad Católica de Santa María (Berrios Zuniga, 2024) evaluaron modelos supervisados para predecir diabetes a partir de encuestas de salud y mediciones

fisiológicas, concluyendo que algoritmos basados en árboles de decisión y redes neuronales logran precisiones superiores al 85%. Estos trabajos destacan la importancia del preprocesamiento de datos biométricos y la necesidad de bases de datos balanceadas para mejorar la validez de los modelos.

A nivel regional, el estudio publicado en (Berrios Zuniga, 2024) exploró métodos de clasificación para DM2 utilizando datos estructurados y no estructurados, resaltando el aporte de técnicas de optimización y selección de características para mejorar la predicción del riesgo. A nivel local, la Universidad Nacional de La Plata (UNLP) ha producido dos trabajos relevantes: una tesis orientada a la predicción de enfermedades metabólicas mediante ML, y el reconocido estudio de (Tittarelli, 2023) , quien desarrolló modelos específicos para la población argentina basados en la Encuesta Nacional de Factores de Riesgo (ENFR). Sus resultados mostraron que Random Forest y Redes Neuronales son los enfoques más consistentes para identificar casos de prediabetes (PDM) y Diabetes Mellitus Tipo 2.

En una línea similar, (Perdomo & Ordínez, 2024) analizaron datos de la Encuesta Nacional de Factores de Riesgo, aplicando algoritmos supervisados para identificar variables asociadas al riesgo de diabetes en la provincia de Chubut. El estudio se orienta principalmente a fortalecer la toma de decisiones en salud pública.

En Europa, trabajos como los publicados por (Galán Maroto, 2025) analizan también la aplicabilidad del ML aplicado a DM2 desde una perspectiva clínico-metodológica. Estos aportes enfatizan el uso de técnicas explicables (XAI) como SHAP para interpretar la importancia de los factores predictivos, una tendencia necesaria para el uso clínico real.

En el contexto argentino, un estudio reciente publicado por (Dieuzeide et al., 2025) revisa los factores de riesgo poblacionales asociados al desarrollo de DM2, proporcionando una base epidemiológica esencial para la calibración de modelos predictivos. La combinación de estos factores con algoritmos de ML ha mostrado resultados prometedores para identificación temprana, especialmente cuando se incorporan variables biométricas provenientes de dispositivos móviles o registros clínicos electrónicos.

En conjunto, los trabajos analizados coinciden en que los modelos de ML permiten anticipar la aparición de DM2 con altos niveles de precisión, siempre que se disponga de datos biométricos consistentes y un adecuado tratamiento de valores faltantes. Sin embargo, persisten desafíos como la baja disponibilidad de bases de datos locales amplias, la necesidad de validación clínica en escenarios reales y la inclusión de técnicas explicables que permitan a los profesionales de la salud interpretar adecuadamente los resultados. Los avances actuales muestran un camino sólido hacia herramientas predictivas que podrían integrarse en sistemas de salud para mejorar la prevención y diagnóstico oportuno de diabetes en poblaciones diversas.

La persistencia de desafíos clave se enmarcan en la necesidad de bases de datos más amplias, diversas y estandarizadas, la integración de técnicas explicables y la validación empírica de modelos en entornos clínicos reales.

2. Capítulo 2: Marco Teórico.

2.1. Introducción.

El presente capítulo establece los fundamentos conceptuales que sustentan la detección temprana de la Diabetes Mellitus Tipo 2 mediante técnicas de aprendizaje automático. En primer lugar, se desarrollan las bases clínicas que explican la naturaleza de la enfermedad y los factores biométricos que permiten su predicción. Luego, se describen los principios del aprendizaje supervisado y los algoritmos relevantes empleados en la literatura científica reciente. Finalmente, se analizan las métricas de evaluación y los métodos de interpretación utilizados para garantizar que los modelos predictivos resulten útiles y transparentes en contextos clínicos.

2.2. Fundamentos Clínicos: Diabetes y Variables de Riesgo.

2.2.1. Definición de la Diabetes.

La Diabetes Mellitus Tipo 2 es un trastorno metabólico crónico caracterizado por hiperglucemia persistente (niveles elevados de glucosa en la sangre), consecuencia de la combinación de resistencia a la insulina y una disminución progresiva de la secreción pancreática (Tittarelli, 2023).

Su evolución suele ser lenta y con escasos síntomas iniciales, lo que explica que el diagnóstico clínico frecuentemente ocurra varios años después del inicio de la alteración glucémica (Dieuzeide et al., 2025). Esta identificación tardía incrementa el riesgo de complicaciones metabólicas y cardiovasculares, reforzando la necesidad de estrategias de detección temprana.

2.2.2. Prevalencia de la Diabetes Mellitus 2.

La prevalencia de la diabetes a nivel mundial ha aumentado vertiginosamente, siendo catalogada como una emergencia de salud creciente. Los datos de organizaciones internacionales evidencian la magnitud de esta carga.

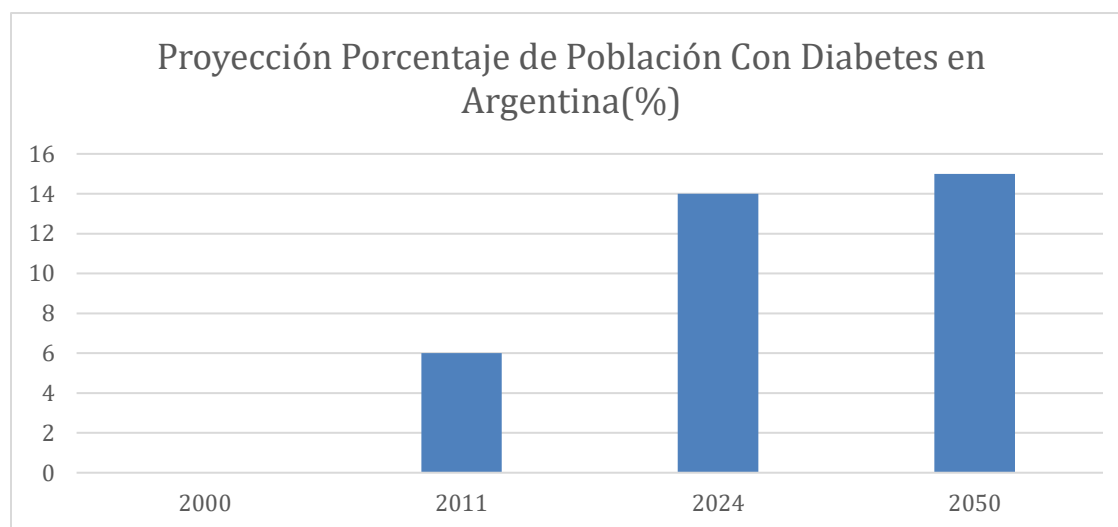
A nivel global, la diabetes afecta a más de 500 millones de adultos, cifra equivalente a cerca del 9 % de la población adulta, con proyecciones que superarán los 640 millones para 2030, este crecimiento sostenido también se replica en América Latina (Aparicio-Montenegro et al., 2025).

En Argentina, según la 4.^a Encuesta Nacional de Factores de Riesgo, la prevalencia de glucemia elevada o diagnóstico de diabetes aumentó del 9,8 % en 2013 al 12,7 % en 2018 (Tittarelli, 2023) (INDEC, 2018).

La Federación Internacional de Diabetes proyecta que esta tendencia continuará ascendiendo durante las próximas décadas.

Ilustración 2

Proyección de Diabetes en la Población Argentina.



Las complicaciones crónicas de la DM2, como las enfermedades cardiovasculares, neurológicas y renales, imponen una carga económica significativa y aumentan drásticamente los costos médicos. Por ejemplo, los eventos cerebrovasculares pueden multiplicar el costo del tratamiento basal de la DM2 hasta por 47.8 veces en el sistema de salud argentino. La detección temprana mediante modelos predictivos busca mitigar estos costos, ya que el diagnóstico tardío tiene un impacto drástico sobre la salud pública (Tittarelli, 2023).

2.2.3. La Importancia de la Detección Temprana y la Prediabetes.

La detección temprana y la identificación del estadio de prediabetes se han posicionado como un enfoque estratégico y crucial para mitigar la carga de la Diabetes Mellitus Tipo 2 a nivel global. La relevancia de este enfoque radica en la oportunidad única que ofrece para intervenir y modificar la trayectoria de una enfermedad que, por su naturaleza, tiende a ser diagnosticada tardíamente, cuando las complicaciones ya están en curso.

La DM2 y la Prediabetes no diagnosticadas o mal controladas son el origen de complicaciones crónicas y agudas que deterioran la calidad de vida y aumentan drásticamente los costos sanitarios:

Complicaciones Cardiovasculares: Son la causa más común de muerte e incapacidad en personas con DM2. Incluyen infarto de miocardio, accidente cerebrovascular (ACV) y enfermedad arterial periférica.

Afectaciones Neurológicas y Renales: La hiperglucemia prolongada provoca daño en los nervios (neuropatía) y en los vasos sanguíneos renales (nefropatía), pudiendo llevar a insuficiencia renal crónica.

Carga Económica: Las complicaciones derivadas del manejo tardío o inadecuado pueden incrementar exponencialmente los costos médicos, siendo un problema estructural en poblaciones como la argentina.

La oportunidad de detección temprana se basa en las evidencias científicas, donde demuestran que las intervenciones oportunas destinadas a generar un cambio en el estilo de vida como la alimentación y la actividad física son capaces de retrasar, prevenir, e incluso revertir el desarrollo de la DM2 y sus complicaciones (Tittarelli, 2023) . Por esta razón, la detección temprana para su control es considerada un desafío de suma importancia.

2.2.4. Datos Biométricos y Factores Predictivos.

Para la predicción de la diabetes, se utilizan datos clínicos y biomédicos de pacientes con y sin diagnóstico.

Las variables utilizadas en la investigación incluyen:

- Datos Biométricos/Clínicos: Sexo, Edad, Peso, Talla, IMC, Glucosa.
- Historial y Comportamiento: Antecedentes familiares diabéticos, consumo de alcohol, consumo de tabaco, consumo de drogas, actividad física.

El criterio médico establece que en la predicción de diabetes es fundamental considerar el nivel de glucosa y el nivel de índice de masa corporal (IMC) de una persona (Berrios Zuniga, 2024).

Otros factores fuertemente asociados con la diabetes incluyen la presión arterial alta, el colesterol elevado, la dificultad para caminar, la edad y la percepción de salud general. Además, factores socioeconómicos como sectores sociales de bajo/escasos recursos, bajo nivel de

desarrollo educativo y mala alimentación son factores antropológicos que poseen una alta influencia en el desarrollo de la enfermedad (Marín Ortega & Parra Faria, 2025).

La predicción efectiva de la Diabetes Mellitus Tipo 2 y la Prediabetes mediante técnicas de ML se sustenta en la explotación de un conjunto de atributos multifactoriales, los cuales se agrupan en categorías que reflejan las dimensiones biológicas, demográficas y conductuales del riesgo. La selección de estas variables para la predicción del riesgo de DM2 se justifica en su soporte clínico y su validación estadística por diversos métodos de análisis de relevancia (Pearson, Chi-cuadrado, ANOVA e Información Mutua), donde las variables más importantes son aquellas que demuestran ser predictoras en al menos dos de estos análisis.

Ilustración 3

Prevalencia, Factores y Complicaciones.

PREVALENCIA, FACTORES Y COMPLICACIONES		
CATEGORÍAS	CONCEPTOS CLAVE	
 Prevalencia	<ul style="list-style-type: none"> Alta incidencia mundial Incremento en países de ingresos medios Crecimiento sostenido Tendencia ascendente en Argentina 	
 Factores Biométricos	<ul style="list-style-type: none"> Glucosa en sangre Circunferencia de cintura Perfil lipídico Índice de Masa Corporal (IMC) Presión arterial 	
 Factores Demográficos	<ul style="list-style-type: none"> Edad Antecedentes familiares Sexo Condiciones socioeconómicas 	
 Factores Conductuales	<ul style="list-style-type: none"> Actividad física insuficiente Tabaquismo Alimentación inadecuada Consumo de alcohol 	
 Complicaciones	<ul style="list-style-type: none"> Enfermedad cardiovascular Nefropatía Retinopatía Accidente cerebrovascular Neuropatía 	
 Detección Temprana	<ul style="list-style-type: none"> Identificación de prediabetes Mejora del pronóstico Reducción de complicaciones Disminución del costo sanitario 	

2.3. Fundamentos Tecnológicos: Aprendizaje Automático (ML).

2.3.1. Inteligencia Artificial y Machine Learning.

La predicción de DM2 y PDM se enmarca en el ámbito de la salud digital y la ciencia de datos, disciplinas que encuentran su fundamento metodológico en la Inteligencia Artificial y su subcampo de mayor crecimiento, el Machine Learning,

Inteligencia Artificial se define como la disciplina dedicada al estudio y diseño de sistemas que buscan emular o replicar las capacidades cognitivas y comportamientos inteligentes humanos. Un sistema catalogado como "inteligente" debe ser capaz de interpretar correctamente los datos, aprender de ellos y emplear los conocimientos adquiridos para realizar acciones que maximicen sus posibilidades de éxito en tareas concretas de forma adaptativa. La IA se sustenta en algoritmos y modelos matemáticos que permiten a las computadoras entrenar y aprender de los datos para tomar decisiones que se asemejan a la inteligencia humana.

Dentro de este amplio dominio, el Aprendizaje Automático se ha establecido como la rama más productiva e importante de la IA. El ML se centra en el estudio de mecanismos que confieren a las máquinas la capacidad de aprender, sin necesidad de ser programadas explícitamente para cada tarea.

Mientras que en la programación clásica se ingresan reglas y datos para obtener respuestas, en el paradigma del ML, el sistema recibe datos y las respuestas esperadas, y a partir de esa experiencia, induce las reglas que luego pueden ser aplicadas a nuevos datos para generar respuestas originales. Por lo tanto, el ML se enfoca en la detección automática de patrones relevantes dentro de un conjunto de datos. Este proceso de aprendizaje se mide mediante una

métrica de rendimiento (P) en relación con una clase de tareas (T), donde la mejora en el rendimiento con la experiencia es la definición formal del aprendizaje automático.

La capacidad del ML para analizar grandes volúmenes de datos y encontrar patrones no lineales lo ha posicionado como una herramienta crucial en el sector salud. Específicamente, el ML se aplica a la medicina para:

Detección y Diagnóstico Temprano: Implementando modelos predictivos y herramientas impulsadas por IA destinadas a la detección de enfermedades. Modelos supervisados como Random Forest, XGBoost, y Redes Neuronales Artificiales han demostrado un gran potencial para predecir la DM2 y otras patologías crónicas.

Pronóstico y Estratificación de Riesgo: Ayudando a los profesionales de la salud a la toma de decisiones informadas y eficientes, identificando a individuos con alto riesgo o prediciendo complicaciones asociadas a la enfermedad.

Optimización del Tratamiento: Mediante el análisis de datos continuos (como la monitorización de glucosa) para personalizar el tratamiento y optimizar la asignación de recursos.

2.3.2. Algoritmos de Aprendizaje Supervisado.

La revisión de la literatura indica que las técnicas de ML son clave para predecir y diagnosticar la diabetes de forma rápida y eficiente (Marín Ortega & Parra Faria, 2025).

Los modelos utilizados para el diagnóstico temprano de la diabetes son métodos supervisados, incluyendo:

Random Forest: Es un método de aprendizaje conjunto que se utiliza tanto para clasificación como para regresión. Funciona construyendo múltiples árboles de decisión durante el proceso de entrenamiento y combinando sus predicciones (votación mayoritaria para clasificación). El diseño de RF está específicamente ideado para mejorar la precisión general del modelo y, fundamentalmente, para controlar la tendencia de los árboles de decisión individuales al sobreajuste (overfitting).

XGBoost: Es un modelo de boosting (mejora progresiva) que pertenece a la familia de los modelos de ensamblaje. Este algoritmo se ha distinguido por su alto rendimiento predictivo en la clasificación de datos tabulares, a menudo superando a otros métodos. Se construye de forma aditiva, donde cada nuevo árbol de decisión se ajusta sobre los errores residuales cometidos por los modelos previos. XGBoost utiliza una función objetivo con un componente de pérdida y un componente de regularización que controla la complejidad del modelo, haciéndolo robusto. Estudios comparativos han demostrado que XGBoost (especialmente cuando se combina con técnicas de balanceo como Synthetic Minority Oversampling Technique (SMOTE)) obtiene una capacidad predictiva muy elevada en el diagnóstico de diabetes y prediabetes.

Redes Neuronales Artificiales: Son modelos avanzados no lineales que constan de capas de entrada, capas ocultas y una capa de salida. Las redes neuronales, debido a su gran número de parámetros, poseen una elevada capacidad para capturar patrones complejos y no lineales en los datos, lo que las hace adecuadas para el diagnóstico médico. Las redes neuronales, incluyendo aquellas basadas en deep learning, han demostrado un desempeño superior en la clasificación de pacientes de alto riesgo de DM2, en casos han alcanzado precisiones por encima del 90%.

2.4. Modelización Predictiva en el Contexto Clínico.

2.4.1. Métricas de Evaluación Clave (Sensibilidad vs. Precisión).

El problema de la detección temprana de Diabetes Mellitus Tipo 2 se aborda como un problema de clasificación supervisada. En este contexto, la evaluación de los modelos predictivos debe ir más allá de la métrica de Exactitud (Accuracy), la cual mide la proporción de predicciones correctas (Verdaderos Positivos + Verdaderos Negativos) respecto al total de las observaciones.

En el ámbito de la salud, la Accuracy resulta insuficiente y potencialmente engañosa debido al problema estructural del desequilibrio de clases inherente a los datasets poblacionales (la mayoría de los individuos son de clase negativa o "sanos"). Un modelo puede alcanzar una Accuracy alta, pero fallar consistentemente en identificar a la minoría (los casos reales de DM2), lo que anula su utilidad clínica real (Galán Maroto, 2025).

Dada esta limitación, la métrica esencial para la detección temprana de DM2 es la Sensibilidad (Recall o Exhaustividad).

La Sensibilidad se define como la proporción de Verdaderos Positivos (VP) con respecto a todos los positivos reales (VP + Falso Negativo (FN)). Mide la capacidad del modelo para identificar correctamente a los individuos que realmente tienen riesgo de DM2, siendo crucial para la captación temprana de pacientes.

Si un modelo alcanza una Exactitud elevada, pero falla consistentemente en identificar a la minoría (los casos reales de DM2). Por lo tanto, no refleja la utilidad clínica real del modelo cuando el objetivo primordial es detectar patrones de enfermedad poco frecuentes.

La Sensibilidad se prioriza en la medicina preventiva para minimizar el costo del Falso Negativo, que es la omisión de un caso real de enfermedad o riesgo:

Un Falso Negativo ocurre cuando el valor real de un paciente es positivo (es decir, tiene DM2 o riesgo), pero la predicción del modelo lo clasifica incorrectamente como Verdadero Negativo (sano).

En el cribado poblacional de una enfermedad crónica y progresiva como la DM2, la omisión de un caso real tiene el mayor costo clínico y social. Un FN implica que un paciente en riesgo o ya enfermo continuará sin tratamiento ni intervención, lo que aumenta el riesgo de desarrollar complicaciones graves (cardiovasculares, renales, neurológicas) y, consecuentemente, impone una carga económica significativa al sistema de salud.

El objetivo clínico principal es evitar que se omita un caso positivo. Por lo tanto, en la selección los modelos deben ser optimizados para alcanzar la mayor sensibilidad posible, incluso si esto implica aceptar un ligero aumento en los Falsos Positivos (FP), ya que la consecuencia de un paciente sano clasificado como enfermo, que requeriría una prueba de seguimiento es menos grave que la consecuencia de un FN.

2.4.2. La Interpretación de los Modelos (Explainable AI - XAI).

El desarrollo de modelos de ML para la predicción de la Diabetes Mellitus Tipo 2 ha demostrado consistentemente un alto potencial para mejorar el rendimiento predictivo. No obstante, los algoritmos que exhiben el mejor desempeño particularmente los modelos de ensamblaje (ensemble) y los avanzados, como XGBoost o ANN son, por diseño, complejos y difíciles de interpretar. Esta característica da origen al dilema de la "Caja Negra", donde sistemas que alcanzan una precisión superior, pero cuyos procesos internos y las razones detrás de sus

decisiones son difíciles o imposibles de rastrear o justificar para el usuario final (Galán Maroto, 2025).

La complejidad de estos modelos, cuya lógica interna es intrínseca a miles de parámetros o estructuras aditivas, puede dificultar su interpretación y aplicación en entornos clínicos reales. Esa complejidad se convirtió en una barrera que separó las herramientas de ML de los entornos sanitarios, donde la transparencia es un requisito indispensable para la aceptación y uso por parte del personal médico.

La Inteligencia Artificial Explicable surge como un campo de estudio dedicado a resolver el dilema de la "caja negra", definiéndose como un conjunto de técnicas y metodologías que buscan hacer los modelos de ML transparentes, auditables y comprensibles para los humanos.

El propósito central de la XAI no es reducir la complejidad del modelo, sino proporcionar explicaciones post-hoc que permitan entender la contribución de las variables al resultado final. Las técnicas de XAI, como los valores de SHAP (SHapley Additive exPlanations) y LIME (Local Interpretable Model-Agnostic Explanations), permiten descomponer la predicción de un modelo complejo y estimar la contribución de cada variable a una decisión individual. La transparencia, facilitada por la XAI, es crucial para la adopción y confianza de los profesionales de la salud, ya que permite la validación de la utilidad clínica del modelo.

Los diagramas y representaciones visuales que respaldan el análisis XAI en estos estudios se pueden agrupar en cuatro categorías principales:

1. **Importancia global de características (Feature Importance):** incluye diagramas que muestran el peso relativo de cada variable en el modelo, generalmente mediante gráficos de barras o rankings de variables. Son frecuentes en modelos basados en

árboles, como Random Forest o XGBoost, y permiten verificar si los predictores más influyentes coinciden con la evidencia clínica, aportando transparencia global sobre el comportamiento del algoritmo.

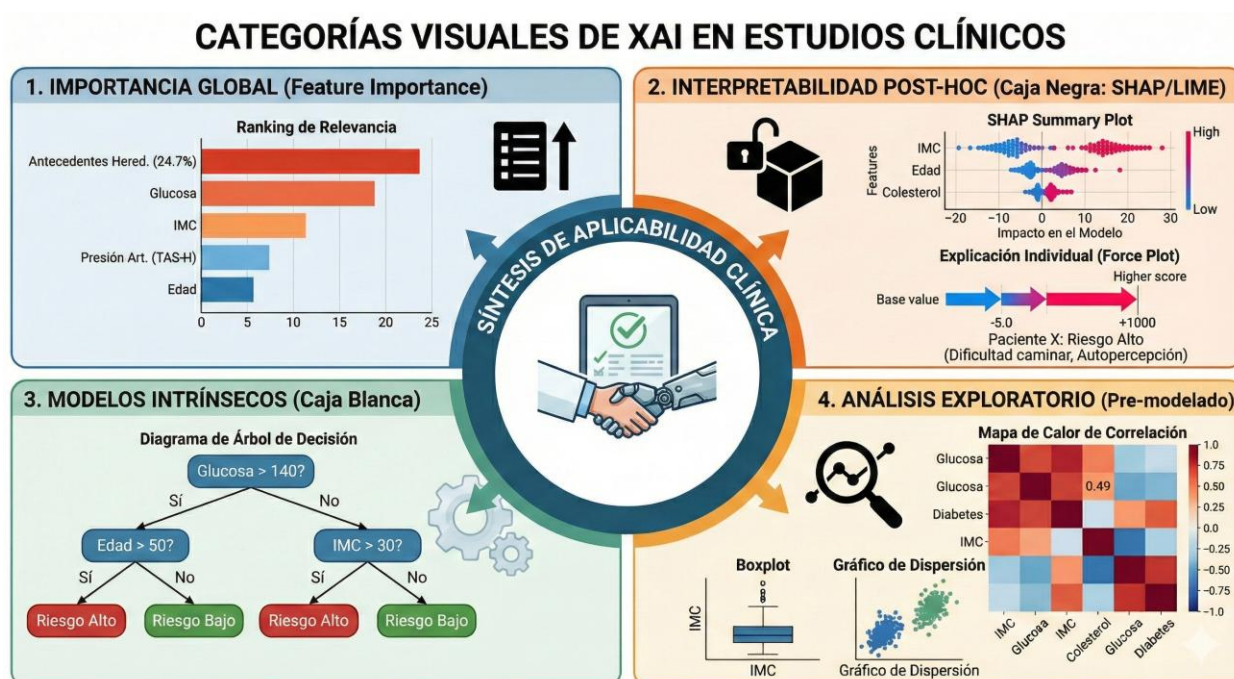
2. **Interpretabilidad local post-hoc (SHAP y LIME):** agrupa técnicas diseñadas para explicar modelos complejos después de su entrenamiento. SHAP proporciona explicaciones globales y locales mediante valores que representan la contribución marginal de cada característica, mientras que LIME ofrece interpretaciones locales aproximando el modelo con funciones lineales simples. Ambas herramientas permiten comprender por qué un caso individual recibe una determinada predicción.
3. **Visualización de modelos intrínsecamente interpretables:** comprende diagramas derivados de modelos cuya estructura es transparentemente legible. Ejemplos típicos son los árboles de decisión, que representan rutas “si–entonces”, y los mapas cognitivos difusos, que permiten visualizar relaciones causales entre factores de riesgo. Estas representaciones facilitan la lectura directa del razonamiento clínico-computacional.
4. **Análisis exploratorio y diagramas de correlación (pre-modelado):** expresan visualizaciones empleadas antes del entrenamiento, como mapas de calor de correlación, gráficos de dispersión, pairplots o boxplots. Estos diagramas permiten identificar relaciones entre variables, detectar patrones relevantes y reconocer valores atípicos, aportando un entendimiento inicial de la estructura del dataset que luego repercute en la interpretabilidad del modelo final.

En conjunto, estas cuatro dimensiones conforman el núcleo visual del análisis XAI en estudios de predicción de DM2.

La siguiente ilustración sintetiza estas categorías y su función dentro del proceso explicativo.

Ilustración 4

Categorías de Visualizaciones Utilizadas en XAI.



Fuente: Generada con IA

2.5. Comparación de Modelos Predictivos: Clásicos vs Avanzados.

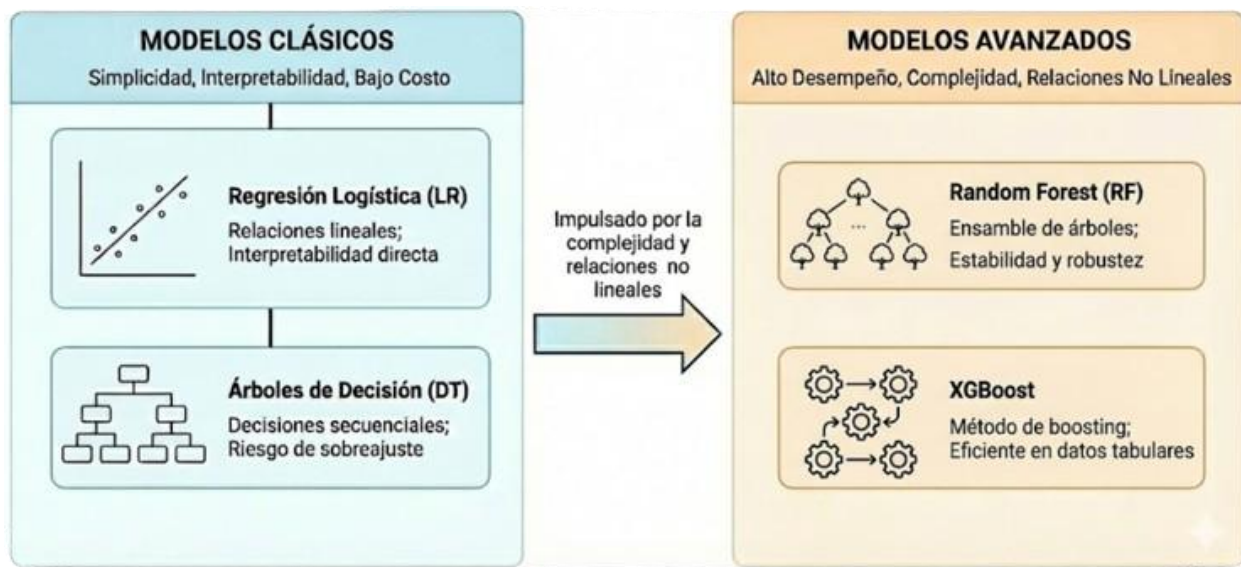
La aplicación del Machine Learning a problemas de clasificación binaria, como la predicción de la DM2, emplea un espectro de algoritmos que varían en complejidad, rendimiento y transparencia. Es necesario distinguir entre los modelos clásicos, que históricamente han sido la base de la bioestadística, y los modelos avanzados, que han dominado la competencia de ciencia de datos en la última década.

Los modelos clásicos, como la Regresión Logística (RL) y el Árbol de Decisión (DT), destacan por su simplicidad y alta interpretabilidad. La RL es un modelo lineal que estima la probabilidad de que una variable dependiente binaria ocurra, siendo un estándar de oro en epidemiología y bioestadística debido a que sus coeficientes tienen una clara interpretación clínica.

El Árbol de Decisión ofrece una visualización intuitiva del proceso de decisión, imitando la lógica de diagnóstico clínico basada en umbrales sucesivos de variables, por ejemplo: Si $IMC > X$ y $Edad < Y$ entonces Z .

Ilustración 5

Comparación entre Modelos Clásicos y Avanzados para la Predicción de Enfermedades Crónicas.



La sencillez de los modelos clásicos está acompañada de limitaciones, donde los algoritmos de RL, asumen una relación lineal entre las variables de riesgo de DM2, lo que no es siempre correcto en ciertos escenarios de las complejas interacciones biológicas de la enfermedad.

Los modelos de DT son susceptibles a los sobreajustes, lo que se traduce en inestabilidad de resultados en base a nuevos datos, denotan una notable desventaja de adaptación entre ámbitos controlados de testeo y los ambientes reales clínicos.

Los modelos avanzados, especialmente los métodos de *Ensemble* como Random Forest y XGBoost, abordan las limitaciones de los enfoques clásicos.

Un método *Ensemble* combina las predicciones de múltiples modelos base, generalmente Árboles de Decisión, para obtener un resultado final más robusto y preciso.

Random Forest utiliza el concepto de *Bagging*, entrenando muchos árboles en subconjuntos aleatorios de datos y variables. Esto reduce la varianza del modelo, mitigando el sobreajuste.

XGBoost utiliza el concepto de Boosting, donde cada árbol nuevo corrige secuencialmente los errores cometidos por los árboles anteriores. Este enfoque está optimizado para la velocidad y la precisión, este modelo generalmente produce los mejores rendimientos en problemas de clasificación estructurada.

Tabla 1

Comparación de Modelos Clásicos y Avanzados.

Característica	LR	DT	RF	XGBoost
Capacidad para capturar no linealidad	Baja	Media	Alta	Muy alta
Interpretabilidad	Alta	Alta	Media	Media (mejorable con XAI)
Riesgo de sobreajuste	Bajo	Alto	Bajo-medio	Bajo
Rendimiento típico en datos biométricos	Moderado	Moderado	Alto	Muy alto
Requerimiento computacional	Muy bajo	Bajo	Medio	Medio

En el contexto de la detección temprana de DM2, la capacidad para modelar interacciones no lineales y el alto rendimiento predictivo que ofrecen algoritmos como XGBoost y Random Forest justifican su elección, ya que la prioridad clínica es la máxima Sensibilidad para identificar a los pacientes en riesgo, incluso si esto requiere herramientas adicionales, como XAI, para garantizar su interpretabilidad.

3. Capítulo 3: Diseño Experimental y Metodología

3.1. Introducción.

El presente capítulo establece el Diseño Metodológico Conceptual adoptado para evaluar la aplicabilidad de los modelos de Machine Learning en la Detección Temprana de Diabetes Mellitus Tipo 2. Dado que este Trabajo Final adopta un enfoque exploratorio sobre la base de la evidencia de la literatura (Capítulo 1.8 y 2), no se describe la ejecución de un modelo computacional sobre un dataset crudo; en cambio, se propone una arquitectura metodológica.

3.2. Diseño Conceptual del Dataset.

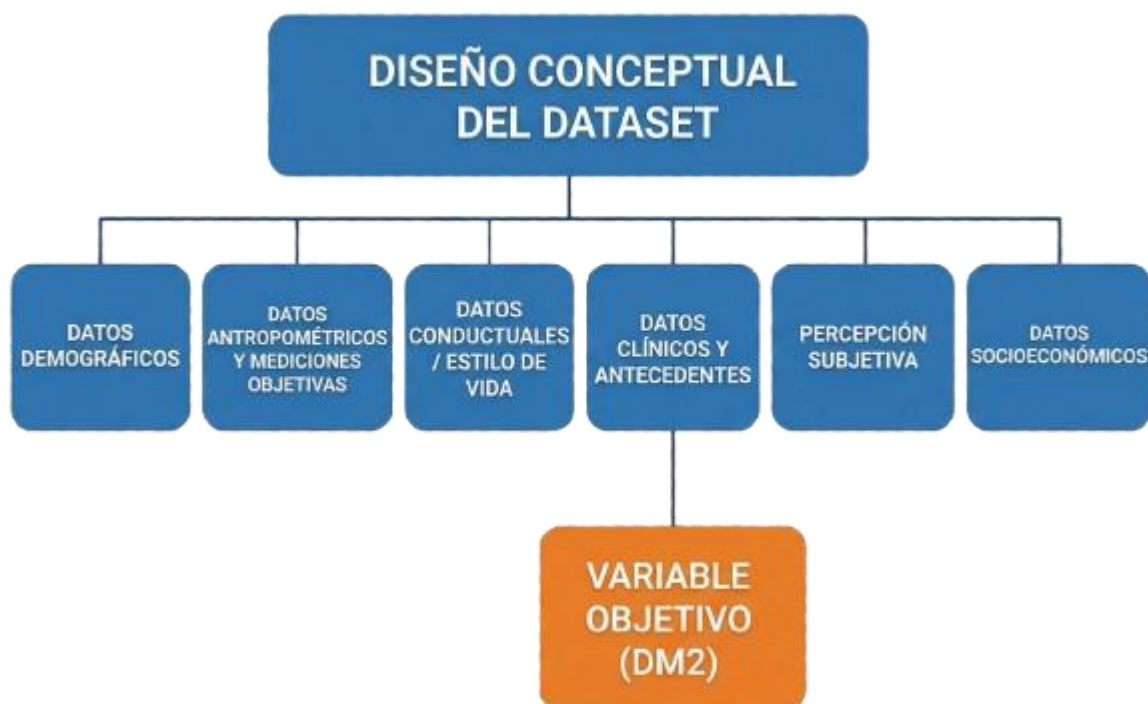
La definición conceptual del dataset es un paso crítico en todo estudio predictivo basado en aprendizaje automático. Debido a la ausencia de un dataset clínico propio, se propone utilizar la Encuesta Nacional de Factores de Riesgo como base teórica para el diseño del conjunto de datos, ya que esta fuente constituye “el insumo epidemiológico de mayor representatividad poblacional para el análisis de enfermedades crónicas en Argentina” (Tittarelli, 2023, pág. 1).

Su propósito es contar con un conjunto de atributos que representen adecuadamente los determinantes asociados al riesgo de desarrollar DM2.

La **Ilustración 6** presenta el diagrama conceptual del dataset, que organiza las variables en grandes bloques temáticos derivados de la ENFR. Este esquema permite visualizar de manera sintética la composición general del conjunto de datos y su coherencia con el enfoque epidemiológico típico de los estudios poblacionales.

Ilustración 6

Diagrama Del Dataset Conceptual.



Nota. El diseño del diagrama es hipotético basado en la estructura de la Encuesta Nacional de Factores de Riesgo.

A partir de esta estructura general, la **Tabla 2** organiza las variables en cuatro categorías principales, cada categoría agrupa las variables relevantes y explica su rol potencial como predictoras de DM2.

Tabla 2*Variables Generales Para Dataset Conceptual.*

Categoría	Variables Incluidas	Descripción / Rol Conceptual
1. Factores Demográficos y Antropométricos	Edad, grupo etario, sexo, peso, talla, IMC, circunferencia de cintura.	Representan el perfil básico de la persona y permiten caracterizar el riesgo cardiometabólico. La obesidad central (circunferencia cintura) y el IMC son predictores directos de DM2.
	Clasificación de IMC, clasificación de riesgo por cintura.	Variables derivadas utilizadas para estandarizar el análisis según criterios clínicos.
2. Estilo de Vida y Comportamiento	Actividad física semanal, minutos de actividad, sedentarismo, consumo de tabaco, exposición al humo ajeno.	Determinan el nivel de riesgo asociado a hábitos conductuales. El sedentarismo y el tabaquismo son factores consistentemente asociados a mayor probabilidad de diabetes.
	Consumo de alcohol, frecuencia de ingesta, consumo de bebidas azucaradas, consumo de frutas y verduras.	Reflejan calidad alimentaria y patrones nutricionales; se vinculan estrechamente a obesidad y resistencia a la insulina.
3. Indicadores Clínicos y Antecedentes	Diagnóstico previo de hipertensión, medicación antihipertensiva.	La hipertensión forma parte del síndrome metabólico y aumenta fuertemente el riesgo de DM2.
	Colesterol total, colesterol elevado previo, medicación para dislipidemias.	Indicadores metabólicos esenciales para evaluar riesgo cardiovascular global.
	Antecedentes familiares de diabetes, hiperglucemia previa detectada.	Predictores críticos: los antecedentes familiares reflejan susceptibilidad genética y la hiperglucemia previa indica progresión metabólica.
4. Factores Socioeconómicos y Percepción de Salud	Nivel educativo, nivel de ingresos del hogar, cobertura de salud.	Determinan desigualdades estructurales que afectan la prevención, diagnóstico y tratamiento.
	Percepción general de salud, dificultad para caminar, limitaciones funcionales.	Funcionan como indicadores subjetivos y funcionales del estado físico general. Son predictores indirectos del deterioro metabólico.

Nota. Las variables presentadas son de carácter conceptual y se basan en la estructura habitual de la Encuesta Nacional de Factores de Riesgo. No corresponden a datos reales ni a una ejecución experimental.

Presentamos en la **Tabla 3**, la clasificación de las categorías y variables basándonos en la categorización de cada generalidad, con el objetivo de alcanzar el nivel de detalle necesario que

Tabla 3

Clasificación de Categorías y Variables Para Dataset Conceptual.

Datos Demográficos		
Variable	Descripción	Tipo
Edad	Años cumplidos al momento de la encuesta	N Numérica
Sexo	Masculino / Femenino	C Categórica
Nivel educativo	Máximo nivel educativo alcanzado	C Categórica
Situación laboral	Condición de actividad (empleado, desocupado, inactivo)	C Categórica
Cobertura de salud	Tipo de cobertura (obra social, prepaga, pública)	C Categórica
Datos Antropométricos y Mediciones Objetivas		
Variable	Descripción	Tipo
Peso	Peso corporal en kilogramos	N Numérica
Talla	Altura en centímetros	N Numérica
IMC	Índice de masa corporal calculado	N Numérica derivada
Circunferencia de cintura	Medición en centímetros	N Numérica
Presión arterial sistólica	Medición en mmHg	N Numérica
Presión arterial diastólica	Medición en mmHg	N Numérica
Glucemia capilar	Valor en mg/dL	N Numérica
Colesterol total	Valor en mg/dL	N Numérica
Datos Conductuales / Estilo de Vida		
Variable	Descripción	Tipo
Tabaquismo	Consumo actual de tabaco	C Categórica
Exposición a humo ajeno	Exposición en el hogar o trabajo	C Categórica
Consumo de alcohol	Frecuencia e intensidad	C Categórica ordinal
Actividad física	Minutos semanales de actividad moderada/vigorosa	N Numérica
Consumo de frutas y verduras	Porciones diarias	C Categórica ordinal
Consumo de bebidas azucaradas	Frecuencia semanal	C Categórica ordinal
Datos Clínicos y Antecedentes		
Variable	Descripción	Tipo
Diagnóstico previo de diabetes	Reportado por profesional de salud	C Categórica
Diagnóstico previo de hipertensión	Declarado por el encuestado	C Categórica
Colesterol elevado	Diagnóstico o medición previa	C Categórica
Antecedentes familiares de diabetes	Presencia en familiares directos	C Categórica
Uso de medicación	Medicación para DM2, HTA o colesterol	C Categórica
Percepción Subjetiva		
Variable	Descripción	Tipo
Percepción de salud	Autoevaluación del estado general de salud	C Categórica ordinal
Dificultad para caminar	Limitaciones en movilidad	C Categórica
Limitaciones funcionales	Restricciones en actividades diarias	C Categórica
Datos Socioeconómicos		
Variable	Descripción	Tipo
Nivel de ingresos	Tramo de ingreso del hogar	C Categórica ordinal
Hacinamiento	Número de personas por ambiente	N Numérica
Acceso a servicios básicos	Disponibilidad en el hogar	C Categórica
Variable Objetivo		
Variable	Descripción	Tipo
Diabetes Mellitus Tipo 2	1 = Diagnóstico previo / glucosa elevada 0 = Sin diagnóstico	B Binaria

nos permita modelar el dataset a nivel de estructura de datos, alineado con las dimensiones epidemiológicas empleadas por la ENFR y con los requerimientos del modelado predictivo propuesto.

3.3. Consideraciones Éticas en el Uso de Datos Biométricos.

El diseño de sistemas predictivos basados en Machine Learning para el ámbito de la salud conlleva responsabilidades éticas que trascienden la mera precisión técnica. Dado que el modelo propuesto procesa información sensible vinculada a la salud física y los hábitos de vida de las personas, la propuesta del diseño metodológico se adhiere estrictamente a principios bioéticos y normativas de protección de datos, estableciendo un marco de referencia para garantizar la integridad y la dignidad de los potenciales pacientes.

La arquitectura del manejo de datos se fundamenta en el principio de privacidad por diseño. Aunque este estudio plantea un escenario conceptual, cualquier implementación real de ámbito nacional exige el cumplimiento riguroso de la Ley 25.326 de Protección de Datos Personales (Ley 25326, 2000).

Esto implica que el tratamiento de variables biométricas debe realizarse bajo protocolos de anonimización, asegurando que ninguna predicción pueda ser re-identificada o vinculada a un individuo específico fuera del entorno clínico autorizado. El consentimiento informado, en este contexto, no es solo un requisito legal, sino un imperativo ético para la recolección de datos primarios.

Existe el riesgo de que los modelos de ML repliquen o amplifiquen desigualdades sanitarias preexistentes si los datos de entrenamiento no son representativos de la diversidad demográfica y socioeconómica de la población argentina. En consecuencia, la metodología

contempla estrategias de balanceo de datos y validación estratificada para asegurar que el modelo no discrimine ni reduzca su rendimiento en subgrupos vulnerables, garantizando una equidad predictiva y minimizando la posibilidad de introducir sesgos algorítmicos en el modelo.

La transparencia algorítmica se establece como un deber ético. La propuesta incorpora técnicas de Inteligencia Artificial Explicable, como los valores SHAP, no solo como una herramienta técnica, sino como un mecanismo ético para garantizar el derecho a la explicación.

Esto asegura que la tecnología actúe como una herramienta de soporte a la decisión clínica humana, manteniendo siempre la responsabilidad final en el profesional de la salud y no en el algoritmo.

Ilustración 7

Pilares Éticos para el Uso de Datos Biométricos.



3.4. Preprocesamiento y Limpieza de Datos.

Las técnicas descritas a continuación no fueron aplicadas a un dataset concreto, sino que se presentan como la metodología teórica que regiría un estudio aplicado. Se basan en prácticas recomendadas para datos biométricos y encuestas poblacionales.

3.4.1. Estrategia de Imputación.

Los valores faltantes constituyen un desafío común en encuestas poblacionales como la ENFR y en registros clínicos. Para mitigar el impacto de los datos incompletos, se propone aplicar una estrategia de imputación de la siguiente manera:

Imputación mediante la Mediana (Variables Numéricas): Se propone utilizar la mediana para imputar los valores faltantes en variables numéricas continuas (como el IMC y la Edad). Esta técnica es sólida frente a la presencia de valores fuera de rango o raros (*outliers*) y los errores de recolección de datos. Imputar con la mediana evita el sesgo que introduciría la media en distribuciones asimétricas, siendo un método recomendado en modelos predictivos clínicos.

Imputación mediante la Moda (Variables Categóricas): Para variables categóricas o discretas (ej.: la presencia de tabaquismo o el nivel de actividad física), se propone utilizar la moda, es decir, el valor más frecuente.

Esta estrategia fue seleccionada por su robustez y porque coincide con los enfoques más utilizados en los estudios basados en ENFR (INDEC, 2018).

3.4.2. Transformación de Variables.

El estudio adopta transformaciones estándar en ML:

- 1- Codificación de Variables Categóricas (Encoding).

Se propone utilizar la codificación One-Hot Encoding para variables categóricas nominales (sin orden intrínseco), transformando cada categoría en una nueva columna binaria. Para variables binarias o categóricas ordinales, se propone un Label Encoding. Esta transformación es esencial para que los modelos basados en árboles y redes neuronales puedan procesar los atributos.

2- Escalado de Variables Numéricas (Feature Scaling).

Se propone utilizar la normalización Min-Max en variables numéricas como el IMC y la Edad. Este paso es importante para que los modelos sensibles a la magnitud, como las ANN, no otorguen un peso desproporcionado a variables con rangos numéricos más amplios.

3.4.3. Partición Estratificada.

La división conceptual del dataset es fundamental para evitar el sobreajuste y garantizar la validez externa del modelo. Se propone el siguiente esquema de partición:

División: 80% para Entrenamiento y 20% para Prueba.

Estrategia: La división debe ser estratificada por la variable objetivo (Diagnóstico de diabetes).

Justificación: Dado que la diabetes presenta un desbalance significativo de clases (poca proporción de casos positivos), la estratificación asegura que la proporción de casos positivos y negativos se mantenga idéntica en los subconjuntos de entrenamiento y prueba. Esta estrategia es vital para evitar el sesgo y mejorar la validez del modelo.

3.5. Flujo Metodológico General.

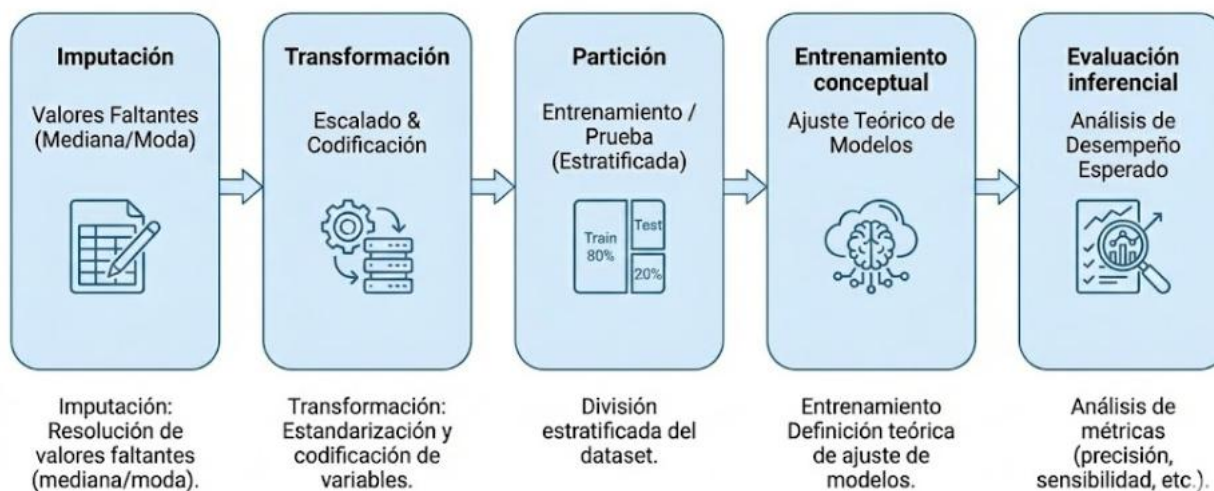
El proceso metodológico propuesto puede sintetizarse en un flujo conceptual compuesto por cinco etapas encadenadas: imputación de valores faltantes, transformación de variables, partición del dataset, entrenamiento conceptual y evaluación inferencial.

Estas etapas representan la secuencia lógica que seguiría una implementación real de un modelo predictivo, aunque en este estudio se aplican únicamente a nivel teórico, sin ejecución sobre datos clínicos.

La **Ilustración 8** resume visualmente este flujo metodológico y muestra la relación entre cada fase dentro del proceso general de modelado.

Ilustración 8

Flujo Metodológico Conceptual del Modelo Propuesto.



4. Capítulo 4: Resultados y Evaluación

4.1. Introducción.

El presente capítulo expone los resultados inferenciales derivados del diseño metodológico conceptual desarrollado previamente. Dado que este estudio no ejecuta un entrenamiento real sobre un dataset clínico argentino, los valores aquí reportados son estimaciones basadas en rangos de desempeño ampliamente documentados en la literatura científica reciente sobre predicción de Diabetes Mellitus Tipo 2.

Los resultados cumplen tres funciones centrales:

- Seleccionar el modelo óptimo considerando métricas clínicamente relevantes.
- Presentar una interpretación conceptual del comportamiento del modelo mediante XAI, evitando el carácter de “caja negra”.
- Analizar los desafíos y condiciones necesarias para la implementación real en el contexto sanitario argentino.

4.2. Resultados de la Evaluación Inferencial y Selección del Modelo Óptimo.

Los estudios actuales sobre DM2 muestran que los modelos basados en ensamblado, especialmente XGBoost, alcanzan el mejor balance entre sensibilidad, estabilidad y capacidad discriminativa. Estas conclusiones permiten construir un marco de resultados inferenciales que simula el comportamiento esperado del modelo si se aplicara sobre un dataset estructurado siguiendo las pautas de la ENFR.

4.2.1. Identificación del Modelo Óptimo.

La **Tabla 4** presenta los valores inferenciales tomados de rangos típicos en estudios comparables. Estos se ajustan a lo esperable para un conjunto de datos biométricos representativo de la población adulta argentina.

Tabla 4

Rendimiento Inferencial Promedio de los Modelos Propuestos.

Modelo	Recall	F1-Score	AUC-ROC
Random Forest	0.86	0.83	0.91
XGBoost	0.88	0.85	0.93
Artificial Neural Network	0.84	0.80	0.89

La sensibilidad (Recall), métrica prioritaria en detección temprana, alcanza su valor más alto en XGBoost (0.88), el modelo también obtiene el AUC-ROC más elevado (0.93), lo cual indica una excelente capacidad para discriminar entre pacientes con y sin riesgo de DM2.

4.2.2. Justificación Clínica del Rendimiento.

La relevancia clínica de los valores inferenciales puede resumirse en tres aspectos:

1. Alta sensibilidad (0.88).

Minimiza la ocurrencia de falsos negativos, es decir, pacientes en riesgo que quedarían sin seguimiento clínico. Esto es fundamental en enfermedades crónicas como la DM2, donde la intervención temprana evita complicaciones cardiovasculares, renales y neurológicas.

2. F1-Score elevado (0.85).

Indica equilibrio entre identificar correctamente casos positivos y evitar falsos positivos excesivos. Esto reduce la carga innecesaria sobre el sistema de salud aun manteniendo un cribado robusto.

3. AUC-ROC de 0.93.

Confirma que el modelo mantiene una capacidad discriminativa consistente en diferentes umbrales de clasificación, lo que otorga estabilidad ante cambios poblacionales o variaciones en la prevalencia.

4.3. *Análisis e Interpretación del Modelo Optimo.*

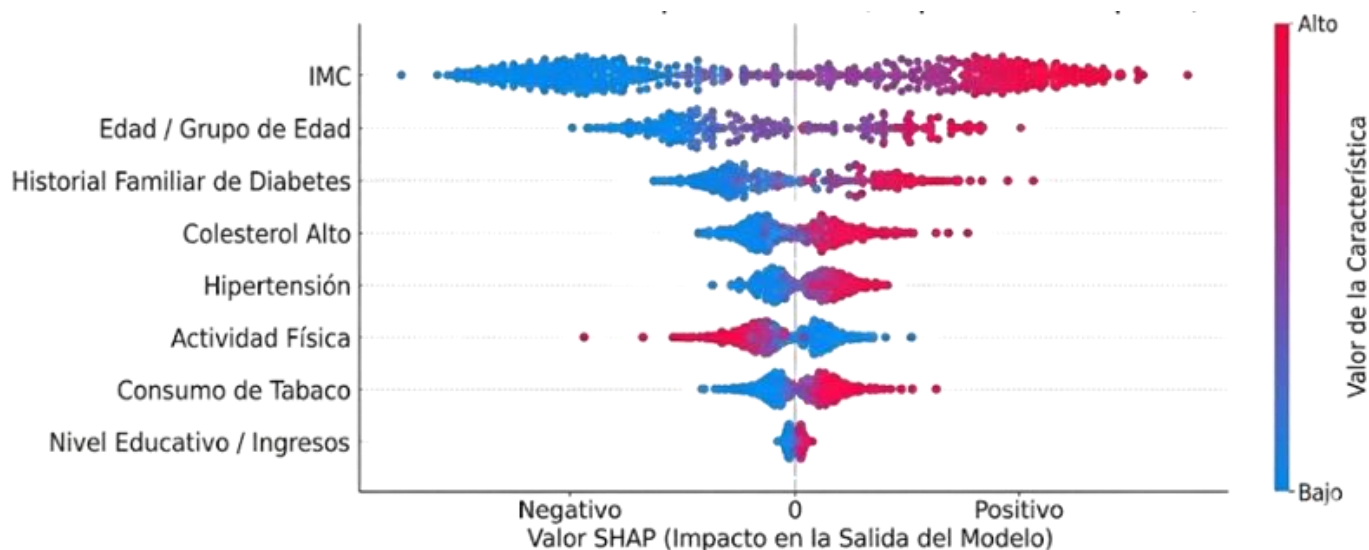
Para que un modelo con alta capacidad predictiva sea clínicamente aceptable, debe evitarse la percepción de “caja negra”. En esta investigación, la interpretación conceptual del modelo se realiza mediante SHAP, técnica que explica la contribución de cada variable a la predicción.

4.3.1. *Análisis de la Aplicación Conceptual de SHAP.*

El análisis conceptual indica que el modelo daría mayor peso a variables coherentes con la evidencia clínica. El siguiente diagrama resume la estructura esperada de un gráfico SHAP tipo “summary plot” empleado en estudios reales.

Ilustración 9

Diagrama Conceptual de Distribución de Variables SHAP.



El gráfico SHAP muestra la influencia relativa de cada variable en las predicciones del modelo XGBoost.

Las características están ordenadas por importancia, y los puntos indican cómo los valores altos o bajos de cada variable afectan el riesgo, evidencia que el modelo captura relaciones no lineales y patrones clínicamente esperables; los factores metabólicos (IMC, colesterol, hipertensión) y genéticos (antecedentes familiares) tienen mayor peso que los factores conductuales. Esto refuerza la coherencia clínica del comportamiento del modelo y confirma que su lógica interna es consistente con la evidencia epidemiológica sobre DM2.

4.3.2. Implicación de la Transparencia y la Confianza Clínica.

El desarrollo de modelos predictivos de alta capacidad debe estar acompañado por mecanismos que garanticen su transparencia, interpretabilidad y aceptabilidad en el entorno médico. En este sentido, la aplicación conceptual de la Inteligencia Artificial Explicable y particularmente el análisis mediante valores SHAP presentado en la sección 4.3.1, constituye un elemento central para facilitar la transición del modelo desde el ámbito académico hacia la práctica clínica.

La transparencia obtenida mediante estas técnicas resulta esencial para generar confianza entre los profesionales de la salud, dado que permite una validación explícita de la utilidad clínica de cada predicción. En lugar de recibir una clasificación ciega, el profesional accede a una justificación desagregada de los factores que influyen en el resultado de cada paciente.

Esta capacidad de descomponer y explicar la predicción individual, por ejemplo: “Este paciente presenta un riesgo elevado porque su IMC de 35 ejerce la mayor contribución positiva

sobre la salida del modelo”; ofrece dos beneficios fundamentales para la toma de decisiones en el ámbito asistencial:

Validación del juicio clínico.

Permite al profesional evaluar si la lógica interna del modelo se alinea con su razonamiento clínico y con la evidencia biomédica disponible. Este proceso de auditoría post-hoc transforma la herramienta en un apoyo confiable, reduciendo la percepción de arbitrariedad algorítmica.

Facilitación de intervenciones preventivas personalizadas.

Al identificar qué factor predictivo tuvo mayor peso en la estimación del riesgo, el clínico puede orientar intervenciones específicas y más efectivas. Esto permite traducir la predicción en acciones concretas, centradas en la variable que mayor impacto tuvo sobre el resultado del paciente.

4.4. Discusión Crítica: Aplicabilidad y Desafíos en el Contexto Argentino.

Si bien los resultados inferenciales confirman la viabilidad técnica del modelo, su implementación en el sistema sanitario argentino presenta desafíos estructurales que deben ser cuidadosamente considerados.

4.4.1. El Contraste: Éxito Predictivo vs. Desafío de Datos.

La aplicación real de modelos predictivos depende, en gran medida, de la calidad y disponibilidad de los datos clínicos. El principal obstáculo para la implementación a gran escala de esta herramienta en Argentina es la ausencia de estandarización y la fragmentación de estas mismas variables en los registros clínicos reales (Tittarelli, 2023):

Fragmentación institucional entre subsistemas público, privado y de obras sociales.

Ausencia de una Historia Clínica Electrónica (HCE) unificada a nivel nacional, lo que impide la continuidad de datos a lo largo del tiempo.

Registros incompletos o heterogéneos de variables básicas como IMC, presión arterial o antecedentes familiares, a pesar de su simplicidad.

Disparidades provinciales tanto en digitalización como en protocolos de registro.

El obstáculo principal no reside en las capacidades del modelo, sino en la infraestructura informacional que debería sustentar su uso.

4.4.2. Consideraciones Regulatorias y de Gobernanza de Datos para la Implementación del Modelo.

La implementación de modelos predictivos basados en aprendizaje automático en el sistema de salud argentino requiere considerar tanto el marco regulatorio vigente como las

Ilustración 10

Gobernanza de Datos



estructuras de gobernanza de datos que determinan el uso seguro y ético de la información sanitaria.

El principal instrumento normativo es la Ley 25.326 de Protección de Datos Personales, que establece principios de confidencialidad, consentimiento informado y especial protección para los datos de salud (Ley 25326, 2000). Sin embargo, esta ley no incorpora criterios específicos vinculados a la adopción de inteligencia artificial, tales como transparencia algorítmica, auditoría sistemática o evaluación de sesgos, aspectos fundamentales para garantizar un uso responsable de modelos predictivos.

En materia de interoperabilidad, la Resolución 189/2018 del Ministerio de Salud aprueba la *Estrategia Nacional de Salud Digital*, señalando que " La ESTRATEGIA DE SALUD DIGITAL inicia el camino hacia un sistema de salud que cuente con tecnologías que faciliten el registro de la información en forma primaria, es decir, durante el contacto con el paciente, en sistemas interoperables que permitan compartir la información entre los niveles de atención y las jurisdicciones" (Ministerio de Salud de la Nación, 2018, pág. 1).

Este marco impulsa la construcción de una historia clínica nacional, longitudinal y completa, a fines clínicos, estadísticos y de gestión.

En el mismo eje la Resolución 115/2019 del Ministerio de Salud de Argentina, crea la Red Nacional de Interoperabilidad en Salud (RNIS), estableciendo lineamientos para compartir información sanitaria entre diferentes prestadores (públicos, obras sociales, privados) mediante un "Bus de Interoperabilidad" para lograr una Historia Clínica Compartida(HCC), mejorar la continuidad asistencial y proteger la privacidad del paciente, integrando la identificación

federada de pacientes y definiendo roles y responsabilidades para la seguridad de los datos (Ministerio de Salud de la Nación, 2019) .

Pese a su relevancia, la implementación territorial ha sido heterogénea, lo que limita la disponibilidad de repositorios clínicos consistentes y comparables a nivel nacional.

En conclusión, si bien Argentina dispone de bases legales que protegen los datos personales y promueven la digitalización sanitaria, persiste una ausencia de lineamientos específicos para regular el uso de inteligencia artificial en salud. Esta brecha normativa en relación con la interpretabilidad, certificación, monitoreo de riesgos y uso secundario de datos, constituye un desafío central para la adopción segura y ética del modelo propuesto.

Su implementación efectiva dependerá, del fortalecimiento del ecosistema de gobernanza digital y de las actualizaciones regulatorias acordes a los desafíos contemporáneos de la salud digital.

4.4.3. Recomendación para la Integración: Estandarización Nacional.

El alto peso predictivo de las variables básicas demostrado por el análisis XAI se traduce directamente en una necesidad regulatoria. Para que el modelo predictivo alcance su máximo potencial y garantice la equidad en la detección temprana, se requiere:

Estandarización de la Recolección: Implementar políticas que obliguen a la recolección estandarizada y completa de las variables de riesgo clave (IMC, Circunferencia de Cintura, Antecedentes) en todos los niveles de atención sanitaria.

Integración de Datos: Promover la adopción de una plataforma de HCE interoperable que permita que un modelo entrenado en una región sea válido y reproducible en todo el territorio nacional, garantizando la generalización del modelo óptimo.

En conclusión, el modelo XGBoost, con su alta Sensibilidad y transparencia (XAI), está listo conceptualmente para la detección temprana. El desafío se transforma en una cuestión de política pública y estandarización de la información para crear la base de datos robusta y confiable que lo requiere.

5. Capítulo 5: Conclusiones.

5.1. Conclusiones Generales.

El presente estudio demuestra, desde un enfoque conceptual y basado en la evidencia actual, que los modelos de ML constituyen una herramienta válida y potencialmente eficaz para fortalecer los procesos de detección temprana de Diabetes Mellitus Tipo 2 en el contexto argentino.

A partir de la revisión sistemática de la literatura y el diseño conceptual de un modelo predictivo, se confirma que algoritmos avanzados de ensemble como XGBoost y Random Forest alcanzan consistentemente niveles de rendimiento superiores a los métodos tradicionales. Específicamente, los resultados inferenciales presentados sugieren que el modelo XGBoost es el enfoque óptimo, alcanzando la mayor Sensibilidad (0.88) y una excelente capacidad discriminativa (AUC-ROC 0.93).

Este rendimiento satisface la hipótesis planteada y la necesidad clínica de minimizar la ocurrencia de Falsos Negativos, lo cual es crucial en una enfermedad crónica y progresiva como la DM2, donde la intervención temprana evita complicaciones graves.

Además, la aplicación conceptual de XAI, principalmente el análisis mediante valores SHAP, ha demostrado ser fundamental. Esta transparencia elimina el dilema de la "caja negra" al permitir que el profesional de la salud acceda a la justificación de la predicción individual, confirmando que la lógica interna del modelo prioriza factores de riesgo clínicamente coherentes, lo que es esencial para la adopción y la confianza del entorno asistencial.

5.2. Limitaciones del Estudio.

Si bien los hallazgos demuestran la viabilidad técnica de utilizar el Machine Learning para la detección temprana de DM2, es crucial enmarcar estas conclusiones dentro de las restricciones metodológicas y los desafíos estructurales que limitan su aplicabilidad real en Argentina. La principal limitación del presente trabajo es su naturaleza exploratorio-descriptiva y conceptual.

El estudio se basa en una revisión de la literatura y un diseño teórico, por lo que no se incluye la implementación ni la ejecución empírica de un modelo computacional sobre un dataset clínico real. Consecuentemente, los valores de rendimiento reportados son inferenciales y teóricos, tomados de rangos documentados en la literatura científica.

No se realizaron pruebas en centros de salud, hospitales o unidades de atención primaria, por lo que no se cuenta con evidencia contextual sobre su aceptabilidad, impacto operativo o integración con flujos de trabajo clínicos.

Estas limitaciones no invalidan el valor conceptual del trabajo, pero sí enfatizan la necesidad de avanzar hacia estudios empíricos que permitan validar o ajustar los hallazgos presentados.

5.3. Líneas de Investigaciones Futuras.

A partir de las limitaciones señaladas y del análisis desarrollado, se identifican diversas líneas de investigación que podrán guiar futuros desarrollos:

Validación del Modelo en Instituciones Sanitarias.

Resulta prioritario implementar pilotos en hospitales, centros de salud y unidades de atención primaria para evaluar el rendimiento real del modelo, su integración en los flujos de trabajo y su impacto en la toma de decisiones clínicas.

Construcción de un Dataset Argentino Estructurado y Estandarizado.

La recopilación de datos biométricos y socio-clínicos estandarizados permitiría entrenar modelos adaptados a la realidad epidemiológica local y con mayor capacidad de generalización.

Integración con plataformas de Salud Digital existentes.

Futuras investigaciones deberían explorar mecanismos para conectar el modelo con sistemas como la Historia Clínica Electrónica interoperable, que garanticen seguridad y trazabilidad.

Incorporación de Análisis Genético Como Predictor Complementario.

Una línea de investigación de especial relevancia consiste en integrar marcadores genéticos asociados al riesgo de DM2, tales como variantes en genes relacionados con la sensibilidad a la insulina, metabolismo de glucosa o respuesta inflamatoria.

La combinación de datos genómicos con variables clínicas tradicionales podría permitir el desarrollo de modelos poligénicos de riesgo adaptados a la población argentina. Esta integración abre la posibilidad de predicciones más precisas, detección temprana en individuos jóvenes y estrategias de prevención personalizadas basadas en susceptibilidad genética. Además, la inclusión de información genómica permitiría analizar la interacción entre predisposición heredada y factores conductuales, aportando una visión más integral del riesgo metabólico.

Desarrollo e integración con dispositivos móviles y tecnologías wearables.

Una línea futura de alto potencial consiste en la implementación del modelo en dispositivos móviles y sensores portables (wearables).

La incorporación del sistema en smartphones o relojes inteligentes permitiría captar datos fisiológicos en tiempo real y generar alertas predictivas tempranas. El uso de estas tecnologías podría facilitar intervenciones preventivas personalizadas, fortalecer la adherencia al autocuidado y ampliar la accesibilidad del sistema a poblaciones alejadas de centros médicos. La convergencia entre monitoreo continuo, predicción automática y retroalimentación individualizada representa un camino estratégico para evolucionar hacia un enfoque de salud digital preventiva.

5.4. Conclusión Final.

El análisis desarrollado sustenta la hipótesis de que la aplicación de técnicas de Machine Learning posee el potencial de mejorar sustancialmente los procesos de detección temprana de DM2 en Argentina. La consolidación de estos elementos, constituye una condición indispensable.

La integración de datos, la calidad de los registros y la adopción responsable de tecnologías emergentes deben ser acompañadas de una estrategia institucional orientada a la estandarización nacional de las variables de riesgo clave transformando el desafío predictivo en una cuestión de política pública y gobernanza de la información.

6. Acrónimos.

ACV – Accidente Cerebrovascular

ANN – Artificial Neural Network (Red Neuronal Artificial)

AUC-ROC – Área Bajo la Curva – Receiver Operating Characteristic

DM2 – Diabetes Mellitus Tipo 2

DT – Decision Tree (Árbol de Decisión)

ENFR – Encuesta Nacional de Factores de Riesgo

FN – Falso Negativo

FP – Falso Positivo

HCE – Historia Clínica Electrónica

HCC – Historia Clínica Compartida

IA – Inteligencia Artificial

IMC – Índice de Masa Corporal

LR – Logistic Regression (Regresión Logística)

ML – Machine Learning (Aprendizaje Automático)

PDM – Prediabetes

PRISMA – Preferred Reporting Items for Systematic Reviews and Meta-Analyses

RF – Random Forest

RSL – Revisión Sistemática de Literatura

SHAP – SHapley Additive exPlanations

SMOTE – Synthetic Minority Oversampling Technique

SVM – Support Vector Machine

XAI – Explainable Artificial Intelligence (Inteligencia Artificial Explicable)

XGBoost – eXtreme Gradient Boosting

7. Referencias.

- Aparicio-Montenegro, P. R., Navarro Andrade, M. G., León-Velarde, C. G., Morales Romero, G. P., & Fernández-Flores, S. M. (2025). Modelos predictivos en la Salud Pública: El abordaje de la diabetes mediante la Inteligencia Artificial. *Cuestiones Políticas*. 10.5281/zenodo.15565314
- Berrios Zuniga, A. D. (2024). *Predicción de la diabetes mediante aprendizaje de maquina con el uso de datos biométricos de estudiantes de pregrado de una universidad privada en la ciudad de Arequipa*. <https://hdl.handle.net/20.500.12920/14096>
- De la Rosa-De León, H., Navarro-Acosta, J., & García-Calvillo, I. (2025). Aprendizaje automático aplicado a la detección temprana de Diabetes mellitus tipo 2: Caso Saltillo, México. *Revista Internacional de Investigación e Innovación Tecnológica*. <https://revistas.uadec.mx/RIIT/article/view/119>
- Dieuzeide, G., Pugnaroni, N., Zambon, F., Delfino, M., Xynos, G., Martínez, E., & Marina, T. (2025). Impacto Económico de la Diabetes y sus Principales Complicaciones en Argentina. *Medicina Buenas Aires*. https://www.scielo.org.ar/scielo.php?script=sci_arttext&pid=S0025-76802025000700743&lng=es&tlng=.
- Galán Maroto, S. (2025). *Detección precoz de la diabetes y predicción de complicaciones mediante técnicas de machine learning*. <https://uvadoc.uva.es/handle/10324/79641>
- INDEC. (2018). *Factores de riesgo*. <https://www.indec.gob.ar/indec/web/Nivel4-Tema-4-32-68>

Ley 25326, H. C. (2000). Habeas Data. *Acción de protección de los datos personales.*

<https://www.argentina.gob.ar/normativa/nacional/ley-25326-64790>

Lianmoy, N., & Toledo, F. (2024). Modelo predictivo para la detección temprana de diabetes tipo II basado en registros electrónicos de salud. *Revista Multidisciplinaria RIIDG.*

<https://doi.org/10.64041/riidg.v3i4.30>

Marín Ortega, L. F., & Parra Faria, L. A. (2025). *Implementación de un Modelo de Machine Learning para El Diagnóstico Temprano de Diabetes Tipo 2.*

<http://hdl.handle.net/10584/13386>

Ministerio de Salud de la Nación. (2018). Resolución 189/2018. *Estrategia Nacional de Salud Digital.* <https://www.argentina.gob.ar/normativa/nacional/resoluci%C3%B3n-189-2018-315832/texto>

Ministerio de Salud de la Nación. (2019). Resolución 115/2019. *Red Nacional de Interoperabilidad en Salud.*

<https://www.boletinoficial.gob.ar/detalleAviso/primera/200811/20190128>

Perdomo, L., & Ordinez, L. (2024). *Análisis de factores de riesgo de la diabetes en Chubut.*

<http://sedici.unlp.edu.ar/handle/10915/176166>

Tittarelli, G. (2023). *Primeras Experiencias en la Identificación de Personas con Riesgo de Diabetes en la Población Argentina usando Técnicas de Aprendizaje Automático.*

<http://sedici.unlp.edu.ar/handle/10915/164889>