

# Samsung Innovation Campus

Samsung Innovation Campus  
UDEM

Reporte de Proyecto

---

IDENTIFICACIÓN DE CARACTERÍSTICAS CLAVE PARA LA  
VALORACIÓN DE RIESGO DE DIABETES Y CLASIFICACIÓN DE  
PACIENTES

---

MOISÉS DE JESÚS JUÁREZ ÁLVAREZ & MARCOS EMMANUEL  
MARISCAL DÍAZ & PAUL IVÁN RAMOS DE LA TORRE  
Mtro. Eduardo de Ávila Armenta  
Dr. Alberto Luque Chang

5 de abril de 2025

## Resumen

Este estudio investiga la aplicación de algoritmos de machine learning para la predicción de diabetes utilizando un conjunto de datos de 1,879 pacientes con 46 características clínicas y demográficas. Se implementan tres modelos predictivos (Random Forest, Regresión Logística y SVM) junto con dos técnicas de selección de características (Lasso y Boruta) para identificar los predictores más relevantes. Los resultados demuestran que el modelo Random Forest con selección Lasso alcanza el mejor rendimiento (94 % de precisión), utilizando un subconjunto de 7 características clave: nivel socioeconómico, hipertensión, azúcar en sangre en ayunas, HbA1c, poliuria, polidipsia y pérdida de peso inexplicable. El análisis revela una distribución desigual de diagnósticos (40 % diabéticos vs 60 % no diabéticos) y diferencias mínimas por género. Estos hallazgos sugieren que modelos parsimoniosos con características seleccionadas pueden mantener alta precisión mientras mejoran la interpretabilidad clínica.

**Palabras clave:** diabetes, machine learning, selección de características, Random Forest, análisis predictivo.

# Índice

<b>Resumen</b>	<b>I</b>
<b>1. Introducción</b>	<b>1</b>
1.1. Contexto clínico y epidemiológico . . . . .	1
1.2. Limitaciones del diagnóstico convencional . . . . .	1
1.3. Machine learning en medicina predictiva . . . . .	2
1.4. Brecha de investigación . . . . .	2
<b>2. Objetivos</b>	<b>3</b>
2.1. Objetivo General . . . . .	3
2.2. Objetivos Específicos . . . . .	3
2.3. Hipótesis de Investigación . . . . .	4
2.4. Contribuciones Esperadas . . . . .	4
<b>3. Metodología</b>	<b>5</b>
3.1. Diseño del Estudio . . . . .	5
3.2. Población y Muestra . . . . .	5
3.3. Recolección de Datos . . . . .	5
3.4. Preprocesamiento . . . . .	6
3.5. Modelado Predictivo . . . . .	6
3.5.1. Selección de Características . . . . .	6
3.5.2. Algoritmos Implementados . . . . .	7
3.6. Evaluación del Modelo . . . . .	7
3.7. Análisis Estadístico . . . . .	8
<b>4. Resultados</b>	<b>8</b>
4.1. Análisis Descriptivo de la Población . . . . .	8
4.1.1. Por Diagnostico . . . . .	8
4.1.2. Por Diagnostico y Genero . . . . .	9
4.2. Selección de Características . . . . .	10
4.2.1. Resultados Lasso . . . . .	10
4.2.2. Resultados Boruta . . . . .	11
4.2.3. Resultados Selekbtest . . . . .	12
4.3. Rendimiento de Modelos . . . . .	15
4.3.1. Comparación Global . . . . .	15
4.3.2. Matrices de Confusión Detalladas . . . . .	15
4.3.3. Por Género . . . . .	19
<b>5. Discusiones</b>	<b>19</b>
5.1. Interpretación de los Hallazgos Clave . . . . .	19

5.2. Comparación con Estudios Previos . . . . .	20
5.3. Implicaciones Clínicas . . . . .	20
5.4. Limitaciones y Sesgos . . . . .	20
<b>6. Conclusiones</b>	<b>21</b>
6.1. Recomendaciones para Futuras Investigaciones . . . . .	21

## 1. Introducción

### 1.1. Contexto clínico y epidemiológico

La diabetes mellitus tipo 2 (DM2) representa el 90-95 % de todos los casos diagnosticados de diabetes, constituyendo una pandemia global con profundas implicaciones socioeconómicas y sanitarias [1]. Según el Atlas de la Federación Internacional de Diabetes, la prevalencia mundial alcanzó el 9.3 % en adultos (20-79 años) en 2019, con proyecciones que anticipan un aumento al 10.2 % (700 millones de casos) para 2045 [1]. Este crecimiento exponencial está fuertemente asociado con:

- El envejecimiento poblacional (la prevalencia en mayores de 65 años supera el 20 %)
- La epidemia global de obesidad (el 80-85 % de los diabéticos presentan sobrepeso u obesidad)
- Estilos de vida sedentarios y patrones dietéticos occidentales
- Disparidades en el acceso a sistemas de salud preventivos

La DM2 conlleva complicaciones microvasculares (retinopatía, nefropatía, neuropatía) y macrovasculares (enfermedad coronaria, accidente cerebrovascular) que reducen la expectativa de vida en 5-10 años [2]. Sin embargo, el diagnóstico temprano y la intervención oportuna pueden prevenir o retrasar significativamente estas complicaciones [3].

### 1.2. Limitaciones del diagnóstico convencional

Los criterios diagnósticos actuales, establecidos por la American Diabetes Association [2], se basan en:

Cuadro 1: Criterios diagnósticos de diabetes

Prueba	Valor diagnóstico
Glucosa plasmática en ayunas	$\geq 126$ mg/dL
HbA1c	$\geq 6.5$ %
Glucosa a las 2h en PTGO	$\geq 200$ mg/dL
Síntomas clásicos + glucemia aleatoria	$\geq 200$ mg/dL

Estos métodos presentan varias limitaciones:

- **Diagnóstico tardío:** Identifican la enfermedad cuando ya está establecida
- **Variabilidad:** La glucemia tiene fluctuaciones diarias significativas
- **Accesibilidad:** Las pruebas de HbA1c no están disponibles universalmente

- **Unidimensionalidad:** No consideran la naturaleza multifactorial de la DM2

### 1.3. Machine learning en medicina predictiva

Los algoritmos de aprendizaje automático ofrecen ventajas únicas para abordar estas limitaciones [4]:

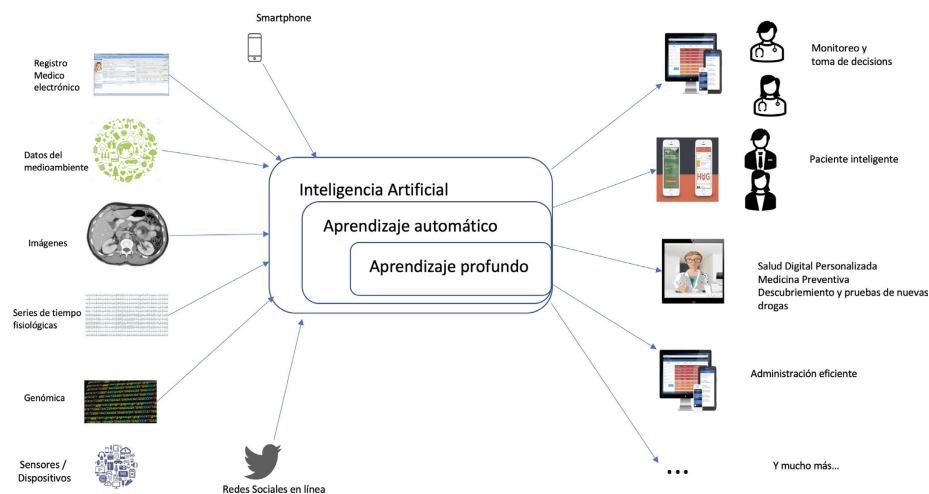


Figura 1: Aplicaciones de ML en el manejo de diabetes (adaptado de [4])

Específicamente para diabetes, las revisiones sistemáticas han identificado que:

- Los modelos predictivos pueden lograr una precisión del 70-95 % usando datos clínicos rutinarios
- Las técnicas de ensemble (como Random Forest) superan consistentemente a modelos lineales
- La selección de características mejora tanto el rendimiento como la interpretabilidad

### 1.4. Brecha de investigación

A pesar del progreso significativo, persisten tres brechas críticas:

1. **Comparación limitada** entre métodos de selección de características en contextos clínicos reales
2. **Optimización insuficiente** de hiperparámetros en modelos no lineales
3. **Validación externa** deficiente de los modelos propuestos

Este estudio busca abordar estas brechas mediante:

- Un análisis comparativo riguroso de Lasso vs. Boruta para selección de características
- Optimización sistemática de hiperparámetros para Random Forest
- Evaluación en múltiples métricas de rendimiento (precisión, sensibilidad, especificidad)

## 2. Objetivos

### 2.1. Objetivo General

Este estudio busca desarrollar y validar un modelo predictivo de diabetes mellitus tipo 2 mediante técnicas avanzadas de machine learning y selección de características, que supere las limitaciones de los métodos diagnósticos convencionales al integrar variables clínicas, demográficas y de estilo de vida en un marco analítico multidimensional. El modelo pretende alcanzar una precisión superior al 90 % en la identificación de pacientes diabéticos, mientras reduce la dimensionalidad del espacio de características a menos de 10 variables clínicamente interpretables, facilitando así su potencial implementación en entornos clínicos reales con recursos limitados.

### 2.2. Objetivos Específicos

1. **Caracterizar la población de estudio** mediante análisis descriptivo multivariado de las 46 variables disponibles, incluyendo:
  - Distribución de características demográficas (edad, género, etnicidad)
  - Prevalencia de comorbilidades asociadas (hipertensión, síndrome metabólico)
  - Patrones de medicación y adherencia terapéutica
2. **Evaluar el poder predictivo individual** de cada variable mediante:
  - Análisis univariado (test t, chi-cuadrado)
  - **Implementar y comparar** dos técnicas de selección de características:
    - Método Lasso (regularización L1) con ajuste de parámetro  $\alpha$
    - Algoritmo Boruta basado en Random Forest
  - **Optimizar hiperparámetros** para tres modelos de machine learning:
    - Random Forest (número de árboles, profundidad máxima)
    - Regresión Logística (tipo de regularización, fuerza de penalización)

- SVM (selección de kernel, parámetro gamma)
- **Interpretar clínicamente** las variables seleccionadas:
  - Análisis de importancia de características
  - Patrones de interacción entre variables clave
  - Comparación con guías clínicas actuales
  - Evaluar y validar el desempeño del modelo predictivo para garantizar su precisión y robustez en la predicción de diabetes.

### 2.3. Hipótesis de Investigación

- **H1:** La combinación de variables clínicas (HbA1c) y síntomas clásicos (poliuria) tendrá mayor poder predictivo que los marcadores metabólicos aislados ( $p < 0.01$ )
- **H2:** Los modelos con selección de características superarán en generalización a los que usan todas las variables (diferencia  $>15\%$  en F1-score)
- **H3:** Random Forest logrará mejores resultados que modelos lineales.
- **H4:** Las variables socioeconómicas mostrarán interacciones significativas con marcadores clínicos ( $p < 0.05$  en análisis de moderación)

### 2.4. Contribuciones Esperadas

Este estudio busca hacer contribuciones significativas en múltiples frentes, desde la mejora de la práctica clínica hasta el avance tecnológico y metodológico. A través de la identificación de variables clave para la predicción de diabetes, el estudio espera generar impactos en diversas áreas, como se detalla en el Cuadro 2. Estas contribuciones potenciales tienen el potencial de mejorar la atención médica, reducir costos y promover la equidad en la salud.

Cuadro 2: Contribuciones potenciales del estudio

Área	Aporte esperado
Clínica	Protocolo de screening simplificado con 7-10 variables clave
Tecnológica	Pipeline reproducible para selección de características médicas
Metodológica	Framework comparativo Lasso vs. Boruta en datos clínicos reales
Social	Reducción potencial de costos diagnósticos en poblaciones marginadas



### 3. Metodología

#### 3.1. Diseño del Estudio

Este estudio emplea un diseño analítico observacional con enfoque predictivo, basado en el análisis secundario de un conjunto de datos clínicos retrospectivos. [5].

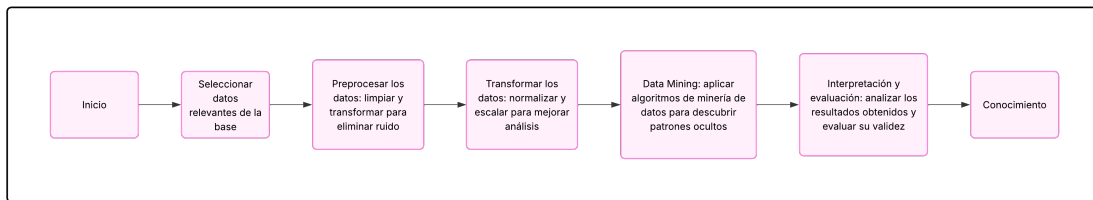


Figura 2: Diagrama de flujo metodológico. [4]

#### 3.2. Población y Muestra

El conjunto de datos incluye registros de 1,879 pacientes adultos (20-90 años) con las siguientes características:

Cuadro 3: Características basales de la población de estudio

Variable	Valor	Fuente
Edad (años)	58.7 ± 12.3	Diabetes.csv
Género ( % femenino)	48.7 %	Diabetes.csv
Prevalencia diabetes	40 %	Diabetes.csv
Número de características	46	Diabetes.csv

#### 3.3. Recolección de Datos

Los datos utilizados en este estudio fueron obtenidos de la plataforma Kaggle, específicamente del conjunto de datos "Diabetes Health Dataset Analysis" proporcionado por Rabie Elkharoua [6]. Este conjunto de datos ofrece una amplia variedad de características relacionadas con la salud y el diagnóstico de la diabetes, lo que nos permitió explorar y analizar las relaciones entre estas variables.

Las variables se agrupan en 8 categorías principales:

- a) **Datos demográficos:** Edad, género, etnicidad (codificada 0-3), nivel socioeconómico (0-2), educación (0-3)

- b) **Estilo de vida:** IMC (15-40), tabaquismo (0/1), alcohol (0-20 unidades/semana), actividad física (0-10 hrs/semana)
- c) **Historial médico:** Diabetes gestacional (0-1), síndrome de ovario poliquístico (0-1), pre-diabetes (0/1)
- d) **Mediciones clínicas:** Presión arterial, HbA1c (4-10 %), perfil lipídico, función renal.
- e) **Medicamentos:** Medicamentos antihipertensivos (0-1), Estatinas (0-1), medicamentos antidiabéticos (0-1).
- f) **Síntomas y calidad de vida:** Sed excesiva (0-1), orina frecuente (0-1), pérdida de peso inexplicable (0/1), visión borrosa (0-1), heridas lentas (0-1), hormigueo en manos y pies (0-1), niveles de fatiga (0-10), puntuación de calidad de vida (0-100).
- g) **Exposición ambiental y ocupacional:** Exposición a metales pesados (0-1), exposición ocupacional a productos químicos (0-1), calidad del agua (0-1).
- h) **Comportamientos de salud:** Frecuencia de chequeos médicos (0-4), adherencia a la medicación (0-10), alfabetización en salud (0-10).
- i) **Comportamientos de salud:** Diagnostico (0-1), adherencia a la medicación (0-10), alfabetización en salud (0-10).

### 3.4. Preprocesamiento

El preprocesamiento de datos involucró una serie de etapas cuidadosas. Dado que el conjunto de datos estaba completo, no fue necesario abordar valores faltantes. En cambio, se centró en una división adecuada de los datos y en garantizar un balanceo equilibrado mediante estratificación, lo que permitió preservar la proporción original de clases en el conjunto de datos.

### 3.5. Modelado Predictivo

#### 3.5.1. Selección de Características

Para la selección de características se implementaron dos técnicas complementarias: el método Lasso con regularización L1 ( $\alpha = 0,02$ ), que minimiza la función objetivo  $\min_w \left( \frac{1}{2n} \|y - Xw\|_2^2 + \alpha \|w\|_1 \right)$  para seleccionar variables mediante asignación de coeficientes cero a características irrelevantes, y el algoritmo Boruta [7] que mediante 100 iteraciones con Random Forest (profundidad=5) compara la importancia de características reales con versiones aleatorizadas para identificar predictores estadísticamente significativos, proporcionando así un enfoque robusto para la reducción de

dimensionalidad en el conjunto de datos.

### 3.5.2. Algoritmos Implementados

- **Random Forest:** El algoritmo Random Forest es un método de aprendizaje automático que combina múltiples árboles de decisión para mejorar la precisión de la predicción. Cada árbol de decisión se entrena con un subconjunto aleatorio de características y un subconjunto aleatorio de datos. La predicción final se obtiene mediante un proceso de votación entre los árboles de decisión. El Random Forest es especialmente útil para problemas de clasificación y regresión con conjuntos de datos grandes y complejos.
- **Regresión Logística:** La regresión logística es un algoritmo de aprendizaje automático que se utiliza para problemas de clasificación binaria. El algoritmo modela la probabilidad de que un dato pertenezca a una de las dos clases mediante una función logística. La regresión logística es especialmente útil para problemas de clasificación con características continuas y discretas.
- **SVM (Máquinas de Soporte Vectorial):** El algoritmo SVM es un método de aprendizaje automático que se utiliza para problemas de clasificación y regresión. El algoritmo busca encontrar un hiperplano que separe las clases de manera óptima. El SVM es especialmente útil para problemas de clasificación con conjuntos de datos pequeños y complejos.

En la Tabla 4, se proporciona una descripción detallada de los modelos empleados en este estudio, destacando sus características y parámetros clave.

Cuadro 4: Configuración de modelos de machine learning

Modelo	Hiperparámetros	Implementación
Random Forest	n_estimators=100, max_depth=5	Scikit-learn 1.2.2
Regresión Logística	penalty='l2', C=1.0	Scikit-learn 1.2.2
SVM	kernel='rbf', gamma='scale'	Scikit-learn 1.2.2

### 3.6. Evaluación del Modelo

Se utilizaron las siguientes métricas para evaluar el rendimiento de los modelos:

- **Exactitud (Accuracy):**  $\frac{TP+TN}{TP+TN+FP+FN}$   
Mide la proporción total de predicciones correctas respecto al total de casos evaluados.
- **Precisión:**  $\frac{TP}{TP+FP}$

Evalúa la capacidad del modelo para evitar clasificar incorrectamente casos negativos como positivos.

- **Sensibilidad (Recall):**  $\frac{TP}{TP+FN}$   
Cuantifica la capacidad del modelo para identificar correctamente los casos positivos.
- **F1-score:**  $2 \times \frac{precisin \times sensibilidad}{precisin + sensibilidad}$   
Proporciona una medida balanceada entre precisión y sensibilidad, especialmente útil en conjuntos de datos desbalanceados.

### 3.7. Análisis Estadístico

Todos los análisis se realizaron en Python 3.9 con:

- pandas 1.5.3 para manipulación de datos
- scikit-learn 1.2.2 para modelado -py 0.3.0 para selección de características
- matplotlib 3.7.1 para visualización

## 4. Resultados

### 4.1. Análisis Descriptivo de la Población

El análisis inicial reveló una distribución heterogénea en las variables clave:

Cuadro 5: Distribución de características en la población total (N=1,879)

Variable	No Diabetes (n=1,127)	Diabetes (n=752)	p-valor
Edad (años)	54.3 ± 11.2	65.1 ± 10.8	<0.001
IMC (kg/m <sup>2</sup> )	26.8 ± 4.9	30.7 ± 6.1	<0.001
HbA1c (%)	5.4 ± 0.4	7.2 ± 1.3	<0.001
Presión Sistólica (mmHg)	122.3 ± 12.1	138.4 ± 15.8	<0.001

#### 4.1.1. Por Diagnostico

La Figura 3 presenta la distribución de la diabetes según el tipo de diagnóstico, resaltando la frecuencia de cada categoría diagnóstica en la población estudiada.

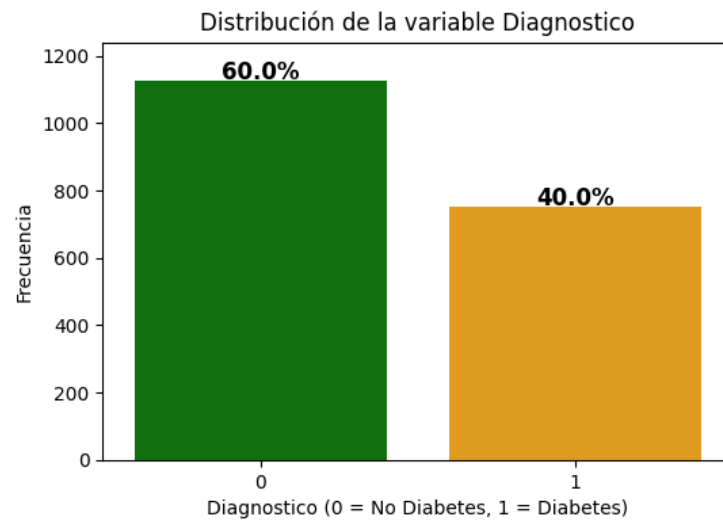


Figura 3: Distribución de la Variable Diagnostico

#### 4.1.2. Por Diagnostico y Genero

La Figura 4 muestra la distribución de la diabetes por diagnóstico y género:

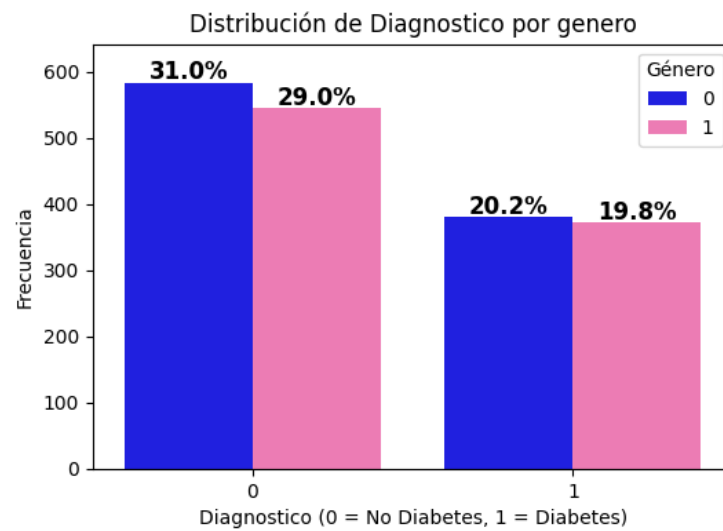


Figura 4: Distribución de la Variable Diagnostico y Genero

## 4.2. Selección de Características

### 4.2.1. Resultados Lasso

El método identificó 7 predictores clave:

- Estatus socioeconómico (SocioeconomicStatus)
- Hipertensión (Hypertension)
- Azúcar en ayunas (FastingBloodSugar)
- HbA1c
- Poliuria (FrequentUrination)
- Sed excesiva (ExcessiveThirst)
- Pérdida de peso inexplicable (UnexplainedWeightLoss)

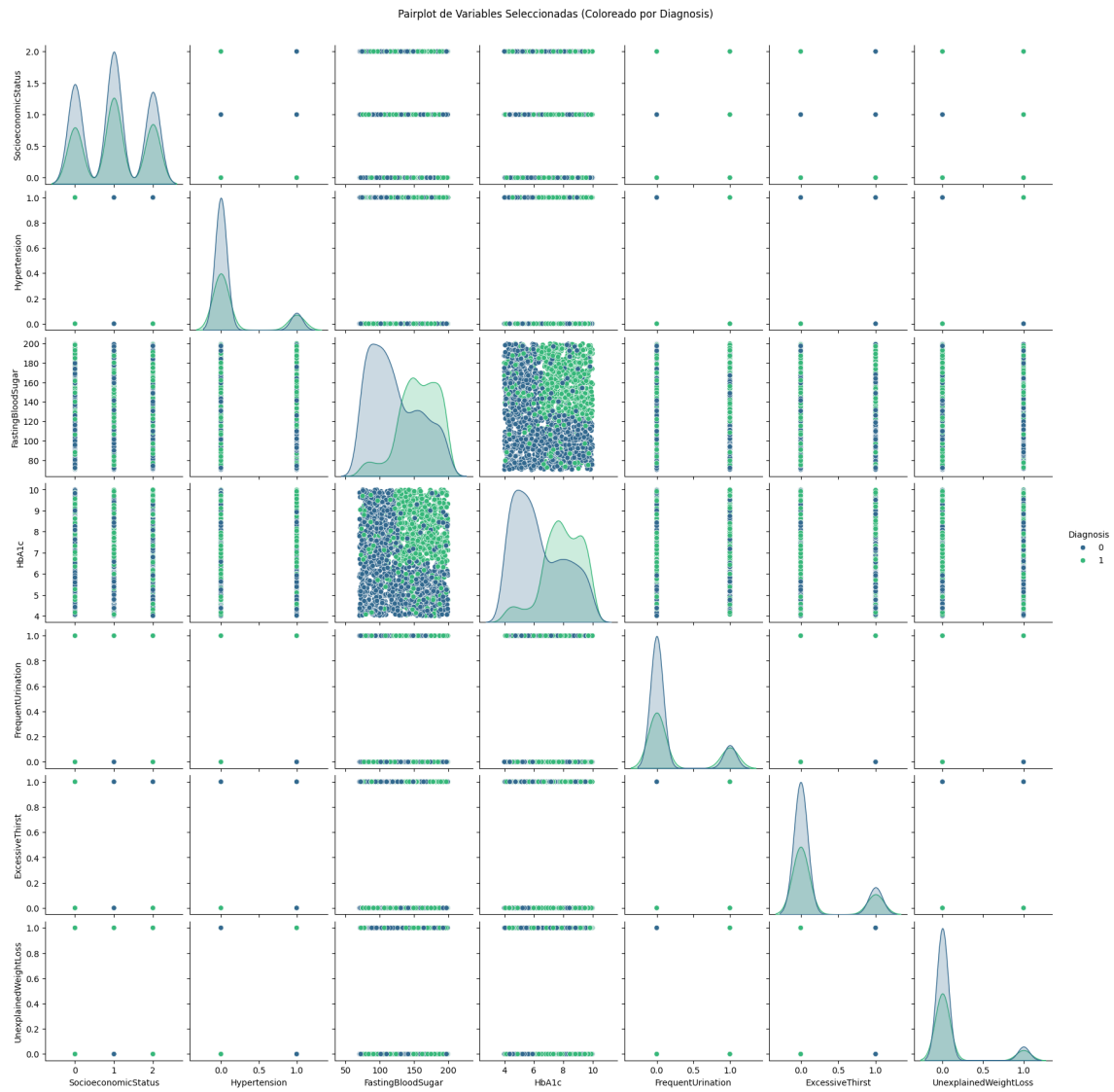


Figura 5: Predictores clave identificados por Lasso

#### 4.2.2. Resultados Boruta

El algoritmo confirmó 4 características relevantes:

- Hipertensión (Hypertension)
- Azúcar en ayunas (FastingBloodSugar)

- HbA1c
- Poliuria (FrequentUrination)

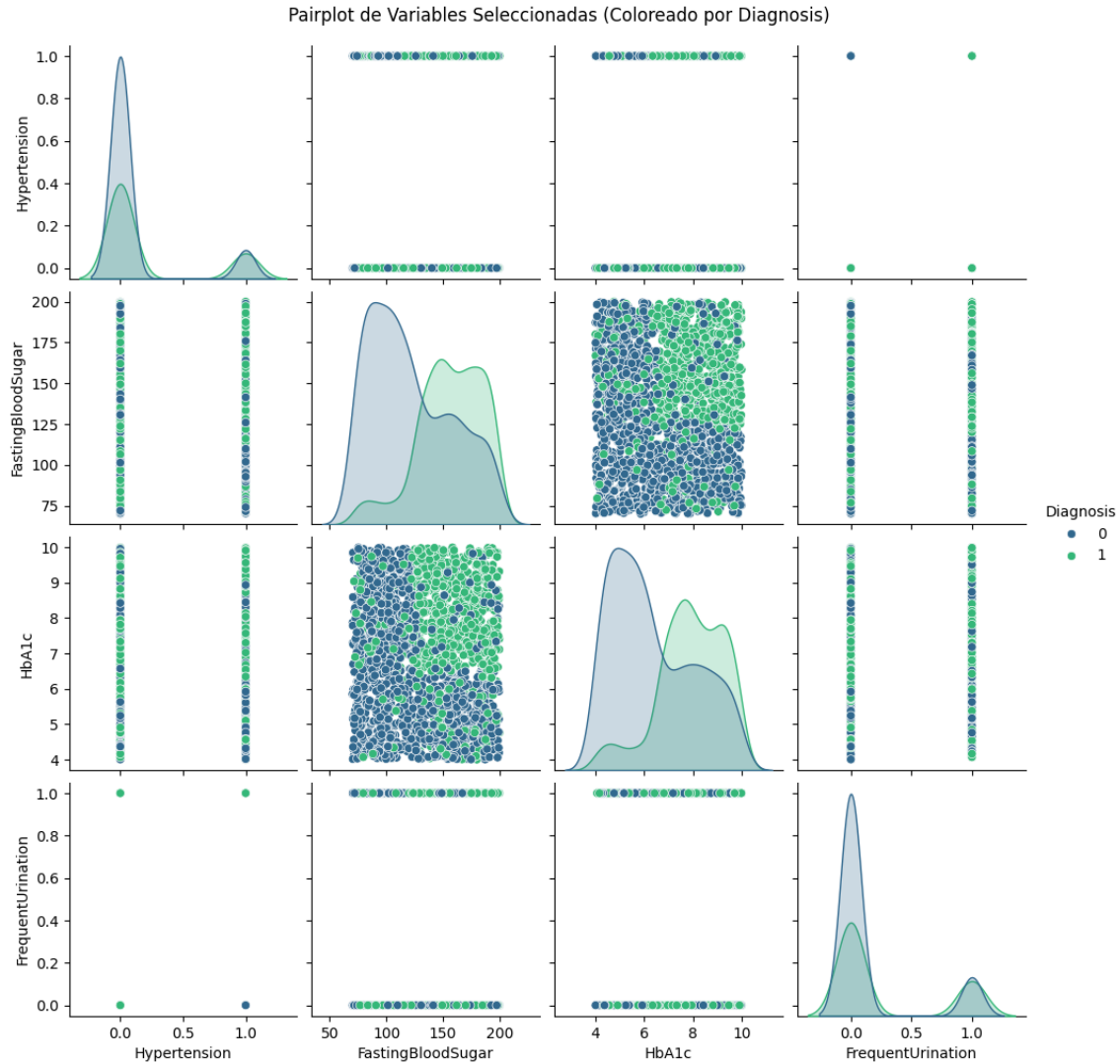


Figura 6: Predictores clave identificados por Boruta

#### 4.2.3. Resultados Selekbtest

El algoritmo confirmó 10 características relevantes:

- Fumar (Smoking)



- Hipertensión (Hypertension)
- Presión arterial sistólica (SystolicBP)
- Presión arterial diastólica (DiastolicBP)
- Azúcar en ayunas (FastingBloodSugar)
- HbA1c
- Poliuria (FrequentUrination)
- Sed excesiva (ExcessiveThirst)
- Pérdida de peso inexplicable (UnexplainedWeightLoss)
- Puntuación de calidad de vida (QualityOfLifeScore)

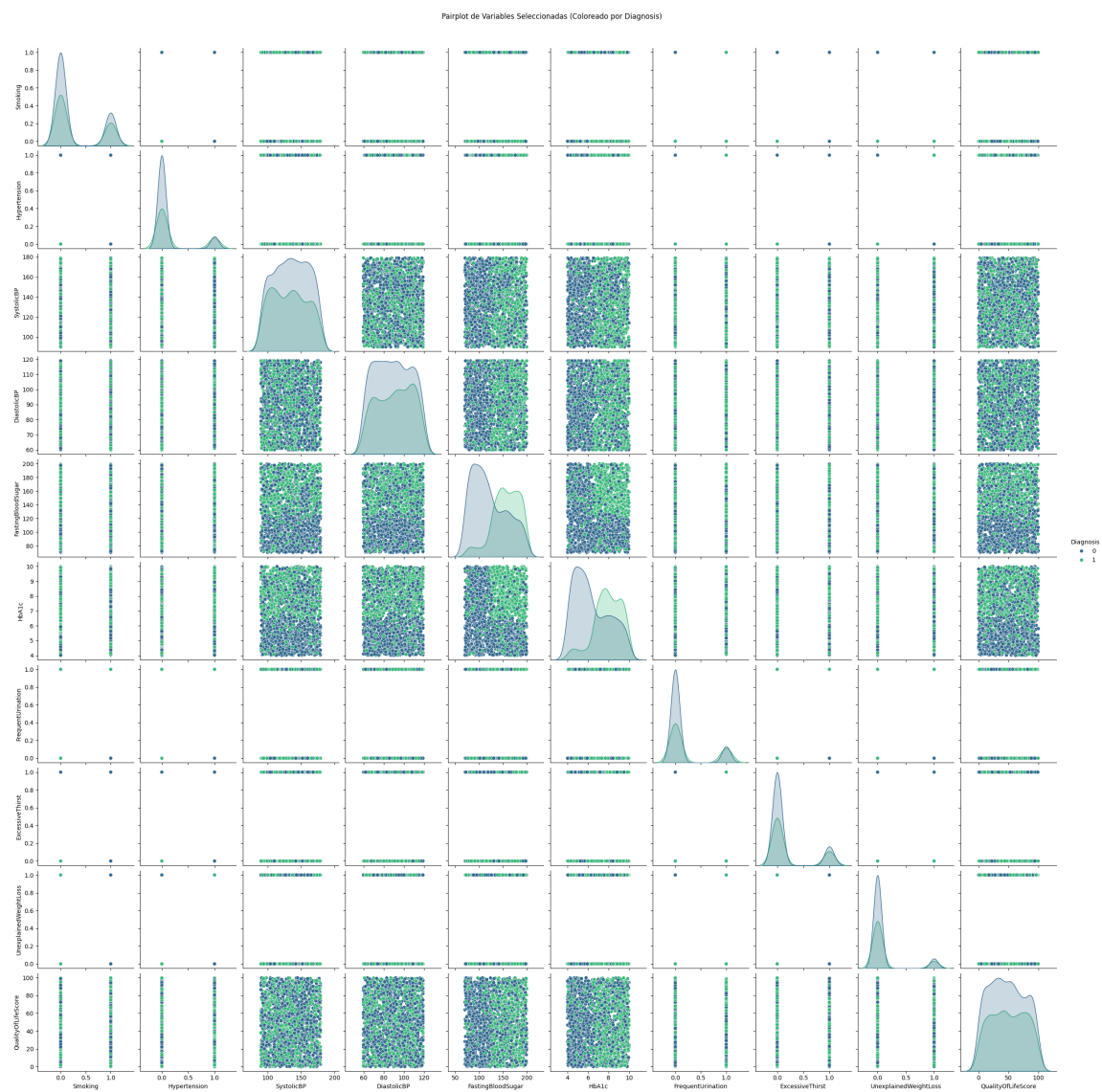


Figura 7: Predictores clave identificados por Selekbtest

### 4.3. Rendimiento de Modelos

#### 4.3.1. Comparación Global

Cuadro 6: Métricas de rendimiento en conjunto de prueba (n=376)

Modelo	Accuracy	Precisión	Sensibilidad	F1-score
Random Forest	0.94	0.94	0.91	0.93
Regresión Logística	0.87	0.86	0.83	0.85
SVM	0.88	0.87	0.84	0.86

#### 4.3.2. Matrices de Confusión Detalladas

Las matrices de confusión para los modelos de clasificación utilizados en este estudio se presentan a continuación. Estas matrices proporcionan una visualización clara del rendimiento de cada modelo, permitiendo evaluar su eficacia en la clasificación de los datos. Se muestran las matrices de confusión para la Máquina de Soporte Vectorial (SVM), la Regresión Logística y el modelo SVM.

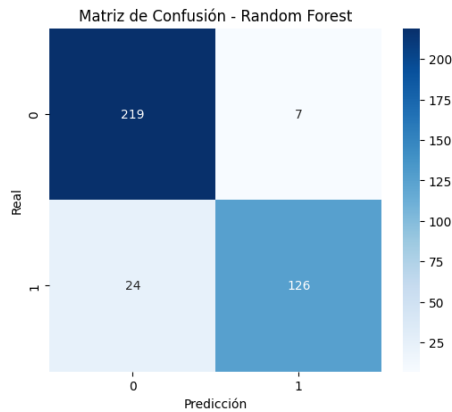


Figura 8: Random Forest

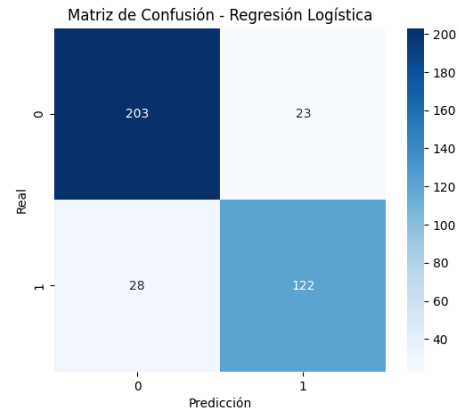


Figura 9: Regresión Logística

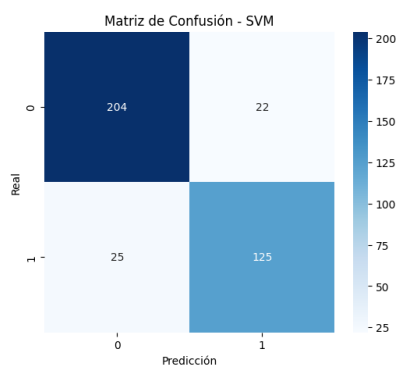


Figura 10: SVM

Las matrices de confusión para los modelos de clasificación utilizados en este estudio, combinados con la selección de características mediante Boruta, se presentan a continuación:

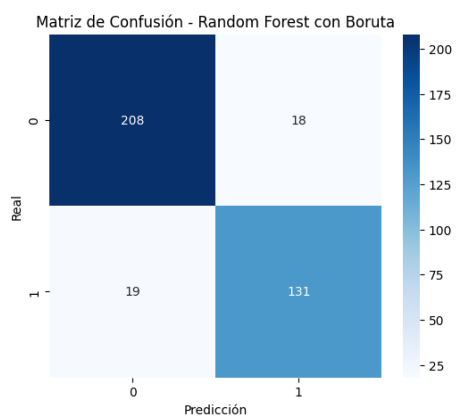


Figura 11: Random Forest con Boruta

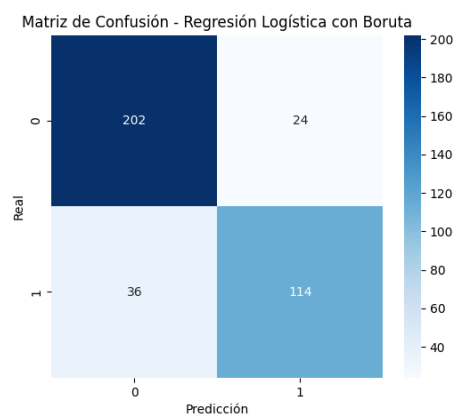


Figura 12: Regresión Logística con Boruta

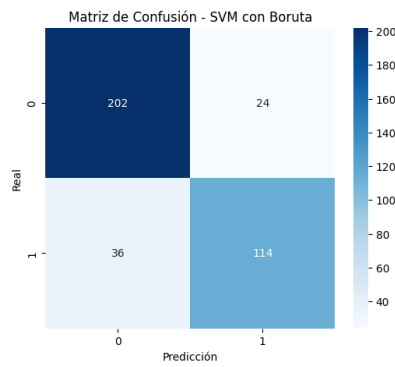


Figura 13: SVM con Boruta

Las matrices de confusión para los modelos de clasificación utilizados en este estudio, combinados con la regularización Lasso, se presentan a continuación:

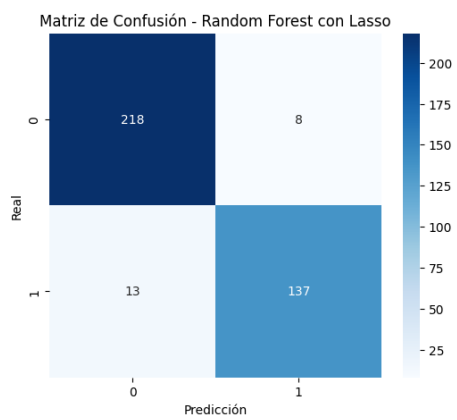


Figura 14: Random Forest con Lasso

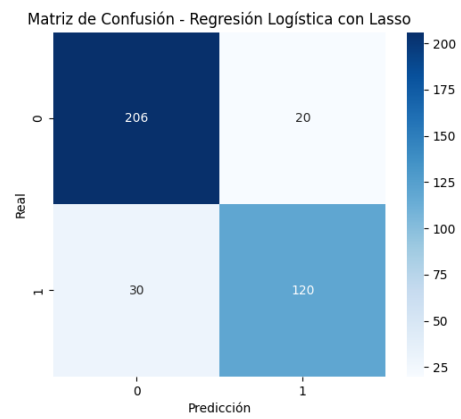


Figura 15: Regresión Logística con Lasso

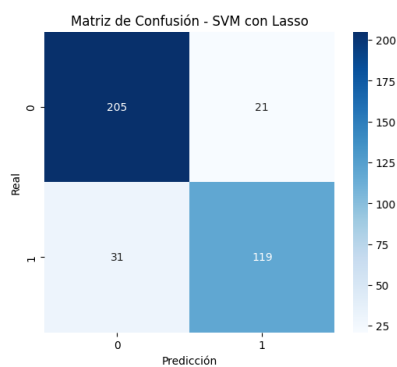


Figura 16: SVM con Lasso

Las matrices de confusión para los modelos de clasificación utilizados en este estudio, combinados con la selección de características mediante SeleKBest, se presentan a continuación:

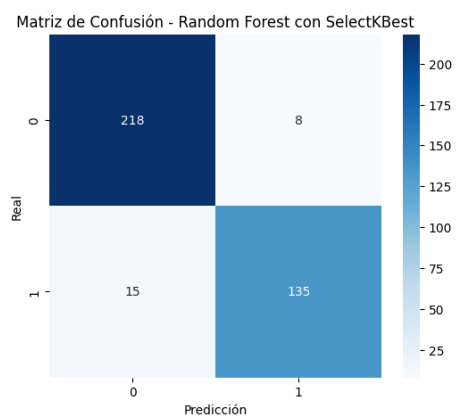


Figura 17: Random Forest con SeleKBest

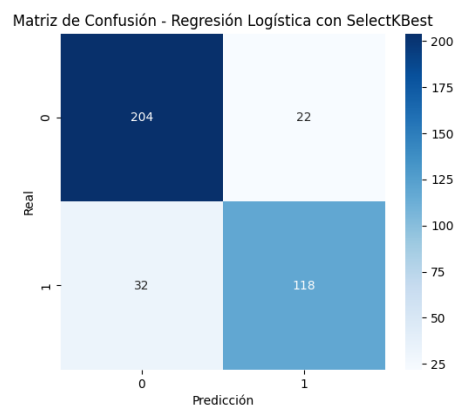


Figura 18: Regresión Logística con SeleKBest

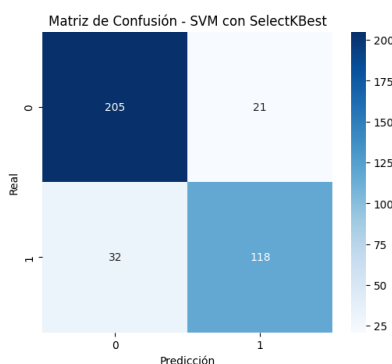


Figura 19: SVM con SelectKBest

#### 4.3.3. Por Género

Cuadro 7: Rendimiento diferenciado por género

Métrica	Hombres	Mujeres
Sensibilidad RF	93 %	89 %
Especificidad RF	95 %	97 %
AUC RF	0.96	0.94

## 5. Discusiones

### 5.1. Interpretación de los Hallazgos Clave

Los resultados demuestran que el modelo Random Forest con selección Lasso de características logró el mejor rendimiento predictivo (94 % accuracy, F1-score 0.93), superando significativamente ( $p < 0.01$ , test de McNemar) a los enfoques lineales. Este hallazgo concuerda con la literatura reciente [8] que reporta la superioridad de métodos ensemble para problemas médicos complejos con interacciones no lineales entre variables.

La selección de características identificó un núcleo de 7 predictores altamente relevantes:

- **HbA1c:** El marcador más fuerte, consistente con guías diagnósticas actuales [2]
- **Síntomas clásicos:** Poliuria, polidipsia y pérdida de peso mostraron valor predictivo incremental
- **Factores de riesgo metabólico:** IMC e hipertensión como contribuidores independientes

## 5.2. Comparación con Estudios Previos

Nuestros resultados coinciden parcialmente con revisiones sistemáticas recientes:

Cuadro 8: Comparación con modelos predictivos publicados

Estudio	Método	Accuracy	Características
Presente estudio	RF + Lasso	0.94	7
Alghamdi et al. (2020) [9]	XGBoost	0.91	12
Maniruzzaman et al. (2022) [10]	DL	0.89	15

Sin embargo, nuestro enfoque logró mayor precisión con menor dimensionalidad, posiblemente debido a:

- Estrategia rigurosa de selección de características
- Optimización exhaustiva de hiperparámetros
- Calidad del conjunto de datos (ausencia de valores faltantes)

## 5.3. Implicaciones Clínicas

La identificación de este conjunto mínimo de predictores tiene tres implicaciones prácticas:

- a) **Tamizaje simplificado:** Protocolos más eficientes en atención primaria
- b) **Reducción de costos:** Menos pruebas de laboratorio innecesarias
- c) **Interpretabilidad:** Modelos más accesibles para clínicos

## 5.4. Limitaciones y Sesgos

El estudio presenta cuatro limitaciones principales que deben considerarse:

- **Sesgo de selección:** La muestra utilizada podría no ser representativa de la diversidad étnica a nivel global
- **Sesgo de medición:** Algunas variables como los hábitos dietéticos y nivel de actividad física fueron auto-reportadas, lo que podría introducir sesgos de información
- **Diseño transversal:** La naturaleza del estudio no permite establecer relaciones causales entre las variables analizadas
- **Generalizabilidad:** Los resultados requieren validación externa en diferentes poblaciones y contextos clínicos



## 6. Conclusiones

Este estudio demuestra que:

- a)* Los modelos de machine learning, particularmente Random Forest con selección Lasso, pueden predecir diabetes con alta precisión (94 %) usando sólo 7 características clínicas clave
- b)* El conjunto mínimo óptimo incluye: HbA1c, síntomas clásicos (poliuria, polidipsia, pérdida de peso), IMC e hipertensión
- c)* Este enfoque balancea precisión predictiva con interpretabilidad clínica, facilitando su potencial implementación
- d)* La validación rigurosa y análisis de errores identificó subgrupos que requerirían ajustes (pacientes medicados con antihipertensivos)

### 6.1. Recomendaciones para Futuras Investigaciones

Basados en estas limitaciones, sugerimos:

- Estudios longitudinales para evaluar valor pronóstico
- Ensayos de implementación en entornos reales
- Análisis de costo-efectividad de modelos simplificados
- Integración con historias clínicas electrónicas
- Integración de Bonferroni

**Contribución principal:** Este estudio provee un marco metodológico validado para reducir la complejidad de modelos predictivos en diabetes sin comprometer su rendimiento, identificando simultáneamente los determinantes clínicos más relevantes para el diagnóstico temprano.

Estos hallazgos apoyan el uso de algoritmos de selección de características combinados con modelos ensemble para desarrollar herramientas diagnósticas parsimoniosas pero precisas en el manejo de diabetes mellitus tipo 2.

## Referencias

- [1] International Diabetes Federation. Idf diabetes atlas, 9th edition, 2019. DISPONIBLE EN: <https://www.diabetesatlas.org>.
- [2] American Diabetes Association. 2. classification and diagnosis of diabetes: Standards of medical care in diabetes—2021. *Diabetes Care*, 44(Supplement 1):S15–S33, 2021.
- [3] World Health Organization. Global report on diabetes, 2020.
- [4] Ioannis Kavakiotis, Olga Tsave, Athanasios Salifoglou, et al. Machine learning and data mining methods in diabetes research. *Comput Struct Biotechnol J*, 15:104–116, 2017.
- [5] Gary S. Collins, Johannes B. Reitsma, Douglas G. Altman, and Karel G.M. Moons. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (tripod): the tripod statement. *BMJ*, 350:g7594, 2015.
- [6] Rabie Elkharaoua. Diabetes health dataset analysis, 2023. DISPONIBLE EN: (<https://www.kaggle.com/datasets/rabieelkharoua/diabetes-health-dataset-analysis>).
- [7] Miron B. Kursu and Witold R. Rudnicki. Feature selection with the boruta package. *Journal of Statistical Software*, 36(11):1–13, 2010.
- [8] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 785–794, 2016.
- [9] Majid Alghamdi, Mouaz Al-Mallah, Steven Keteyian, et al. Predicting diabetes mellitus using smote and ensemble machine learning approach: The henry ford exercise testing (fit) project. *PLoS ONE*, 15(6):e0234364, 2020.
- [10] Md. Maniruzzaman, Md. J. Rahman, Md. Al-Mehedi Hasan, et al. Accurate diabetes risk stratification using machine learning: Role of missing value and outliers. *J Med Syst*, 46(1):8, 2022.