

TD VI: Inteligencia Artificial

Trabajo Práctico I

Universidad Torcuato Di Tella

Participantes:

Catalina Brusco, Catalina Chab López y Belén Chen

Fecha: Abril del 2025

1. Introducción al problema: origen y variables principales del conjunto de datos

El conjunto de datos elegido proviene de Kaggle y se utiliza para predecir si un préstamo solicitado por una persona será aprobado o rechazado, basándose en distintas características del solicitante. El conjunto fue enriquecido con variables adicionales basadas en *Riesgo Financiero* para la Aprobación de Préstamos. Además, se aplicó *SMOTENC* (*Synthetic Minority Over-sampling Technique for Nominal and Continuous*) para generar nuevos puntos de datos y ampliar el conjunto de instancias. El conjunto de datos contiene 45,000 registros y 14 variables. A continuación se describen las variables clave:

Table 1: Resumen de los atributos del dataset

Columna	Descripción	Tipo.de.Dato
person_age	Edad de la persona	Float
person_gender	Género de la persona	Catégorico
person_education	Nivel de educación más alto alcanzado	Catégorico
person_income	Ingreso anual de la persona	Float
person_emp_exp	Años de experiencia laboral	Entero
person_home_ownership	Estado de propiedad de la vivienda (ej., alquiler, propia)	Catégorico
loan_amnt	Monto del préstamo solicitado	Float
loan_intent	Propósito del préstamo	Catégorico
loan_int_rate	Tasa de interés del préstamo	Float
loan_percent_income	Monto del préstamo como porcentaje del ingreso anual	Float
cb_person_cred_hist_length	Longitud del historial de crédito en años	Float
credit_score	Puntaje de crédito de la persona	Entero
previous_loan_defaults_on_file	Indicador de los anteriores incumplimientos de préstamo	Catégorico
loan_status (objetivo)	Estado del préstamo: 1 = aprobado; 0 = rechazado	Entero (Binario)

Utilizar un árbol de decisión para modelar este problema es adecuado por varias razones. Primero, permite realizar tareas de clasificación binaria, como predecir si un préstamo será aprobado o rechazado, ya que divide los datos en función de las características más importantes. Segundo, maneja eficientemente datos mixtos, es decir, tanto numéricos (como ingresos y puntajes de crédito) como catégoricos (como género y estado civil), sin necesidad de transformaciones complicadas. Tercero, la estructura de un árbol de decisión permite dividir los datos en subgrupos homogéneos, lo que ayuda a identificar clientes con características similares que pueden tener un comportamiento parecido, como los que tienen un puntaje de crédito bajo o altos ingresos. Esto es importante porque las relaciones entre características, como ingresos y puntaje de crédito, no siempre son lineales, y los árboles pueden manejar estas interacciones. Finalmente, su gran ventaja es la interpretabilidad; los árboles de decisión son fáciles de entender y explican claramente las razones detrás de cada clasificación, lo cual es fundamental en el ámbito bancario, donde se necesitan justificar las decisiones de aprobación o rechazo de préstamos.

2. Preparación de los datos

Es necesario preprocesar los datos, para lo cual realizaremos un análisis que incluya la verificación de las variables y su clasificación, la normalización o escalado de las variables, la detección de valores faltantes, la identificación de posibles anomalías y balanceado de datos.

Como se mencionó anteriormente, las variables están correctamente clasificadas según su tipo (catégorico o numérico). Dado que vamos a utilizar árboles de decisión, no es necesario normalizar ni escalar los datos, ya que este modelo no depende de las magnitudes de las variables. Los árboles de decisión dividen los datos basándose en los valores específicos de las características, por lo que la escala de las variables no influye en

el desempeño del modelo. Asimismo, detallaremos que no existen valores faltantes en el conjunto de datos, como se puede observar en la Tabla 2.

Table 2: Resumen de valores nulos

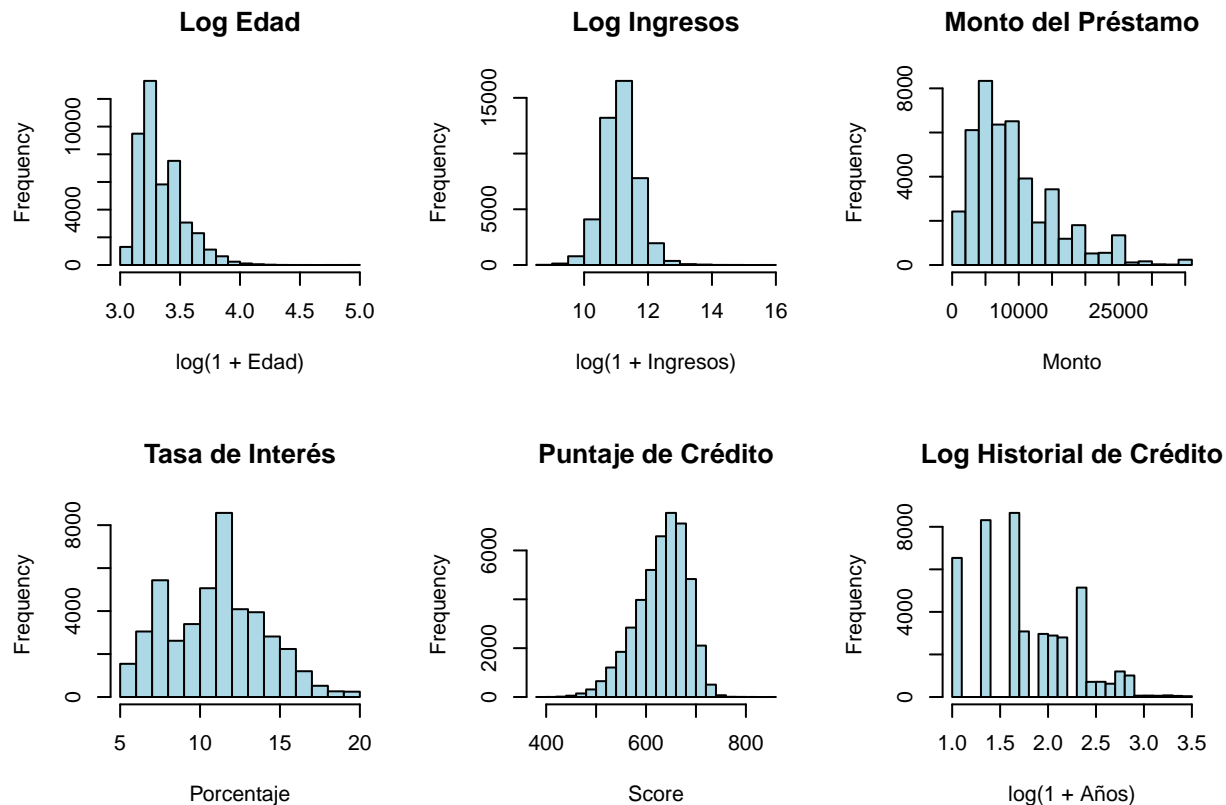
Cantidad.de.atributos.verificados	Atributos.con.0.nulos
14	14

También corroboramos que no hubiese una cantidad excesiva de filas repetidas, de hecho, nos dió cero.

Habiendo mencionado lo anterior, continuaremos con el análisis de las distribuciones de los distintos atributos. Como se puede observar en la Figura 1, se aplicó una transformación logarítmica a tres atributos clave: Edad, Ingresos e Historial de Crédito. Esta transformación es relevante porque, al ser creciente, no altera la distribución subyacente de los datos, pero sí comprime el rango de los valores grandes y expande el de los valores pequeños, lo que facilita su análisis.

Es importante destacar la relevancia de analizar la distribución de cada atributo. En este caso, observamos que atributos como Edad, Monto del Préstamo e Ingresos presentan una distribución sesgada a la derecha, mientras que el Puntaje de Crédito está sesgado a la izquierda. Esta asimetría podría generar predicciones erróneas para los valores extremos en cada uno de estos atributos. Además, notamos que el Historial de Crédito no sigue una distribución bien definida, lo que añade complejidad a su análisis.

Distribución de Variables Numéricas



Distribución de Variables Categóricas

Luego de analizar la distribución de las variables categóricas, no se identifican patrones relevantes ni desequilibrios significativos que requieran mención.

Anomalías

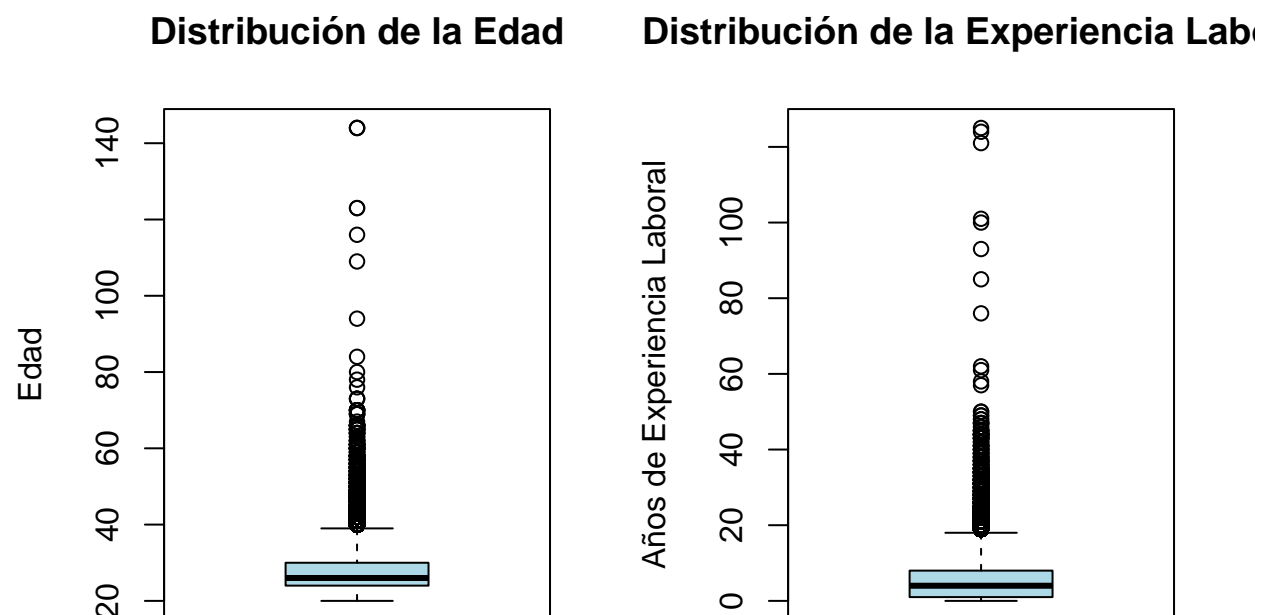


Figure 1: Boxplots de Experiencia Laboral y Edad

En el análisis preliminar de este conjunto de datos, se han identificado algunos valores atípicos (outliers) en los atributos de “Edad” (como personas con más de 120 años) y “Experiencia Laboral” (casos donde había más de 100 años de experiencia laboral). Sin embargo, su cantidad es chica (7 casos en un total de 45,000), lo que sugiere que no tendrán un impacto significativo en el entrenamiento del modelo. También, la distribución de la media en ambos Boxplots es bastante concentrada, lo que indica que la presencia de estos outliers no afecta la distribución general de los datos. En resumen, debido a que su número es tan pequeño, su impacto es prácticamente nulo, por lo que no es necesario eliminar estas filas.

Como resultado de los análisis realizados, no se eliminará ninguna información del conjunto de datos para evitar introducir sesgos en el proceso.

Chequeo de datos balanceados

Para evaluar el balance de nuestro conjunto de datos, analizamos las proporciones de las categorías de la variable de interés, en este caso, el estado del préstamo (aceptado o rechazado). A continuación, creamos un gráfico de barras que muestra la proporción de “Aceptado” y “Rechazado” en el conjunto de datos.

Al examinar el gráfico y las proporciones obtenidas, observamos que aproximadamente el 22% de los registros corresponden a “Aceptados” y el 77% a “Rechazados”. Esto indica que el conjunto de datos no está completamente balanceado, aunque la desproporción no es extremadamente alta. La desproporción en el conjunto de datos podría generar un sesgo hacia la clase mayoritaria (“Rechazados”), lo que podría resultar en un mejor desempeño del modelo para predecir esta clase y un desempeño inferior para la clase minoritaria (“Aceptados”).

Sin embargo, decidimos no intervenir en el conjunto de datos eliminando registros o utilizando técnicas de balanceo como sobremuestreo o submuestreo. Esto se debe a que tales técnicas podrían introducir sesgos

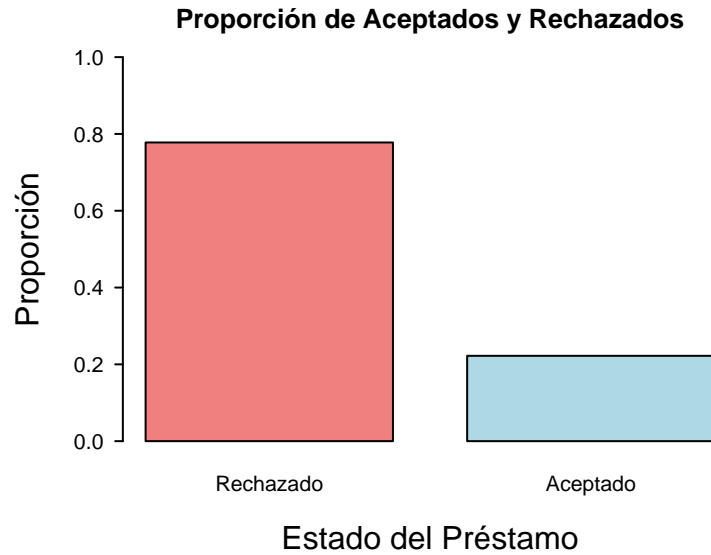


Figure 2: Proporción de Aceptados y Rechazados

adicionales o alterar la representatividad del conjunto de datos. En cambio, optamos por mantener la estructura original y tener en cuenta las proporciones de las clases al entrenar, validar y evaluar el modelo.

3. Construcción de árbol de decisión básico

Al realizar el análisis, mantuvimos la proporción de otorgamiento de préstamos (22%) y rechazo (78%) a lo largo de las etapas de entrenamiento, validación y testeo. En la descripción del árbol de decisión, observamos que cada nodo se compone de tres filas:

- 1.La primera fila indica la categoría de “no” (rechazo del préstamo).
- 2.La segunda fila muestra la proporción de casos dentro del subgrupo.
- 3.La tercera fila indica cuántos datos se concentran en ese nivel del árbol.

El árbol de decisión muestra que el primer factor que se utiliza para determinar si se le otorga un préstamo es la existencia de defaults previos (`previous_loan_defaults_on_file`), donde la presencia de antecedentes implica una altísima probabilidad de no pago. En ausencia de defaults, el siguiente factor relevante es la relación entre el monto solicitado y el ingreso (`loan_percent_income`), siendo los valores altos indicativos de mayor riesgo. Además, la tasa de interés del préstamo (`loan_int_rate`) cumple un rol clave: tasas superiores al 14% incrementan significativamente la probabilidad de incumplimiento. El nivel de ingresos (`person_income`) también influye, ya que a menor ingreso, el riesgo aumenta, especialmente combinado con tasas altas. La intención del préstamo (`loan_intent`), cuando es para educación, fines personales o emprendimientos, aparece asociada a mayor riesgo. Finalmente, la tenencia de vivienda propia o bajo hipoteca (`person_home_ownership`) contribuye a elevar la probabilidad de default en combinación con tasas altas y un elevado porcentaje préstamo/ingreso. En conjunto, estas variables permiten identificar perfiles de alto riesgo de manera clara.

Como conclusión, utilizando esta instancia del modelo entrenado con el conjunto de datos train se le rechazaría al 83% el préstamo y al resto se le otorgaría.

Hiperprámetros por defecto del árbol

A continuación se detallan los hiperparámetros utilizados por defecto en la construcción del árbol y el impacto que tiene cada uno en la estructura del modelo:

Árbol de Decisión

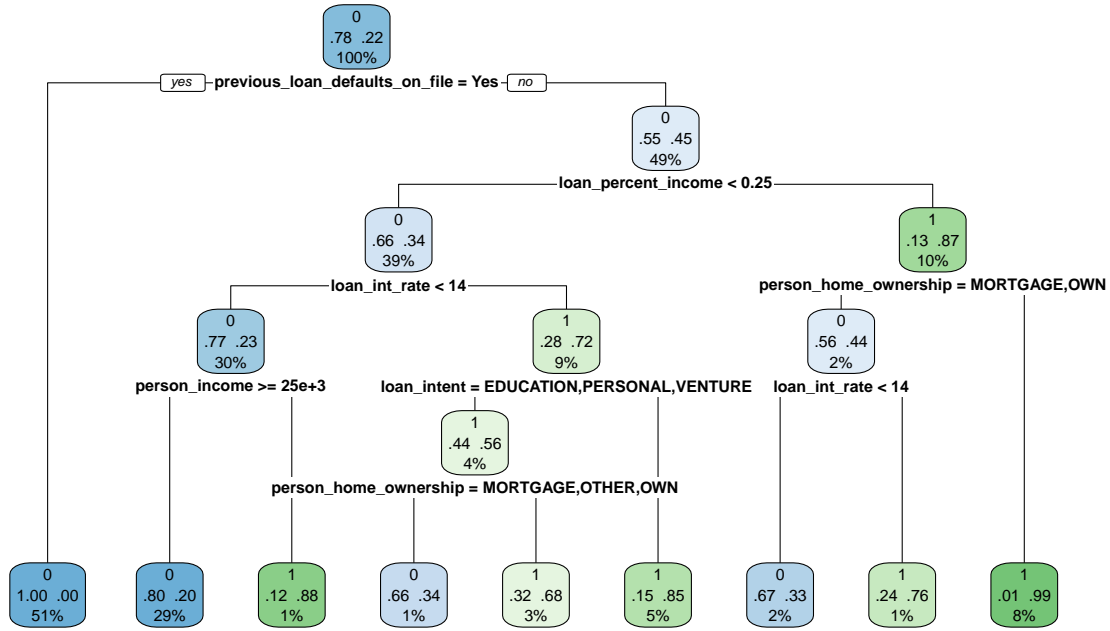


Figure 3: Árbol de Decisión del Modelo

`minsplit = 20`: Define el tamaño mínimo de observaciones en un nodo para que el árbol considere realizar una partición. Un valor de 20 limita las divisiones a grupos con un tamaño suficientemente grande, evitando sobreajuste en nodos con pocos datos.

`minbucket = 7`: Indica el tamaño mínimo que pueden tener las hojas terminales. Esto asegura que cada hoja final contenga al menos 7 observaciones, promoviendo estabilidad en las predicciones.

`cp = 0.01`: Es el parámetro de complejidad que regula el proceso de poda. Solo se aceptan divisiones que logren mejorar la calidad del ajuste en al menos un 1%. Un valor de 0.01 representa un control intermedio, balanceando entre un árbol complejo y uno demasiado simple.

`maxcompete = 4`: Almacena hasta 4 splits “competidores” cercanos al mejor split en cada nodo. Esto es útil para analizar qué otras variables casi logran ser seleccionadas.

`maxsurrogate = 5`: Controla la cantidad máxima de variables sustitutas utilizadas cuando hay datos faltantes en la variable principal de división. La presencia de 5 surrogates mejora la capacidad del modelo para manejar datos incompletos. En nuestro caso, al no haber datos faltantes, este hiperparámetro pasa a ser irrelevante.

`usesurrogate = 2`: Define cómo se utilizan los surrogates. Con un valor de 2, el árbol emplea surrogate splits incluso si la variable principal está disponible, siempre que aporten mejora en el ajuste.

`surrogatestyle = 0`: Indica que la selección de variables sustitutas se realiza utilizando un índice de concordancia simple, priorizando velocidad y simplicidad.

`maxdepth = 30`: Fija la profundidad máxima que puede alcanzar el árbol. Un valor alto de 30 permite que, si los datos y el pruning lo permiten, el árbol tenga una gran profundidad.

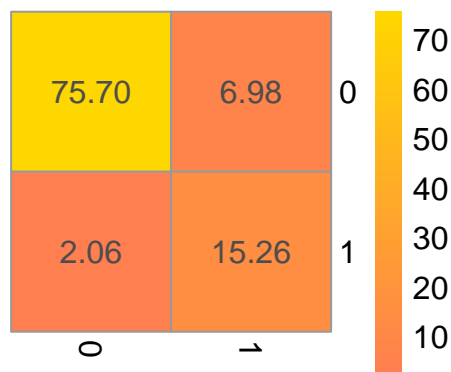
`xval = 10`: Determina que la validación cruzada interna para el proceso de poda se realice en 10 particiones, aportando mayor robustez a la selección del tamaño óptimo del árbol.

4. Evaluación del árbol de decisión básico

Se construye una matriz de confusión para analizar el rendimiento del modelo de clasificación. Inicialmente, se determina el número de predicciones positivas y negativas en el conjunto de prueba. A continuación, se define la función `gen_m_c`, que recibe como parámetros las etiquetas reales y las predicciones del modelo. Finalmente, los valores de la matriz se ajustan para representar la proporción de aciertos y errores, obteniendo así una versión normalizada de la matriz de confusión.

1. *TP* (Verdaderos Positivos): 75.7%
2. *TN* (Verdaderos Negativos): 15.26%
3. *FP* (Falsos Positivos): 6.98%
4. *FN* (Falsos Negativos): 2.06%

Confusion Matrix

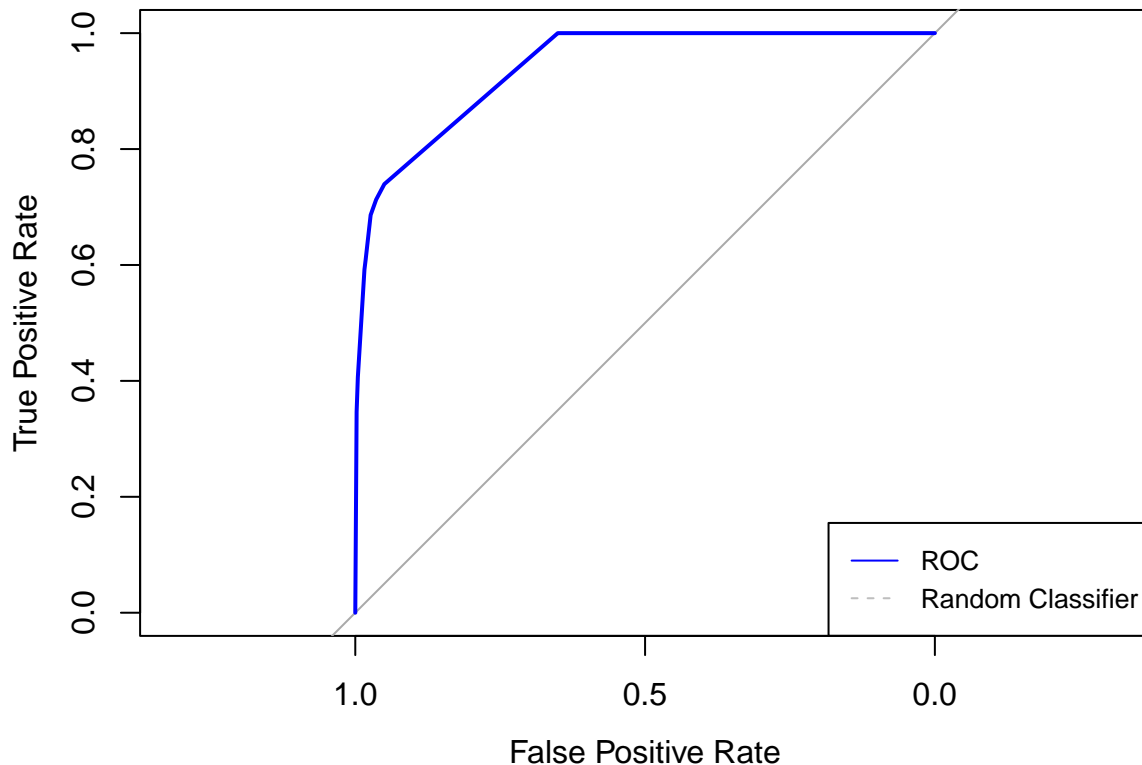


Los resultados obtenidos de las métricas de evaluación se detallan a continuación:

1. **Accuracy (Precisión Global):** 0.9096
 - Esto significa que el modelo clasificó correctamente el 90.96% de todas las instancias en el conjunto de datos (tanto positivas como negativas).
2. **Precision (Precisión):** 0.9156
 - Este valor indica que, de todas las predicciones que el modelo identificó como positivas (FP and TP), el 91.56% fueron efectivamente positivas. Como es un valor alto de precisión, esto significa que hay pocos falsos positivos.
3. **Recall (Sensibilidad o Exhaustividad):** 0.9735
 - Esto significa que el modelo fue capaz de identificar correctamente el 97.35% de todas las instancias positivas reales (FN y TP) en el conjunto de datos. Como es un valor alto de recall, esto significa que el modelo tiene una buena capacidad para identificar casos positivos, aunque pueda tener algunos falsos positivos.
4. **F1-score:** 0.9437
 - El F1-score es 0.9437 y es el promedio ponderado de la precisión y el recall. Este valor es especialmente útil cuando hay un desbalance en las clases, ya que captura tanto la exactitud como la exhaustividad del modelo. Como es un F1-score alto, el modelo tiene un buen balance entre precisión y recall. Si precision y recall tienen valores similares, significa que el modelo no favorece en exceso una métrica sobre la otra, es decir, no sacrifica la detección de positivos (recall) ni la precisión en sus predicciones.
5. **AUC-ROC:** 0.9415
 - El AUC-ROC (Área bajo la curva ROC) es 0.9415, lo que indica un rendimiento muy bueno del modelo (es cercano a 1) ya que la curva ROC compara la tasa de verdaderos positivos con la tasa de falsos positivos.

AUC-ROC: 0.9415207

Figura 1: Curva ROC – Árbol de Decisión



En términos de la curva ROC que obtenemos, efectivamente mostramos que nuestro modelo cuenta con alta sensibilidad (True Positive Rate) y especificidad (False Positive Rate), corroborando el hecho de que la curva azul se aproxime a la esquina de la parte superior izquierda, donde el modelo clasificaría correctamente a muchos de los positivos y a los negativos.

Si dicha curva ocupase valores más próximos a la diagonal gris, significaría que el modelo no se comporta mejor que el de una clasificación aleatoria. En el caso ideal, donde el modelo es perfecto, tendríamos la curva pegada a los ejes, de forma que obtendríamos una AUC de 1, es decir, la curva podría ir hacia el destino óptimo de manera lineal.

Además, la curva ROC en los datos de test indica que el modelo mantiene un buen desempeño fuera de los datos de entrenamiento, lo que sugiere que no está sobreajustado. Si el modelo estuviera fuertemente sobreajustado, esperaríamos una alta performance en entrenamiento pero un deterioro significativo en test. Dado que la curva en test sigue mostrando una buena separación entre clases, se puede inferir que el modelo generaliza correctamente. En la siguiente parte, veremos cómo para mejorar el rendimiento de este modelo se podrían buscar mejores hiperparámetros que los default.

5. Optimización del modelo

Mejores Hiperparámetros:

```
##      maxdepth_value minsplit_value minbucket_value    aucroc
## 106              17             15              21 0.9680951
```

Hiperparámetros default:

```
## $minsplit
## [1] 20
##
## $minbucket
```



```

## [1] 7
##
## $cp
## [1] 0.01
##
## $maxcompete
## [1] 4
##
## $maxsurrogate
## [1] 5
##
## $usesurrogate
## [1] 2
##
## $surrogatestyle
## [1] 0
##
## $maxdepth
## [1] 30
##
## $xval
## [1] 10
##
## Valores de AUC de arbol con hiperparametros default:
## Area under the curve: 0.9415
##
## Valores de AUC de arbol con mejores hiperparametros:
## Area under the curve: 0.9636

```