

Misael Redrejo Fernández, Yerson Caleb Yarhui Sarate



La toma de decisiones en la **aprobación de préstamos** constituye un desafío crítico en el ámbito financiero. La capacidad de evaluar con precisión la viabilidad de un solicitante para reembolsar un préstamo se ha vuelto cada vez más crucial para mitigar riesgos y garantizar la estabilidad financiera.

Conjunto de Datos

- Atributos numéricos = no_of_depends, income_annum, loan_amount, loan_term, cibil_score, residential_assets_value, commercial_assets_value, luxury_assets_value, bank_assets_value
- Atributos categóricos binarios = loan_status(salida), education, self_employed.

El gráfico de barras apiladas muestra la distribución de los préstamos aprobados y rechazados en función de los años para solventarlos. El eje horizontal (X) representa los años (2, 4, 6, 8, 10, 12, 14, 16, 18, 20) y el eje vertical (Y) representa el número de préstamos, con una escala de 10^9 . Las barras verdes representan los préstamos aprobados y las barras rojas representan los préstamos rechazados.

Años para solventar el préstamo	Préstamos Aprobados (Aproximado)	Préstamos Rechazados (Aproximado)
2	3.6	1.2
4	4.5	1.1
6	0.9	3.2
8	0.5	2.6
10	0.4	3.3
12	2.0	2.6
14	1.1	2.6
16	1.0	2.6
18	1.6	2.4
20	0.9	2.7

Category	Percentage	Number of Applicants
Approved	62.22%	2656
Rejected	37.78%	1613

Box plot showing the distribution of credit bureau score (cibil_score) for two loan status categories: Aprobado and Rechazado.

The y-axis represents the cibil_score, ranging from 300 to 900. The x-axis represents the loan status.

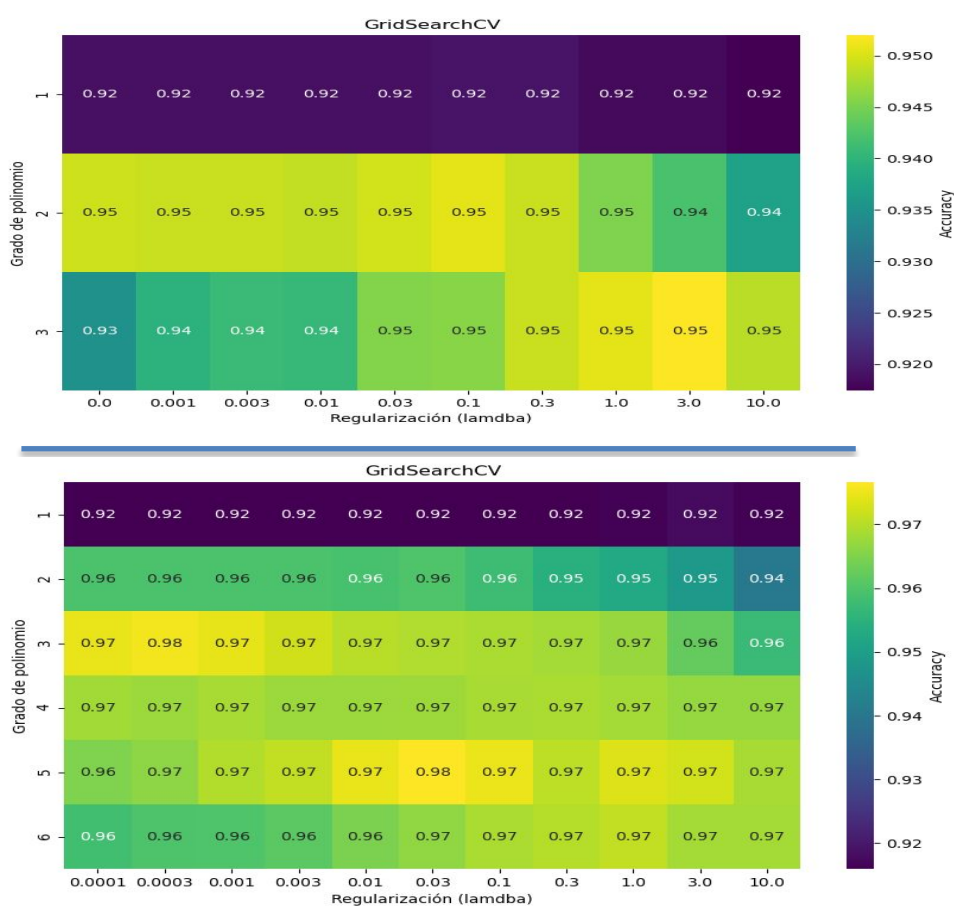
For the Aprobado group, the median score is approximately 710. The interquartile range (IQR) is from approximately 620 to 800. The whiskers extend from approximately 340 to 900. There are several outliers below the lower whisker, around 300-320.

For the Rechazado group, the median score is approximately 430. The IQR is from approximately 380 to 490. The whiskers extend from approximately 300 to 680. There are many outliers above the upper whisker, ranging from 680 to 900.

UPNA

La investigación se enfoca en comparar técnicas de aprendizaje automático, incluyendo regresión logística, RFE, Naive Bayes, redes neuronales, K-means, ensambles (Bagging, Decision Tree, Random Forest, Boosting), OVA y OVO. Se evalúan en dimensiones críticas del conjunto de datos, considerando la cantidad y relevancia de las características, así como el contexto socioeconómico de los solicitantes de préstamos.

Para la regresión logística se probaron utilizando diferentes grados regularización y variables polinomiales, utilizando la regresión normal y la regresión con RFE (5 atributos).

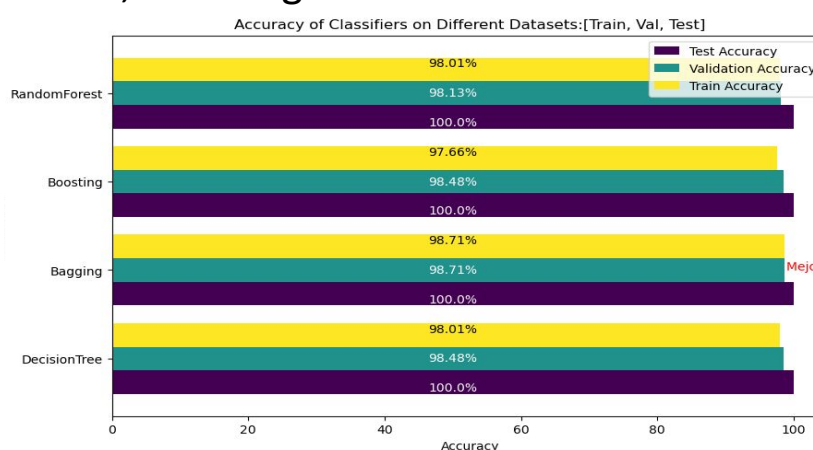


Con Naive Bayes se probaron diferentes técnicas de cálculo de probabilidad para variables categóricas y numéricas.

Con Redes Neuronales las siguientes combinaciones de capas ocultas: (11,),(50,),(11,50),(50,50); activación: logistic, relu, tanh; alpha: 0.001, 0.01, 0.1. Consiguiendo más éxito con capas ocultas (11, 50),activación relu y alpha 0.001.

K-Means, se realizó con un número fijo clusters de $K=2$ puesto que solo habían 2 posibles clases.

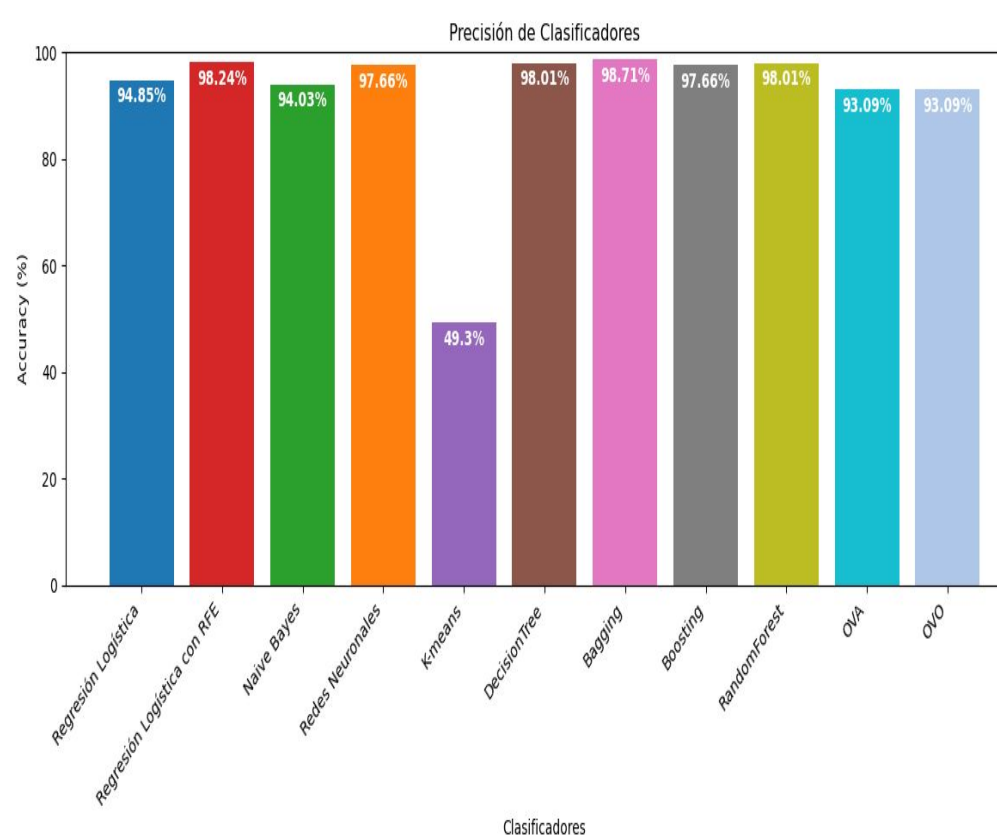
Con los Ensembles se evaluaron los clasificadores Bagging, Decision Tree, Random Forest, Boosting



Con OVA y con OVO se evaluaron los diferentes tipos de Kernel posibles linear, poly, rbf y sigmoid.

Nuestro análisis concluye que el modelo basado en la técnica de Bagging logra el mayor accuracy en los datos de prueba con un 98.71%. Esto significa que ha demostrado ser el más robusto y efectivo para generalizar, es decir, en realizar predicciones precisas en datos no vistos, lo cual es idóneo para aplicaciones del mundo real.

En la siguiente gráfica se muestran el accuracy obtenido en el conjunto de pruebas para los diferentes modelos entrenados.



En la selección de los mejores hiperparámetros para los diferentes modelos de aprendizaje automático se ha tenido en cuenta el accuracy en los datos de validación.

Al seleccionar el mejor modelo, observamos el rendimiento en el conjunto de prueba para obtener el modelo que mejor clasifique los datos.

- o Exploración de técnicas de normalización.
- o Mitigación de falsos positivos o falsos negativos: depende de los objetivos y riesgos específicos del negocio.
- o Optimización de hiper parámetros específicos: utilizar Google Colab para una búsqueda más exhaustiva.
- o Exploración de nuevas técnicas: es posible que existan técnicas no estudiadas que permitan ajustarse a los datos mejorando la generalización del modelo.

Yerson Caleb Yarhui Sarate
Misael Redrejo Fernandez