

# Análisis y Modelado Predictivo de Indicadores de Salud para la Detección de Diabetes mediante Machine Learning

1<sup>st</sup> Emanuel Vasquez Yepes  
Facultad de Ingeniería  
Universidad de Antioquia  
Medellín, Colombia  
emanuel.vasquea@udea.edu.co

2<sup>nd</sup> Jose Andres Echavarria Rios  
Facultad de Ingeniería  
Universidad de Antioquia  
Medellín, Colombia  
jose.echavarria1@udea.edu.co

## I. DESCRIPCIÓN DEL PROBLEMA

### 1. Contexto del problema

La diabetes constituye una de las enfermedades crónicas más prevalentes a nivel mundial, afectando a millones de personas y generando un impacto económico y sanitario significativo. Este padecimiento se caracteriza por la incapacidad del organismo para regular adecuadamente los niveles de glucosa en sangre, lo cual puede derivar en complicaciones severas como enfermedades cardiovasculares, pérdida de visión o daño renal. La detección temprana es un factor crítico para mejorar la calidad de vida y reducir la mortalidad asociada.

En este contexto, el desarrollo de modelos predictivos basados en *Machine Learning* (ML) permite identificar individuos en riesgo y apoyar la toma de decisiones médicas y de salud pública, optimizando recursos y facilitando estrategias preventivas más eficaces.

Esto puede ser de gran utilidad para la detección temprana y prevención de la enfermedad, así como para apoyar tomas de decisiones médicas basadas en datos.

### 2. Composición de la base de datos

El conjunto de datos utilizado proviene del *Behavioral Risk Factor Surveillance System* (BRFSS) 2015, recopilado por el *Centers for Disease Control and Prevention* (CDC). La base procesada contiene 22 variables, de las cuales 21 son predictoras y una corresponde a la variable objetivo (Diabetes\_012), con tres clases: 0 para no diabetico, 1 para prediabetico y 2 para diabetico. Inicialmente se contaba con 253,680 registros, observandose un marcado desbalance de clases, predominando la categoría sin diabetes. Para mitigar este problema y mejorar la capacidad de generalización del modelo, se aplicó una estrategia de submuestreo, resultando en un conjunto final de 38,052 instancias. Todas las variables son de tipo numérico, pero algunas variables tienen valores numéricos que en realidad representan categorías, por lo que necesitan una codificación adicional antes de ser utilizadas en los modelos.

El conjunto de datos “Diabetes Health Indicators Dataset” (extraído de la base Behavioral Risk Factor Surveillance

System — BRFSS 2015) contiene información proveniente de encuestas realizadas en Estados Unidos, enfocadas en hábitos de vida y condiciones de salud de los participantes.

A continuación, presentaremos el dataset antes y después del proceso de submuestreo donde explicaremos a fondo la distribución del mismo y haremos algunos énfasis en los resultados obtenidos en el análisis exploratorio.

#	Variable	Tipo de dato	Valores no nulos
0	Diabetes_012	float64	253,680
1	HighBP	float64	253,680
2	HighChol	float64	253,680
3	CholCheck	float64	253,680
4	BMI	float64	253,680
5	Smoker	float64	253,680
6	Stroke	float64	253,680
7	HeartDiseaseorAttack	float64	253,680
8	PhysActivity	float64	253,680
9	Fruits	float64	253,680
10	Veggies	float64	253,680
11	HvyAlcoholConsump	float64	253,680
12	AnyHealthcare	float64	253,680
13	NoDocbcCost	float64	253,680
14	GenHlth	float64	253,680
15	MentHlth	float64	253,680
16	PhysHlth	float64	253,680
17	DiffWalk	float64	253,680
18	Sex	float64	253,680
19	Age	float64	253,680
20	Education	float64	253,680
21	Income	float64	253,680

Cuadro I: Estructura del conjunto de datos BRFSS 2015. Todas las variables son de tipo float64 y no presentan valores nulos.

A continuación, explicaremos cada una de estas variables y como se encuentran codificadas actualmente en el dataset, esto nos ayudara a contextualizar como se tomarán decisiones y dará un mejor contexto del problema.

Variable	Tipo	Descripción
<b>Diabetes_012</b>	Catagórica (0, 1, 2)	<b>Variable objetivo.</b> Indica el estado de diabetes: 0: No diabetes, 1: Prediabetes, 2: Diabetes.
<b>HighBP</b>	Binaria (0, 1)	Diagnóstico de presión arterial alta.
<b>HighChol</b>	Binaria (0, 1)	Diagnóstico de colesterol alto.
<b>CholCheck</b>	Binaria (0, 1)	Revisión de colesterol en los últimos 5 años.
<b>BMI</b>	Númerica (float)	Índice de Masa Corporal (peso / altura <sup>2</sup> ).
<b>Smoker</b>	Binaria (0, 1)	Ha fumado al menos 100 cigarrillos en su vida.
<b>Stroke</b>	Binaria (0, 1)	Ha sufrido un accidente cerebrovascular.
<b>HeartDiseaseorAttack</b>	Binaria (0, 1)	Enfermedad cardíaca o ataque al corazón previo.
<b>PhysActivity</b>	Binaria (0, 1)	Realizó actividad física en los últimos 30 días.
<b>Fruits</b>	Binaria (0, 1)	Consume frutas al menos una vez al día.
<b>Veggies</b>	Binaria (0, 1)	Consume vegetales al menos una vez al día.
<b>HvyAlcoholConsump</b>	Binaria (0, 1)	Consumo excesivo de alcohol: más de 14 bebidas/semana (hombres) o más de 7 (mujeres).
<b>AnyHealthcare</b>	Binaria (0, 1)	Tiene cobertura médica o acceso a servicios de salud.
<b>NoDocbcCost</b>	Binaria (0, 1)	No pudo acudir al médico en el último año por costo.
<b>GenHlth</b>	Ordinal (1–5)	Salud general autoevaluada: 1 = Excelente, 5 = Mala.
<b>MentHlth</b>	Númerica (0–30)	Días con mala salud mental en el último mes.
<b>PhysHlth</b>	Númerica (0–30)	Días con mala salud física en el último mes.
<b>DiffWalk</b>	Binaria (0, 1)	Dificultad para caminar o subir escaleras.
<b>Sex</b>	Binaria (0, 1)	Sexo biológico: 0 = Mujer, 1 = Hombre.
<b>Age</b>	Catagórica (1–13)	Grupo de edad: 1 = 18–24, 2 = 25–29, ..., 13 = 80+.
<b>Education</b>	Ordinal (1–6)	Nivel educativo: 1 = Ninguno, 6 = Título universitario.
<b>Income</b>	Ordinal (1–8)	Nivel de ingreso: 1 = <10k USD, 8 = 75k USD.

Cuadro II: Descripción de las variables del conjunto de datos de indicadores de salud (BRFSS).

Es importante aclarar que, aunque todas las variables del conjunto de datos se encuentran almacenadas con tipo `float64`, la mayoría representan categorías discretas o respuestas binarias (0 o 1). Por tanto, se consideran variables categóricas codificadas numéricamente. La única variable verdaderamente continua es `BMI`, que corresponde al índice de masa corporal y presenta una distribución de valores reales en el rango aproximado de 12 a 98.

El siguiente paso será validar y revisar la distribución de los datos en la variable objetivo:

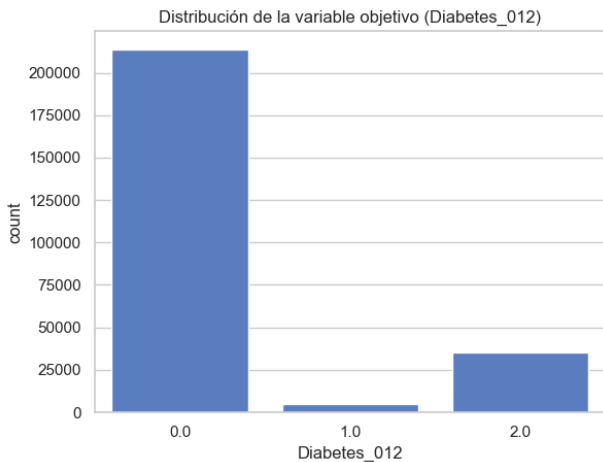


Figura 1: Distribución de la variable objetivo `Diabetes_012`. Se observa un claro desbalance entre las clases, predominando la categoría 0 (no diabético).

La variable objetivo `Diabetes_012` presenta un marcado desbalance entre clases, lo que puede afectar negativamente el desempeño de los modelos de clasificación. Como se observa en la Tabla IV, la mayoría de los registros corresponden a individuos sin diagnóstico de diabetes, mientras que las clases de prediabetes y diabetes representan una proporción mucho menor.

Clase ( <b>Diabetes_012</b> )	Proporción (%)	Cantidad
0 (No diabético)	84.24	213703
1 (Prediabético)	1.83	4631
2 (Diabético)	13.93	35346

Cuadro III: Distribución de clases en la variable objetivo `Diabetes_012`.

La matriz de correlación muestra cómo se relacionan entre sí las variables numéricas del dataset. Cada valor en la matriz representa un coeficiente de correlación pearson entre dos variables. A continuación se presenta la matriz obtenida para el conjunto de datos

La matriz de correlación permite identificar relaciones lineales entre las variables del conjunto de datos. Se observan correlaciones positivas moderadas entre variables relacionadas con condiciones de salud, como `HighBP` y `HeartDiseaseorAttack`, lo que sugiere que la hipertensión puede coexistir con enfermedades cardíacas. Por otro lado, variables como `Age`, `BMI` y `GenHlth` presentan correlaciones más bajas, indicando que su relación con el diagnóstico de diabetes puede ser más compleja y no estrictamente lineal.

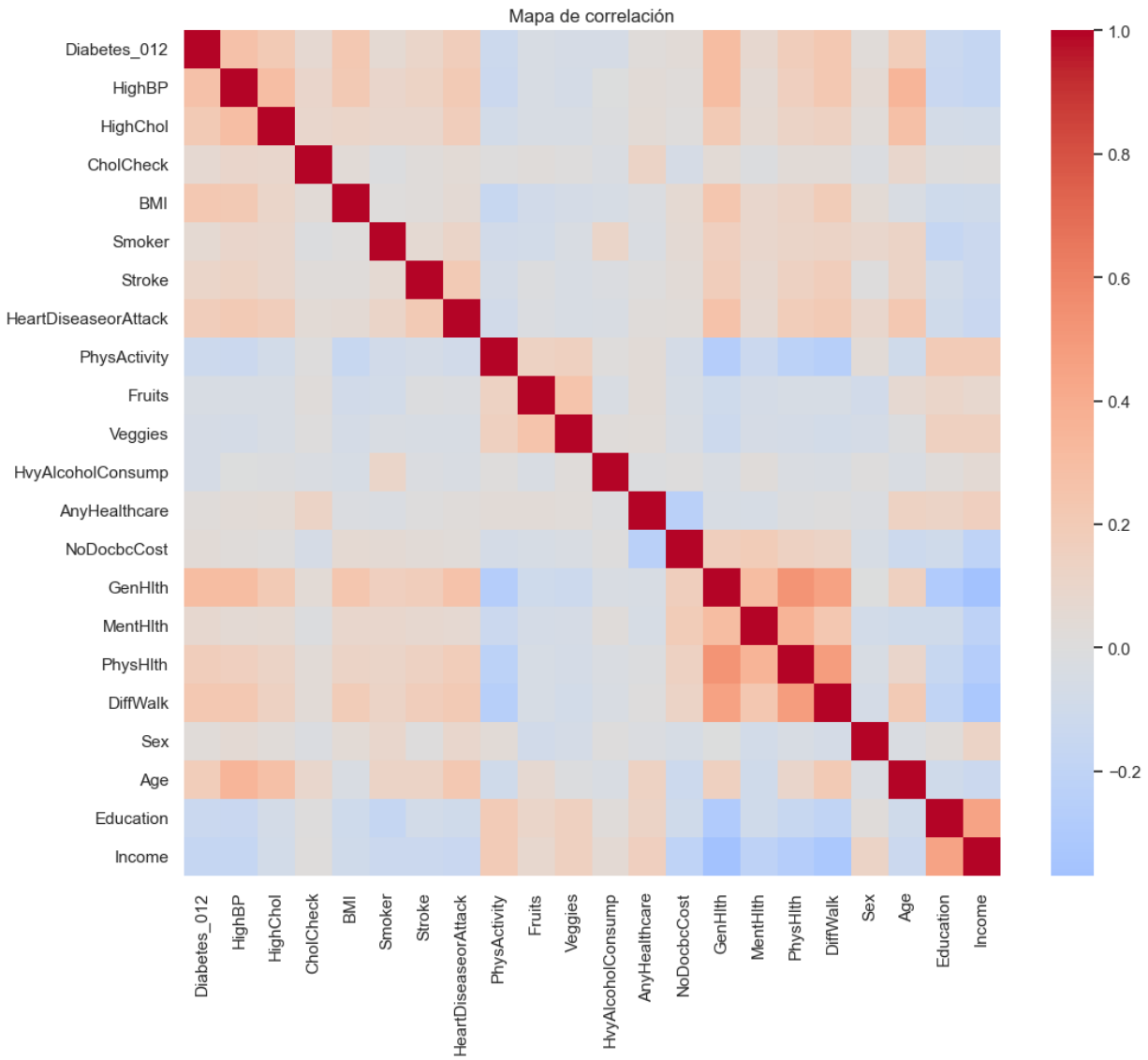


Figura 2: Matriz de correlación Diabetes\_012.

Para finalizar la exploracion del dataset se presentara de manera analoga al dataset completo los datos obtenidos mediante el submuestreo y se dara una breve justificación del por que de este:

Para reducir el costo computacional sin modificar la estructura estadística del conjunto de datos, se aplicó un submuestreo estratificado, el cual selecciona aleatoriamente una fracción del conjunto original manteniendo las proporciones de la variable objetivo. De esta forma, se preserva el desbalance natural del problema, lo que permite evaluar posteriormente el impacto de dicho desbalance en el rendimiento del modelo.

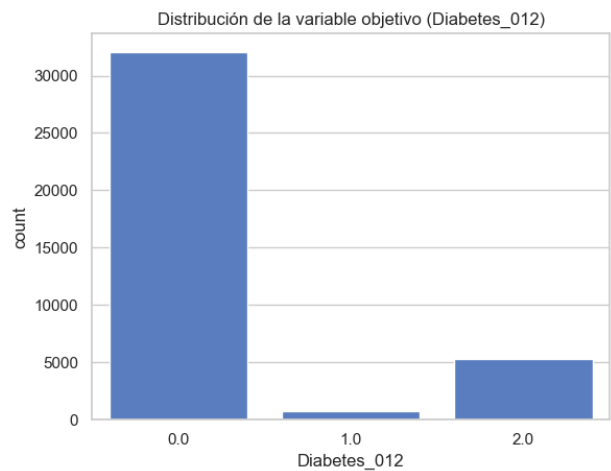


Figura 3: Distribución de la variable objetivo despues del submuestreo Diabetes\_012. Se observa que se mantiene el desbalance en las tres clases

Clase (Diabetes_012)	Proporción (%)	Cantidad
0 (No diabético)	84.24	32055
1 (Prediabético)	1.82	695
2 (Diabético)	13.93	5302

Cuadro IV: Distribución de clases en la variable objetivo Diabetes\_012 después del submuestreo

Podemos observar que después del submuestreo los datos conservan las mismas características del dataset original, manteniendo la correlación igual y el desbalance. Con esto garantizamos que, aunque estemos usando una porción de los datos el problema original que consiste en enfrentarse al efecto del desbalance se mantiene.

En resumen, el análisis exploratorio del conjunto de datos permitió comprender la naturaleza y composición de las variables empleadas para la predicción de diabetes. Se identificó que la mayoría de las variables corresponden a características categóricas codificadas numéricamente, con excepción del índice de masa corporal (BMI), que presenta valores continuos. El análisis de la variable objetivo evidenció un marcado desbalance entre clases, predominando los casos sin diagnóstico de diabetes.

Asimismo, la matriz de correlación permitió observar asociaciones moderadas entre variables relacionadas con la salud cardiovascular, como HighBP y HeartDiseaseorAttack, lo que respalda su relevancia en la detección del riesgo de diabetes. Finalmente, mediante el submuestreo estratificado se logró reducir el tamaño del conjunto de datos sin alterar su distribución original, conservando el desbalance inherente del problema.

### 3. Paradigma de aprendizaje y justificación

Dado que la variable objetivo representa tres categorías mutuamente excluyentes, el problema se aborda como una tarea de clasificación supervisada multiclase. Este paradigma permite al modelo aprender patrones asociados con diferentes niveles de riesgo de diabetes a partir de las características sociodemográficas y de comportamiento reportadas en la encuesta. La elección de este enfoque se fundamenta en la necesidad de discriminar entre distintos estados de salud (sin diabetes, prediabetes y diabetes) de manera precisa, con el objetivo de contribuir al desarrollo de herramientas predictivas que puedan integrarse en sistemas de apoyo a la decisión médica y en programas de monitoreo poblacional.

Además, para abordar la naturaleza multiclase del problema y la mezcla de variables numéricas y categóricas, se evaluaron cinco modelos representativos de distintos paradigmas de aprendizaje: un modelo paramétrico (Regresión Logística), un modelo no paramétrico (KNN), un ensamble de árboles de decisión (Random Forest), una red neuronal artificial (MLP) y una máquina de vectores de soporte (SVM). La comparación se realizó mediante una malla de hiperparámetros diseñada para cada algoritmo, con el fin de analizar su impacto en el desempeño y seleccionar la configuración más adecuada.

## II. ESTADO DEL ARTE

### Artículo 1: Diabetic Mellitus Prediction with BRFSS Data Sets [1]

- **Paradigma:** Aprendizaje supervisado. El estudio aplica distintos algoritmos de clasificación para predecir la diabetes tipo 2 en función de variables de salud obtenidas de los conjuntos de datos BRFSS 2014 y 2015.
- **Técnicas:** Se utilizaron múltiples modelos de Machine Learning: Árbol de Decisión, Regresión Logística, SVM (lineal, RBF y polinomial), Naïve Bayes, Random Forest y Redes Neuronales. Además, se aplicó SMOTE para balancear las clases (diabético, prediabético y no diabético).
- **Validación:** Se realizó una división del conjunto de datos en entrenamiento (dos tercios) y prueba (un tercio). No se especifica validación cruzada, pero se comparan los modelos usando los mismos subconjuntos.
- **Métricas:** Exactitud (Accuracy), Sensibilidad (Sensitivity), Especificidad (Specificity) y Área Bajo la Curva (AUC).
- **Resultados:** En el conjunto BRFSS 2014, la Red Neuronal obtuvo la mayor exactitud (82.4 %) y especificidad (90.2 %). El Árbol de Decisión logró la mayor sensibilidad (51.6 %). Para 2015, los modelos CatBoost y XGBoost alcanzaron 86 % de exactitud, superando a los demás clasificadores.
- **Fuente:** M. H. Mohamed et al., "Diabetic Mellitus Prediction with BRFSS Data Sets," *Journal of Theoretical and Applied Information Technology*, vol. 102, no. 3, 2024. Disponible en: [www.jatit.org/volumes/Vol102No3/10Vol102No3.pdf](http://www.jatit.org/volumes/Vol102No3/10Vol102No3.pdf)

### Artículo 2: An investigation of machine learning algorithms and data augmentation techniques for diabetes diagnosis using class imbalanced BRFSS dataset [2]

- **Paradigma:** Aprendizaje supervisado enfocado en el manejo del desbalance de clases. El estudio utiliza el conjunto de datos BRFSS (Behavioral Risk Factor Surveillance System) para diagnosticar diabetes, aplicando técnicas de aumento y balanceo antes del entrenamiento de los modelos.
- **Técnicas:** Se emplearon métodos de muestreo como oversampling (SMOTE-N, una versión adaptada de SMOTE para datos nominales), undersampling (Edited Nearest Neighbour, ENN) y enfoques híbridos (SMOTE-Tomek, SMOTE-ENN). Los algoritmos de Machine Learning utilizados incluyen Regresión Logística, Random Forest, Gradient Boosting y AdaBoost, todos evaluados sobre los diferentes conjuntos de datos generados.
- **Validación:** Los autores enfatizan la prevención de *data leakage* (fuga de datos) al aplicar las técnicas de muestreo antes de la división de datos de entrenamiento y prueba. Esto garantiza que los datos aumentados no contaminen el conjunto de validación. La comparación de modelos se realiza bajo las mismas condiciones experimentales.

- **Métricas:** Se usaron métricas estándar de clasificación, principalmente Exactitud (Accuracy) y Sensibilidad (Recall), para evaluar el efecto del balanceo y comparar el rendimiento entre algoritmos.
- **Resultados:** Antes del balanceo los modelos alcanzaron exactitudes cercanas al 83 %. Tras aplicar distintas estrategias, la técnica *ENN (Edited Nearest Neighbour)* con Gradient Boosting obtuvo el mejor desempeño (Accuracy 70.3 %, AUC 0.791). Entre los factores más influyentes se identificaron la edad, el índice de masa corporal (IMC) y la presión arterial alta, coherentes con hallazgos clínicos previos.
- **Fuente:** M. M. Chowdhury, R. S. Ayon, M. S. Hossain, “An investigation of machine learning algorithms and data augmentation techniques for diabetes diagnosis using class imbalanced BRFS dataset,” *Healthcare Analytics*, vol. 5, p. 100297, 2024. DOI: doi.org/10.1016/j.health.2023.100297

### III. ENTRENAMIENTO Y EVALUACIÓN DE LOS MODELOS

#### 1. Configuración experimental

Para el entrenamiento y evaluación de los modelos se utilizó un esquema de validación basado en *validación cruzada*, lo que permitió estimar el desempeño de cada configuración de hiperparámetros de manera más confiable.

Se trabajó directamente con los datos originales disponibles, sin aplicar técnicas adicionales de modificación en la distribución de clases. De esta manera, los modelos fueron entrenados y evaluados sobre la información tal como fue proporcionada.

Posteriormente, todas las variables fueron sometidas a un proceso de *escalado* mediante la técnica de *StandardScaler*, con el fin de normalizar las características y facilitar el entrenamiento de los modelos, especialmente aquellos sensibles a la magnitud de los datos como las redes neuronales y las máquinas de soporte vectorial.

#### 2. Conjunto de Hiperparámetros Analizados

En esta sección se presentan los hiperparámetros considerados durante el entrenamiento y evaluación de los modelos de aprendizaje automático. El objetivo fue identificar configuraciones que ofrecieran un mejor desempeño en términos de la métrica de *Balanced Accuracy*.

Se definieron distintas combinaciones de parámetros y se evaluaron mediante *validación cruzada*, lo que permitió comparar el rendimiento de manera confiable sin depender de un único conjunto de datos.

Los modelos incluidos fueron: Regresión Logística, KNN, Balanced Decision Forest, MLP y SVM. Con esta selección se cubrieron diferentes enfoques del aprendizaje automático y se obtuvo una visión amplia sobre cómo cada tipo de modelo responde a las condiciones del problema y a la variación de sus hiperparámetros.

Modelo	Hiperparámetros	Malla de valores
Regresión Logística	C, penalty, solver, max_iter, class_weight	C = [0.001, 0.01, 0.1, 1, 10, 100]; penalty = [l2]; solver = [lbfgs, saga]; max_iter = [1000]; class_weight = [balanced]
KNN	n_neighbors, weights, p	n_neighbors = [3, 5, 7, 9, 11, 15, 21, 31]; weights = [uniform, distance]; p = [1, 2]
MLP	hidden_layer_sizes, activation, solver, alpha, learning_rate_init	hidden_layer_sizes = [(64,),(128,),(64,64),(128,64,32),(32,16,8)]; activation = [relu, tanh]; solver = [adam]; alpha = [0.0001, 0.001]; learning_rate_init = [0.001, 0.01]
SVM	kernel, C, gamma	kernel = [linear, rbf]; C = [0.01, 0.1, 1, 10]; gamma = [scale, auto]
Balanced Decision Forest	max_depth, min_samples_leaf, criterion, n_estimators	max_depth = [3, 5, 7, 8]; min_samples_leaf = [20, 30, 50, 100]; criterion = [gini, entropy]; n_estimators = [50, 100]

Cuadro V: Conjunto de hiperparámetros analizados en cada modelo de aprendizaje automático.

Como se observa en la tabla, cada modelo fue configurado con parámetros específicos que permitieron explorar distintos niveles de complejidad. Esto garantizó que la comparación entre modelos se realizara en condiciones similares y de manera ordenada.

#### 3. Métricas de Desempeño

Para la evaluación del sistema se utilizarán métricas que permitan medir de manera adecuada el rendimiento en un contexto de clases desbalanceadas. La métrica principal será la *Balanced Accuracy*, ya que ofrece una visión más justa del desempeño al considerar el promedio de las tasas de acierto en cada clase, evitando que el resultado se vea dominado por la clase mayoritaria. De manera complementaria, se analizarán otras métricas como la *precision* y el *F1-score*, además de la métrica de *accuracy* tradicional únicamente como referencia, reconociendo que en escenarios con desbalance puede resultar poco representativa.

Este conjunto de métricas permitirá obtener una evaluación más completa y confiable del sistema, ofreciendo una visión equilibrada del rendimiento de los modelos bajo diferentes condiciones.

#### 4. Resultados del entrenamiento de Modelos

En esta sección se presentan los resultados de la experimentación para los modelos evaluados. Se incluyen métricas de *Balanced Accuracy* en entrenamiento, validación y prueba, junto con intervalos de confianza estimados en la fase de validación. El análisis permite observar el efecto de los hiperparámetros en el desempeño y comparar el rendimiento global de cada modelo.

La Figura 4 muestra la comparación de desempeño en los conjuntos de entrenamiento, validación y prueba, así como el ranking final de los modelos en el conjunto de prueba. Se observa que el modelo SVM alcanza la mejor *Balanced Accuracy* en test (0.5376), seguido por Logistic Regression

y Balanced Decision Forest. En contraste, KNN y la Red Neuronal MLP presentan menor capacidad de generalización.

En conclusión, el modelo SVM se selecciona como el mejor candidato para el sistema final, dado que logra el mayor desempeño en el conjunto de prueba.

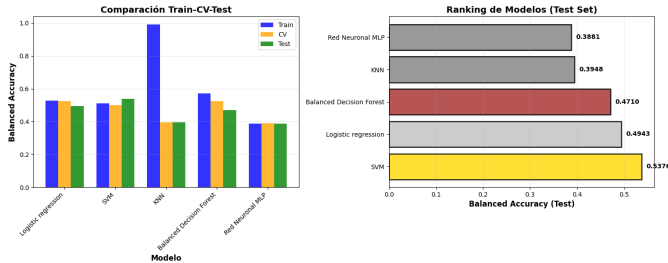


Figura 4: Comparación de desempeño de los modelos.

Modelo	BA Train	BA CV	BA Test
Logistic Regression	0.527378	0.523217	0.494334
SVM	0.509800	0.501200	0.537600
KNN	0.990063	0.395228	0.394795
Balanced Decision Forest	0.570369	0.525200	0.470964
Red Neuronal MLP	0.386501	0.389275	0.388096

Cuadro VI: Desempeño de los modelos.

Analizando más a fondo en general todos los modelos presentan márgenes de predicción en los rangos del 40 – 54 por ciento aproximadamente, lo que constituye en una aproximación no tan adecuada para el mercado bajo el cual se fundamenta y se aproxima el modelo, ya que podemos considerar que puede producir una amplia cantidad de falsos positivos que perjudicarían integralmente la vida y el tratamiento oportuno de un paciente.

#### IV. REDUCCIÓN DE DIMENSIÓN

##### 1. Análisis individual de variables

Se evaluó la capacidad discriminativa de cada característica frente a la variable objetivo mediante pruebas estadísticas según su tipo (Chi-cuadrado para binarias y ANOVA F-test para ordinales y continuas). Los resultados permiten identificar variables con baja relevancia que podrán ser eliminadas en la reducción de dimensión.

Cabe resaltar que para regresión logística se evalúan los pesos para clasificar si una variable puede ser eliminada donde si al calcular el coeficiente es número bastante grande indica que la variable es importante y en caso contrario para coeficientes muy pequeños

Variable	Tipo	Método	Score	p-valor
DiffWalk	Binaria	Chi <sup>2</sup>	1633.71	< 0,001
HighBP	Binaria	Chi <sup>2</sup>	1630.82	< 0,001
HeartDiseaseorAttack	Binaria	Chi <sup>2</sup>	1142.66	< 0,001
HighChol	Binaria	Chi <sup>2</sup>	987.08	< 0,001
Stroke	Binaria	Chi <sup>2</sup>	359.49	< 0,001
PhysActivity	Binaria	Chi <sup>2</sup>	144.35	< 0,001
HvyAlcoholConsump	Binaria	Chi <sup>2</sup>	112.38	< 0,001
Smoker	Binaria	Chi <sup>2</sup>	81.64	< 0,001
NoDocbcCost	Binaria	Chi <sup>2</sup>	38.41	< 0,001
Fruits	Binaria	Chi <sup>2</sup>	23.21	< 0,01
Sex	Binaria	Chi <sup>2</sup>	17.54	< 0,01
Veggies	Binaria	Chi <sup>2</sup>	17.42	< 0,01
CholCheck	Binaria	Chi <sup>2</sup>	7.22	< 0,05
AnyHealthcare	Binaria	Chi <sup>2</sup>	0.57	0.75
GenHlth	Ordinal	ANOVA F-test	1893.89	< 0,001
Age	Ordinal	ANOVA F-test	665.71	< 0,001
Income	Ordinal	ANOVA F-test	551.17	< 0,001
Education	Ordinal	ANOVA F-test	334.62	< 0,001
BMI	Continua	ANOVA F-test	1167.94	< 0,001
PhysHlth	Continua	ANOVA F-test	593.23	< 0,001
MentHlth	Continua	ANOVA F-test	108.15	< 0,001

Cuadro VII: Resultados del análisis individual de variables según el método estadístico aplicado.

Del análisis se concluye que la mayoría de las variables presentan alta capacidad discriminativa frente a la clase objetivo. Sin embargo, variables como AnyHealthcare, CholCheck, Fruits, Veggies, Sex, así como Education y MentHlth, muestran menor relevancia estadística y se consideran candidatas a ser eliminadas en la fase de reducción de dimensión.

##### 2. Regresión Logística

En el caso de la regresión logística inicialmente se realiza la selección de variables donde obtenemos los siguientes hallazgos:

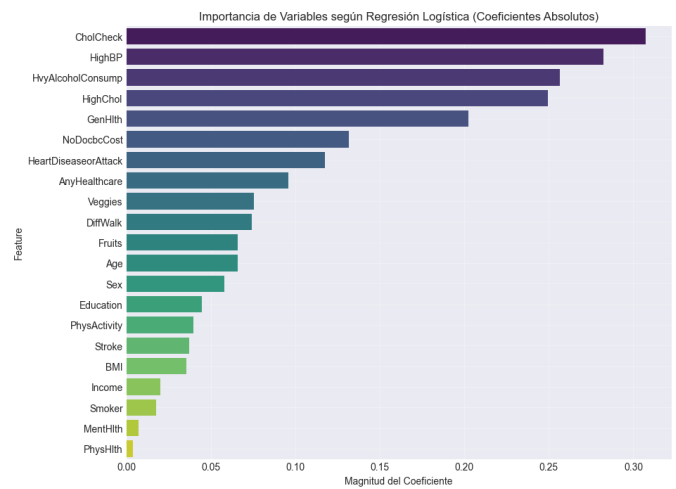


Figura 5: Selección de variables.

posteriormente procedemos a analizar las distribuciones generadas con por los métodos PCA y UMAP donde podemos

observar como estan las clases distribuidas despues de ser reorganizadas y llevadas aun nuevo espacio:

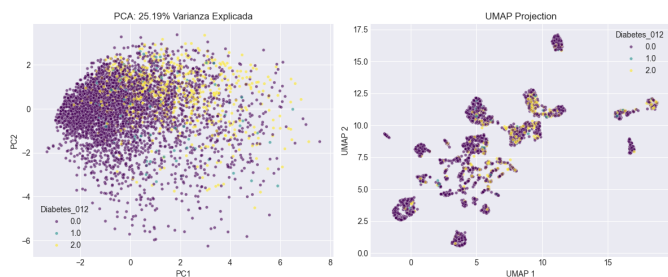


Figura 6: Muestras en PCA Y UMAP.

y obtuvimos los siguientes resultados:

Método	N_Features	% Reducción	Balanced Accuracy (Test)
Original	21	0.0 %	0.4943
PCA	19	9.5 %	0.5023
UMAP	2	90.5 %	0.3807

Cuadro VIII: Comparación de Desempeño Tras Reducción de Dimensión.

Podemos observar que tenemos una ganancia de en PCA donde se logró mejorar el rendimiento de la Regresión Logística, subiendo la precisión balanceada por casi un 1 por ciento, sin embargo, podemos ver que UMAP tuvo una perdida significativa cayendo por debajo del 40 por ciento.

UMAP es una técnica no lineal que agrupa los datos basándose en topología y vecindarios. Los grupos resultantes en 2D pueden tener formas complejas, curvas o islas

### 3. Bosque aleatorio equilibrado

En el caso de del Random Forest el proceso fue analogo al de la regresión logistica, sin embargo, los resultados son diferentes en torno a la seleccion de variables y los valores obtenidos por la reevaluación del modelo.

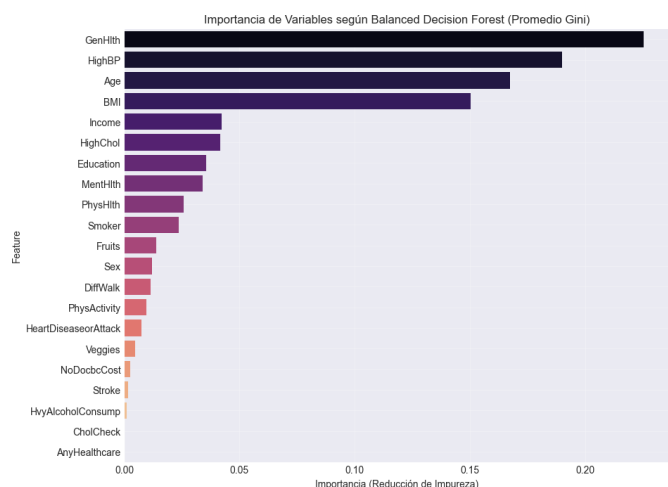


Figura 7: Selección de variables.

y obtuvimos los siguientes resultados:

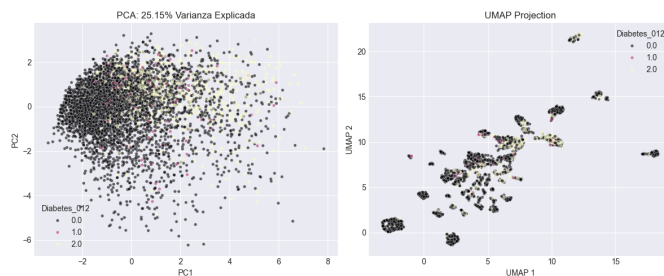


Figura 8: Muestras en PCA Y UMAP.

Método	N_Features	% Reducción	Balanced Accuracy (Test)
Original	21	0.0 %	0.4709
PCA	19	9.5 %	0.4907
UMAP	2	90.5 %	0.4595

Cuadro IX: Comparación de Desempeño Tras Reducción de Dimensión.

PCA solo pudo eliminar 2 variables (bajando de 21 a 19). Esto indica que tus datos originales tienen baja redundancia lineal; es decir, casi todas las variables aportan información única y no están fuertemente correlacionadas entre sí ademas tuvimos una ligera ganancia.

en el caso de UMAP pesar de esta reducción masiva, el rendimiento solo cayó muy ligeramente (de 0,4710 a 0,4596).

Aunque pierde frente a PCA en precisión pura, UMAP demuestra que el problema tiene una estructura no lineal fuerte que puede ser visualizada en un plano 2D sin perder mucha información crítica.

### REFERENCIAS

- [1] Chowdhury, M. M., Ayon, R. S., & Hossain, M. S. (2024). *An investigation of machine learning algorithms and data augmentation techniques for diabetes diagnosis using class imbalanced BRFSS dataset*. *Healthcare Analytics*, 5, 100297. ISSN: 2772-4425. Disponible en: doi.org/10.1016/j.health.2023.100297
- [2] Mohamed, M. H., Kamel, M. H. K., Said, W., & Mohamed, N. (2024). *Diabetic Mellitus Prediction with BRFSS Data Sets*. *Journal of Theoretical and Applied Information Technology*, 102(3), 15th February 2024. Little Lion Scientific. ISSN: 1992-8645, E-ISSN: 1817-3195. Disponible en: www.jatit.org/volumes/Vol102No3/10Vol102No3.pdf
- [3] Kaggle. (2015). *Diabetes Health Indicators Dataset (BRFSS 2015)*. Recuperado de: www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset
- [4] G. E. Melo-Acosta, F. Duitama-Muñoz and J. D. Arias-Londoño, "Fraud detection in big data using supervised and semi-supervised learning techniques," *2017 IEEE Colombian Conference on Communications and Computing (COLCOM)*, Cartagena, Colombia, 2017, pp. 1–6, doi: 10.1109/ColComCon.2017.8088206.