



Recibido: 4 de febrero 2026 / Aceptado: 18 de mayo 2026

SELECCIÓN DE ATRIBUTOS MEDIANTE FEATURE IMPORTANCE PARA LA CLASIFICACIÓN DE CASOS DE DENGUE EN MÉXICO

SELECTION OF ATTRIBUTES USING FEATURE IMPORTANCE FOR THE CLASSIFICATION OF DENGUE CASES IN MEXICO

Miguel Alcaraz Vázquez¹, Marco Antonio Adame Rodríguez², Dan Salvador García Guevara³, Edgardo Tomas Martínez⁴, Gustavo Adolfo Alonso Silverio⁵

Resumen:

El dengue es un importante reto para la salud pública en México, con una incidencia creciente y recursos diagnósticos limitados en las regiones endémicas. Este estudio propone el uso de técnicas de aprendizaje automático combinadas con un algoritmo personalizado de importancia de características para mejorar la clasificación de los casos de dengue usando variables con el fin de apoyar la vigilancia epidemiológica. El dataset proviene de los datos del conjunto de datos «Enfermedades Transmitidas por Vector» del Ministerio de Salud de México (febrero de 2024-febrero de 2025) y se definió la variable DICTAMEN (casos confirmados frente a casos negativos) como objetivo de clasificación. El método de importancia de las características basado en árboles de decisión redujo el conjunto de datos de 22 a 7 atributos clave, eliminando las variables redundantes y menos informativas. Se probaron algoritmos más usados en la literatura (Random Forest, Naive Bayes, MLP, entre otros) tanto en el conjunto de datos completo como en el reducido. Los resultados mostraron mejoras en la precisión y el equilibrio, especialmente en el caso de MLP y Naive Bayes. La vigilancia epidemiológica puede implementarse en dispositivos móviles, lo que permite un uso más amplio en sistemas de salud con recursos limitados.

Palabras claves: Dengue, Selección de características, Aprendizaje automático, Epidemiología, Salud pública

Abstract

Dengue fever is a major public health challenge in Mexico, with increasing incidence and limited diagnostic resources in endemic regions. This study proposes the use of machine learning techniques combined with a customized feature importance algorithm to improve the classification of dengue cases using variables to support epidemiological surveillance. The dataset comes from the “Vector-Borne Diseases” dataset of the Mexican Ministry of Health (February 2024-February 2025), and the DICTAMEN variable (confirmed cases versus negative cases) was defined as the classification target. The feature importance method based on decision trees reduced the dataset from 22 to 7 key attributes, eliminating redundant and less informative variables. The most commonly used algorithms in the literature (Random Forest, Naive Bayes, MLP, among others) were tested on both the complete and reduced datasets. The results showed improvements in accuracy and balance, especially in the case of MLP and Naive Bayes. Epidemiological surveillance can be implemented on mobile devices, allowing for wider use in health systems with limited resources.

Keywords: Dengue, Feature selection, Machine learning, Epidemiology, Public health

1 Licenciatura en Ingeniería en Computación, Universidad Autónoma de Guerrero, México, Chilpancingo de los Bravo. Afiliación institucional: Universidad Autónoma de Guerrero, México. Correo electrónico personal e institucional e-mail: 19275295@uagro.mx, vazquezmicky@gmail.com ORCID: <https://orcid.org/0009-0006-8355-0287>

2 Licenciatura en Ingeniería en Computación, Universidad Autónoma de Guerrero, México, Chilpancingo de los Bravo. Afiliación institucional: Universidad Autónoma de Guerrero, México. Correo electrónico personal e institucional e-mail: 18451084@uagro.mx, ragnarockgames86@gmail.com ORCID: <https://orcid.org/0009-0003-2361-4945>

3 Licenciatura en Ingeniería en Computación, Universidad Autónoma de Guerrero, México, Chilpancingo de los Bravo. Afiliación institucional: Universidad Autónoma de Guerrero, México. Correo electrónico personal e institucional e-mail: 16282142@uagro.mx, danlachele@gmail.com ORCID: <https://orcid.org/0009-0000-3751-4804>

4 Licenciatura en Ingeniería en Computación, Universidad Autónoma de Guerrero, México, Chilpancingo de los Bravo. Afiliación institucional: Universidad Autónoma de Guerrero, México. Correo electrónico personal e institucional e-mail: 24600063@uagro.mx, edgardo_tm@hotmail.com ORCID: <https://orcid.org/0009-0004-9689-8392>

5 Doctorado en Ciencias con en Ingeniería Eléctrica por parte del CINVESTAV-IPN, Universidad Autónoma de Guerrero, México, Chilpancingo de los Bravo. Afiliación institucional: Universidad Autónoma de Guerrero, México. Correo electrónico institucional e-mail: gsilverio@uagro.mx ORCID: <https://orcid.org/0000-0002-2699-140X>

1. Introducción

El dengue es el arbovirus más importante transmitido por mosquitos, principalmente *Aedes aegypti*, que afecta a humanos [1]. Es considerado una enfermedad reemergente con impactos significativos en la salud pública global, especialmente en regiones tropicales y subtropicales de Asia y América Latina [2]. De acuerdo con la Organización Mundial de la Salud, anualmente se registran aproximadamente 100 millones de infecciones sintomáticas y 10,000 muertes a nivel mundial, lo que la convierte en una de las principales amenazas sanitarias en zonas endémicas.

En México, la incidencia de dengue ha mostrado un incremento considerable en los últimos años. Entre 2022 y 2024, la tasa nacional de incidencia pasó de 29.4 a 279.0 casos por cada 100,000 habitantes, lo que representa un aumento de casi diez veces. Además, el número de municipios afectados se elevó de 38.0 % en 2022 a 68.6 % en 2024, evidenciando una clara expansión territorial de la transmisión del virus [3].

La enfermedad presenta un amplio espectro clínico, que va desde cuadros asintomáticos hasta formas graves caracterizadas por fiebre, náuseas, vómitos, exantema y mialgias, con posibilidad de evolucionar hacia hemorragias y desenlaces fatales [4]. Esta variabilidad clínica, junto con la similitud de síntomas con otras enfermedades infecciosas como chikungunya, zika o influenza, complica el diagnóstico oportuno y la instauración de un manejo clínico adecuado.

Aunque el diagnóstico puede realizarse mediante criterios clínico-epidemiológicos o confirmarse por pruebas de laboratorio, en países como México donde existe una alta incidencia y limitaciones en recursos diagnósticos la confirmación por laboratorio no siempre es factible para todos los casos sospechosos. Esta situación resalta la necesidad de desarrollar estrategias complementarias, tales como el uso de modelos predictivos basados en inteligencia artificial [5] y algoritmos de aprendizaje automático [6] para mejorar la detección temprana, la clasificación automatizada y el análisis de riesgo de los casos.

Por lo que resulta fundamental evaluar no solo la capacidad predictiva de los

algoritmos, sino también su desempeño en conjuntos de datos desbalanceados, donde medidas de desempeño tradicionales como la exactitud pueden resultar engañosas. Por ello, este estudio plantea un enfoque que combina técnicas de reducción de atributos mediante feature importance [7] con la aplicación de modelos de machine learning del estado del arte, priorizando la medida de desempeño de balanced accuracy como indicador principal de desempeño. Esto se debe a que la exactitud refleja únicamente el porcentaje de predicciones correctas sin distinguir entre clases mayoritarias y minoritarias. En escenarios epidemiológicos como el del dengue, donde los casos negativos suelen superar en número a los positivos, un modelo que clasifique todos los registros como negativos podría alcanzar valores de exactitud elevados, pero sería inútil para fines de salud pública, ya que no detectaría a los pacientes que realmente requieren atención. En este contexto, métricas como el Balanced Accuracy ofrecen una evaluación más justa y equilibrada del desempeño.

En escenarios con conjuntos de datos desbalanceados, como es el caso del conjunto de datos del dengue, la medida de desempeño de exactitud puede ser engañosa. Por ejemplo, si el 80% de los casos corresponde a diagnósticos negativos y un modelo clasifica todos los registros como negativos, obtendría una exactitud aparente del 80%, a pesar de no identificar ningún caso positivo. Esta situación es crítica en salud pública, ya que los falsos negativos representan un riesgo considerable al impedir una detección temprana de pacientes que requieren atención oportuna. Para evitar esta distorsión, resulta más adecuado emplear medidas de desempeño como Balanced Accuracy que pondera de manera equitativa la sensibilidad y la especificidad, o el Coeficiente de Correlación de Matthews (MCC), que considera de forma conjunta verdaderos y falsos positivos y negativos, ofreciendo una visión más equilibrada del desempeño del modelo en escenarios de desbalance.

2. Datos y métodos

Este es un estudio cuantitativo realizado con datos secundarios provenientes del conjunto "Enfermedades Transmitidas por Vector" de los Datos Abiertos de la Dirección General de Epidemiología (DGE) de la Secretaría de Salud de México. El

conjunto de datos contiene registros individuales de casos notificados entre febrero de 2024 y febrero de 2025, en las 32 entidades federativas del país. El conjunto de datos está anonimizado y contiene información sobre edad, sexo, residencia, signos clínicos, comorbilidades, estado de defunción y resultado diagnóstico.

La investigación se estructuró en un conjunto de etapas metodológicas que abarcan desde la recopilación y preparación de los datos hasta la evaluación de los modelos de aprendizaje automático. Estas etapas se sintetizan en la figura 1.

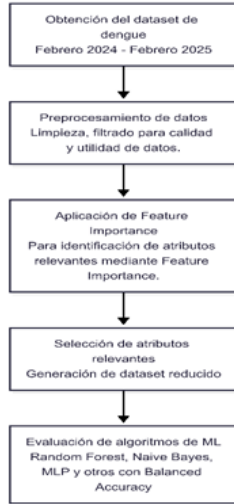


Figura 1. Etapas del proceso metodológico aplicado en la investigación.

2.1 Dataset y preprocesamiento

Con el apoyo de especialistas en salud pública, se definió la columna DICTAMEN como nuestra clase. Esta variable fue tratada como binaria, considerando únicamente los casos confirmados de dengue y los negativos, lo que permitió acotar el universo de estudio y garantizar un marco analítico más coherente. Con estas acciones se obtuvo el dataset final de trabajo, compuesta por 187 patrones y 22 atributos depurados y listos para análisis.

2.2 Algoritmo y análisis

Para determinar los atributos más relevantes en la clasificación de casos de dengue, se desarrolló un algoritmo de feature importance inspirado en la estructura de los árboles de decisión. Estos modelos, ampliamente documentados en la literatura, son reconocidos por su capacidad para dividir de manera recursiva un conjunto de datos en subgrupos homogéneos, basándose en medidas de impureza como el índice de Gini [8]. Esta métrica evalúa la heterogeneidad de las clases dentro de un nodo, de modo que una menor impureza indica una mejor separación de los datos.

El algoritmo implementado construyó un árbol binario evaluando, en cada nodo,

todos los atributos y sus posibles umbrales para identificar la división que maximizara la reducción de impureza. Cada atributo recibió una puntuación acumulativa proporcional a la disminución de impureza que generaba en el proceso de partición, y al finalizar la construcción del árbol dichas puntuaciones fueron normalizadas para generar un vector de importancia de atributos. Valores más altos reflejan una mayor relevancia en la tarea de clasificación [9].

El desarrollo del algoritmo se realizó en Python, evitando dependencias directas de librerías como scikit-learn, lo que permitió controlar la lógica de cálculo y adaptar el método a las necesidades específicas de la investigación. Se aplicó un umbral de 0.05 en la puntuación de importancia para seleccionar únicamente las variables más relevantes.

Para formalizar este procedimiento, en la figura 2 se presenta el pseudocódigo del cálculo de la importancia de atributos. Dicho algoritmo resume de manera estructurada las etapas del proceso: Inicialización del vector de importancias, evaluación de nodos mediante la impureza de Gini, búsqueda del mejor Split, acumulación de reducciones de impureza, normalización y sección de atributos significativos.

Algorithm 1: Cálculo de la importancia de atributos

Input: Matriz de datos X , etiquetas y
Output: Vector normalizado de importancias \hat{S}

- 1 Inicializar $S \leftarrow 0$
- 2 Función $\text{GrowTree}(X, y, \text{depth})$
- 3 if *nodo puro* o *profundidad máxima* then
- 4 | retornar hoja
- 5 Calcular $\text{imp_current} \leftarrow \text{Gini}(y)$
- 6 Buscar el mejor split (j^*, t^*)
- 7 Actualizar $S_j \leftarrow S_j + (\text{imp_current} - \text{imp_split})$
- 8 Llamar recursivamente a GrowTree en ramas izquierda y derecha
- 9 Llamar $\text{GrowTree}(X, y, 0)$
- 10 Normalizar $\hat{S}_j = S_j / \sum_k S_k$
- 11 Seleccionar atributos con $\hat{S}_j > \tau$

Figura 2. Cálculo de importancia de atributos

El flujo del algoritmo propuesto para la obtención de la importancia de características se resume en la figura 3.

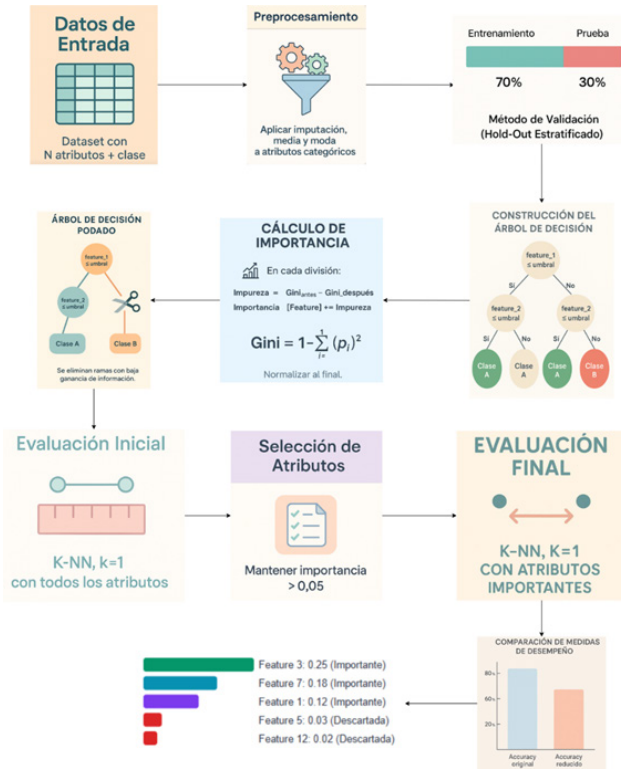


Figura 3. Diagrama de flujo del algoritmo de importancia de atributos propuesto.

2.2.1 Modelamiento matemático del algoritmo de importancia de atributos.

Sea un conjunto de datos:

$$D = \{(x_i, y_i)\}_{i=1}^n, \quad x_i \in \mathbb{R}^d, y_i \in \{0,1\}$$

Un árbol de decisión divide recursivamente el conjunto D en nodos $S \subseteq D$ mediante reglas de la forma $(x_j \leq t)$, donde j es el índice del atributo y t un umbral.

Ecuación 1 Índice de Gini (Impureza de un nodo)

$$G(S) = 1 - \sum_{c \in \{0,1\}} p_c(S)^2$$

Donde $P_c(S)$ es la proporción de observaciones de la clase c en el nodo S .

Ecuación 2 Impureza ponderada de una división

$$I(S; j, t) = \frac{|S_L|}{|S|} G(S_L) + \frac{|S_R|}{|S|} G(S_R)$$

Donde S_L y S_R son los subconjuntos izquierdo y derecho generados al aplicar el umbral t en el atributo j .

Ecuación 3 Disminución de impureza (ganancia)

$$\Delta(S; j^*, t^*) = G(S) - I(S; j^*, t^*)$$

Donde (j^*, t^*) corresponde a la división que minimiza la impureza ponderada.

Ecuación 4 Puntuación acumulada de un atributo

$$S_j = \sum_{u \in N: j(u)=j} \Delta u$$

Donde N es el conjunto de nodos internos y Δu la ganancia en el nodo u .

Ecuación 5 Normalización de la importancia

$$\hat{S}_j = \frac{S_j}{\sum_{k=1}^d S_k}$$

De modo que $\sum_j \hat{S}_j = 1$. Solo se conservan variables con $S_j > \tau$, siendo $\tau=0.06$ el umbral definido en este estudio.

Ecuación 6 Predicción en nodos hoja

$$\hat{y}(L) = \arg \max_{c \in \{0,1\}} P_c(L)$$

Donde L es una hoja del árbol y $P_c(L)$ la proporción de instancias de la clase c .

2.3 Validación de algoritmos

Para evaluar de manera justa el rendimiento de los algoritmos de machine learning y del algoritmo de feature importance, se aplicó un esquema de validación Hold-Out estratificado, con un 70 % de los datos

para entrenamiento y un 30 % para prueba. El uso de estratificación garantizó que la distribución de las clases en ambos conjuntos fuera representativa, lo que permitió obtener medidas de desempeño más consistentes, especialmente en un conjunto de datos desbalanceado.

3. Resultados

A continuación se muestran los resultados de la metodología una vez aplicado el algoritmo feature importance, en la Tabla 1, se pueden observar los resultados de los algoritmos mayormente utilizados en el estado del arte.

Tabla 1. Conjunto de datos original (22 atributos)

Algoritmo	Medidas de desempeño					
	Recall	Specificity	Balanced Accuracy	Precision	F1-Score	MCC
1-NN [10]	0.92	0.73	0.83	0.90	0.91	0.67
3-NN	0.66	0.53	0.60	0.80	0.72	0.18
5-NN	0.73	0.40	0.56	0.77	0.75	0.13
Naive Bayes [11]	0.54	0.40	0.47	0.71	0.62	-0.04
J48 [12]	0.80	0.66	0.73	0.87	0.83	0.45
Random Forest [13]	0.90	0.53	0.71	0.84	0.87	0.47
SVM [14]	1.0	0.0	0.50	0.73	0.84	0.0
Regresión Logística [15]	0.71	0.40	0.55	0.76	0.74	0.10
MLP [16]	1.0	0.0	0.50	0.73	0.84	0.0

En ella se reportan de igual manera las medidas de desempeño más reportadas en el estado del arte, como lo son recall, specificity y balanced accuracy, precision, MCC y F1-Score, esto para darnos un mejor panorama del comportamiento, robustez y capacidad de generalización de los algoritmos en el conjunto de datos utilizado.

De acuerdo con la metodología, al aplicar el algoritmo de feature importance nos quedamos con los atributos que presentan una mayor representación de los datos que conforman el conjunto de datos, es decir, solo se presentan aquellos atributos que tienen mayor grado de significancia con la que puede trabajar los algoritmos para poder realizar la tarea de clasificación. Entonces, como el objetivo del trabajo es demostrar que pueden mantenerse consistentes o mejorar los valores de las medidas de desempeño de los algoritmos, se presenta la Tabla 2 con las mismas medidas de desempeño de la Tabla 1 pero con el conjunto de datos conformado únicamente con N atributos (N = 7 después de aplicar feature importance). Las variables seleccionadas fueron: la edad del paciente (0.3325), la entidad de residencia (0.1220), la presencia de hipertensión (0.1081), la entidad donde se ubica la unidad médica que notificó el caso (0.1042), el municipio de la unidad médica notificante (0.0812), la entidad de asignación médica del paciente (0.0769) y el municipio de residencia del paciente (0.0726).

Tabla 2. Conjunto de datos reducido (7 atributos)

Algoritmo	Medidas de desempeño					
	Recall	Specificity	Balanced Accuracy	Precision	F1-Score	MCC
1-NN [10]	0.92	0.73	0.83	0.90	0.91	0.67
3-NN	0.66	0.53	0.60	0.80	0.72	0.18
5-NN	0.73	0.40	0.56	0.77	0.75	0.13
Naive Bayes [11]	0.54	0.40	0.47	0.71	0.62	-0.04
J48 [12]	0.80	0.66	0.73	0.87	0.83	0.45
Random Forest [13]	0.90	0.53	0.71	0.84	0.87	0.47
SVM [14]	1.0	0.0	0.50	0.73	0.84	0.0
Regresión Logística [15]	0.71	0.40	0.55	0.76	0.74	0.10
MLP [16]	1.0	0.0	0.50	0.73	0.84	0.0

El conjunto de datos utilizado es desbalanceado (IR>2.00) por lo tanto, reportar la medida de desempeño no es correcto, esto es porque al usar exactitud, el resultado del algoritmo respecto de este valor puede verse enormemente beneficiado de resultar la clase mayoritaria como “bien clasificada”, es por eso por lo que en la Tabla 3, se presenta a manera de comparativa, los resultados que obtuvieron los algoritmos, respecto de balanced accuracy, que se pusieron a prueba.

Tabla 3. tabla comparativa (Balanced Accuracy)

Algoritmo	Balanced Accuracy (22 atributos)	Balanced accuracy (7 atributos)
1-NN	0.83	0.87
3-NN	0.60	0.66
5-NN	0.56	0.64
Naive Bayes	0.47	0.65
J48	0.73	0.71
Random Forest	0.71	0.69
SVM	0.50	0.52
Regresión logística	0.55	0.61
MLP	0.50	0.74

Los resultados presentados en la Tabla 3 muestran una clara tendencia a la mejora en la medida de desempeño Balanced Accuracy para la mayoría de los algoritmos al trabajar con el conjunto de datos reducido de atributos. La Balanced Accuracy es una medida de desempeño diseñada para evaluar modelos en escenarios de clases desbalanceadas, ya que calcula el promedio entre la sensibilidad y la especificidad, otorgando el mismo peso a las clases mayoritarias y minoritarias. Su cálculo se expresa mediante la siguiente ecuación:

Ecuación 7. Fórmula de Balanced Accuracy

$$\text{Balanced Accuracy} = \frac{\text{Sensitivity} + \text{Specificity}}{2}$$

Al observar los resultados, algoritmos como 1-NN, 3-NN, 5-NN y Naive Bayes muestran incrementos notables en su desempeño, lo que sugiere que se beneficiaron de un

espacio de características más compacto y representativo, favoreciendo su capacidad de generalización. Destaca el caso de Naive Bayes, que pasó de un Balanced Accuracy de 0.47 a 0.65, y MLP, que logró un aumento significativo de 0.50 a 0.74, demostrando que la reducción de atributos optimizó su rendimiento.

Por otro lado, algoritmos como Random Forest y J48 mostraron valores relativamente estables, lo cual es consistente con su naturaleza robusta ante variables irrelevantes, mientras que SVM y Regresión Logística presentan mejoras más discretas.

4. Discusión

Los resultados de este estudio aportan evidencia sobre la relevancia de integrar técnicas de feature importance en el desarrollo de modelos de aprendizaje automático para la clasificación de casos de dengue. Al aplicar un algoritmo inspirado en árboles de decisión, fue posible reducir el conjunto de datos de 22 a 7 atributos clave sin comprometer la capacidad predictiva de los modelos, e incluso mejorando el desempeño en métricas críticas como balanced accuracy, lo que es especialmente significativo considerando el desbalance de clases presente en los datos. Un hallazgo notable es que la mayoría de las variables seleccionadas como más influyentes están relacionadas con la localización geográfica del paciente o con la ubicación de la unidad médica notificante (por ejemplo, entidad y municipio de residencia, así como entidad y municipio de la unidad médica). Esta tendencia sugiere que la distribución territorial de los casos podría ser un factor determinante en la clasificación y que los patrones espaciales de los datos encierran información valiosa sobre la dinámica de la enfermedad.

A partir de este hallazgo, se recomienda complementar el enfoque actual con un análisis geoespacial que permita mapear los casos confirmados y negativos de dengue, así como sus relaciones con variables clínicas y ambientales.

Los resultados refuerzan la utilidad de los modelos de aprendizaje automático apoyados en técnicas de selección de características, al tiempo que abren la puerta a integrar metodologías de geovisualización en futuras investigaciones. De esta forma, los modelos predictivos no solo ganarían precisión y eficiencia, sino que también aportarían una dimensión espacial esencial para comprender y abordar la problemática

del dengue en México.

5. Conclusiones

Este estudio demuestra que la integración de técnicas de selección de atributos mediante feature importance en modelos de aprendizaje automático puede optimizar significativamente la clasificación de casos de dengue en México. Ya que se logró simplificar el conjunto de datos sin comprometer, e incluso mejorando, métricas clave como la balanced accuracy, especialmente en algoritmos sensibles al ruido como Naive Bayes y MLP.

La reducción de variables no solo incrementa la interpretabilidad de los modelos, aspecto fundamental en el ámbito de la salud pública, sino que también disminuye la carga computacional, lo que facilita su implementación en sistemas con recursos limitados. Esto permite que las instituciones de salud cuenten con herramientas más ligeras, explicables y escalables, que pueden integrarse a plataformas digitales y aplicaciones móviles para fortalecer la vigilancia epidemiológica.

Los resultados aportan una solución práctica a la necesidad de modelos predictivos más eficientes y accesibles, contribuyendo a la detección temprana y clasificación de casos de dengue, y ofreciendo una base sólida para el desarrollo de sistemas de apoyo a la toma de decisiones en contextos de alta incidencia y recursos diagnósticos limitados.

6. Referencias

- [1] World Health Organization, "Dengue and severe dengue," WHO Fact Sheet, Oct. 2023. [Online]. Available: <https://www.who.int/news-room/fact-sheets/detail/dengue-and-severe-dengue>
- [2] S. W. Huang, H. P. Tsai, S. J. Hung, W. C. Ko, and J. R. Wang, "Assessing the risk of dengue severity using demographic information and laboratory test results with machine learning," PLoS Neglected Tropical Diseases, vol. 14, no. 12, p. e0008960, Dec. 2020. [Online]. Available: <https://doi.org/10.1371/journal.pntd.0008960>

- [3] O. Mendoza-Cano et al., "Spatial patterns and clustering of dengue incidence in Mexico: Analysis of Moran's index across 2,471 municipalities from 2022 to 2024", PLOS One, vol. 20, núm. 5, p. e0324754, may 2025, doi: 10.1371/journal.pone.0324754.
- [4] S. B. Halstead, "Dengue," The Lancet, vol. 370, no. 9599, pp. 1644–1652, Nov. 2007. [Online]. Available: [https://doi.org/10.1016/S01406736\(07\)6160](https://doi.org/10.1016/S01406736(07)6160)
- [5] S. W. Huang, H. P. Tsai, S. J. Hung, W. C. Ko, and J. R. Wang, "Assessing the risk of dengue severity using demographic information and laboratory test results with machine learning," PLoS Neglected Tropical Diseases, vol. 14, no. 12, p. e0008960, Dec. 2020. [Online]. Available: <https://doi.org/10.1371/journal.pntd.0008960>
- [6] C. Carvajal, C. Benavides, and P. Parra, "Machine learning models to predict dengue outbreaks: A comparison of approaches in Colombia," International Journal of Medical Informatics, vol. 117, pp. 62–73, 2018. [Online]. Available: <https://doi.org/10.1016/j.ijmedinf.2018.04.4>
- [7] I. Guyon y A. Elisseeff, "An introduction to variable and feature selection," Journal of Machine Learning Research, vol. 3, pp. 1157–1182, 2003.
- [8] L. Breiman, J. Friedman, R. Olshen y C. Stone, Classification and Regression Trees. Belmont, CA: Wadsworth, 1984.
- [9] G. Louppe, L. Wehenkel, A. Suter, y P. Geurts, "Understanding variable importances in forests of randomized trees," Advances in Neural Information Processing Systems, vol. 26, 2013.
- [10] D. W. Aha, Lazy Learning. Boston, MA, USA: Springer Science+Business Media, 1997.
- [11] T. M. Mitchell, Machine Learning. New York, NY, USA: McGraw-Hill, 1997.
- [12] J. R. Quinlan, C4.5: Programs for Machine Learning. San Mateo, CA, USA: Morgan Kaufmann Publishers, 1993.
- [13] L. Breiman, "Random forests," Machine Learning, vol. 45, no. 1, pp. 5–32, 2001, doi: 10.1023/A:1010933404324.
- [14] C. Cortes and V. Vapnik, "Support-vector networks," Machine Learning, vol. 20, no. 3, pp. 273–297, 1995, doi: 10.1007/BF00994018.
- [15] D. R. Cox, "The regression analysis of binary sequences," Journal of the Royal Statistical Society: Series B (Methodological), vol. 20, no. 2, pp. 215–242, 1958.
- [16] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," Nature, vol. 323, no. 6088, pp. 533–536, 1986, doi: 10.1038/323533a0.