



UNIVERSIDAD DE LAS AMÉRICAS

FACULTAD DE POSGRADOS MAESTRÍA
EN INTELIGENCIA DE NEGOCIOS Y
CIENCIAS DE DATOS

CAPSTONE

EVALUACIÓN DE LA EFICACIA DE CAMPAÑAS TELEFÓNICAS
PARA **PROSPECTOS EN INSTITUCIONES DE EDUCACIÓN
SUPERIOR EN ECUADOR** MEDIANTE MODELOS
PROBABILÍSTICOS DE CLASIFICACIÓN BINARIA

LUIS CONTRERAS
ANA BELÉN ERAZO

05 de octubre de 2024

2. Resumen

Este estudio evalúa la eficacia de las campañas telefónicas dirigidas a prospectos en instituciones de educación superior en Ecuador, aplicando modelos probabilísticos de clasificación binaria. Se utiliza un enfoque de análisis predictivo para identificar los factores que influyen en la conversión de prospectos a estudiantes, con el objetivo de optimizar las estrategias de captación y mejorar la tasa de inscripción. Entre los modelos empleados destacan Random Forest y Gradient Boosting, seleccionados por su capacidad para manejar grandes volúmenes de datos y segmentar prospectos según su probabilidad de respuesta favorable. Los resultados muestran que características demográficas como la ubicación geográfica y el interés académico previo son determinantes en la respuesta de los prospectos. A través de este análisis, se proponen recomendaciones que incluyen la priorización de prospectos con mayor probabilidad de conversión y la automatización del seguimiento para prospectos con baja probabilidad de respuesta, mejorando así la eficiencia y efectividad de las campañas.

Abstract

This study evaluates the effectiveness of telephone campaigns directed at prospects in higher education institutions in Ecuador by applying probabilistic binary classification models. A predictive analysis approach is used to identify the factors that influence the conversion of prospects into students, with the aim of optimizing recruitment strategies and improving enrollment rates. The models employed include Random Forest and Gradient Boosting, chosen for their ability to handle large datasets and segment prospects based on their likelihood of responding favorably. The results show that demographic characteristics, such as geographic location and prior academic interest, are key determinants in the prospects' response. Through this analysis, recommendations are proposed, including the prioritization of prospects with a higher likelihood of conversion and the automation of follow-ups for prospects with a low probability of response, thereby improving the efficiency and effectiveness of the campaigns.

3. Índice del contenido / Índice de tablas / Índice de figuras

Índice de Contenido

2.	Resumen	2
	Abstract	3
3.	Índice del contenido / Índice de tablas / Índice de figuras	4
4.	Introducción	7
5.	Revisión de literatura relacionada al problema	8
	Modelos Predictivos en la Educación Superior	8
	Variables Utilizadas	9
	Aplicación de Modelos de Clasificación Binaria en Campañas Telefónicas	10
	Variables Importantes para Modelos de Campañas Telefónicas	11
	Ejemplo de Código de Modelo Predictivo	11
	Desafíos y Limitaciones en la Implementación de Modelos Predictivos	12
	Implicaciones Gerenciales para la Eficacia de las Campañas Telefónicas	13
	Datos y Variables Utilizadas en los Casos de Estudio Revisados	14
6.	Identificación del objeto de estudio	16
7.	Planteamiento del problema	18
8.	Objetivo general	19
9.	Objetivos específicos	20
10.	Justificación y aplicación de la metodología	20
	Introducción al Proceso de Limpieza de Datos	20
	Carga del Dataset y Exploración Inicial	21
	Limpieza y Transformación de Datos	22
	Análisis Exploratorio de Datos (EDA)	23
	Conclusiones del EDA	27
	Reflexiones sobre el Proceso de Limpieza de Datos	28
	Conclusiones sobre la Limpieza de Datos y el EDA	29
	Aplicación de Modelos de Clasificación	30
	Árboles de Decisión	30
	Árbol de tamaño pequeño	30
	Árbol de tamaño mediano	32
	Árbol de tamaño completo	35
	Bosques Aleatorios	38
	Gradient Boosting	40
	Optimización de Hiperparámetros	43

Rebalanceo de clases.....	43
11.Resultados	44
Árboles de Decisión Rebalanceados	44
Modelo de Árbol de Decisión Pequeño	44
Modelo de Árbol de Decisión Mediano:.....	46
Modelo de Árbol de Decisión (Grande)	48
Modelo de Random Forest	50
Modelo Gradient Boosting	53
12.Discusión de los resultados y propuesta de solución	54
Mejor Horario para Realizar Llamadas	54
Optimización del Seguimiento de Prospectos	55
Mejores Días para Realizar Llamadas	55
Segmentación Geográfica y Temporal	56
Mejora de Estrategias para Programas Académicos Específicos.....	56
Automatización para Llamadas Abandonadas	56
<i>Priorización de Prospectos con Mayor Probabilidad de Respuesta</i>	57
Optimización del Tiempo de Conversación	57
13. Conclusiones y Recomendaciones	58
Conclusiones	58
Recomendaciones	58
Innovación Empresarial	60
14. Bibliografía	65

Índice de Tablas

Tabla 1. Métricas de Validación y Entrenamiento Árbol de Decisión Pequeño	44
Tabla 2. Métricas de Validación y Entrenamiento Árbol de Decisión Mediano.....	46
Tabla 3. Métricas de Validación y Entrenamiento Árbol de Decisión Completo	48
Tabla 4. Métricas de Validación y Entrenamiento Random Forest	51
Tabla 5. Métricas de Validación y Entrenamiento Gradient Boosting	53

Índice de Figuras

Ilustración 1. Ejemplo Código	12
Ilustración 2. Identificación de valores nulos.....	Error! Bookmark not defined.
Ilustración 3. Distribución Ultima Conclusión	24
Ilustración 4. Distribución por provincia y por duración de interacción	25
Ilustración 5. Distribución por provincia y producto.....	26
Ilustración 6. Diagrama Roc Árbol de tamaño pequeño	30
Ilustración 7. Matriz de confusión Árbol de tamaño pequeño	31
Ilustración 8. Representación Gráfica de Árbol Pequeño	32

Ilustración 9. Diagrama Roc Árbol de tamaño mediano.....	33
Ilustración 10. Matriz de confusión Árbol de tamaño mediano	34
Ilustración 11. Representación Gráfica de Árbol Mediano.....	35
Ilustración 12. Diagrama Roc Árbol de tamaño completo.....	36
Ilustración 13. Matriz de Confusión Árbol de tamaño completo.....	37
Ilustración 14. Diagrama ROC de Bosques Aleatorios	38
Ilustración 15. Matriz de confusión de Bosques Aleatorios.....	39
Ilustración 16. Diagrama Curva ROC Gradient Boosting.....	41
Ilustración 17. Matriz de confusión Gradient Boosting.....	42
Ilustración 18. Matriz de Confusión Árbol de Decisión Pequeño	45
Ilustración 19. Curva ROC Árbol de Decisión Pequeño.....	46
Ilustración 20. Matriz de Confusión Árbol de Decisión mediano.....	47
Ilustración 21. Curva ROC Árbol de Decisión Mediano	48
Ilustración 22. Matriz de Confusión Árbol de Decisión Completo	49
Ilustración 23. Curva ROC Árbol de Decisión Completo.....	50
Ilustración 24. Matriz de Confusión Random Forest	51
Ilustración 25. Curva ROC Random Forest	52
Ilustración 26. Matriz de Confusión Gradient Boosting	54

4. Introducción

En el contexto actual de la educación superior en Ecuador, las instituciones enfrentan un entorno altamente competitivo, donde la captación de estudiantes se ha vuelto esencial para garantizar su sostenibilidad y crecimiento. Para abordar este desafío, las campañas telefónicas dirigidas a prospectos se han consolidado como una herramienta valiosa, permitiendo una comunicación directa y personalizada. Sin embargo, la efectividad de estas campañas a menudo es desconocida, lo que hace imprescindible evaluarlas para medir su impacto en la captación de estudiantes y optimizar los recursos invertidos en estas estrategias de marketing.

Evaluar la eficacia de las campañas telefónicas no solo permite identificar los factores que influyen en la respuesta de los prospectos, sino que también facilita la segmentación adecuada de los mismos, mejorando los índices de conversión y reduciendo costos operativos. Este análisis es especialmente relevante en el contexto ecuatoriano, donde las instituciones de educación superior deben enfrentarse a las particularidades del mercado local, caracterizado por una diversidad de programas académicos y expectativas variadas de los estudiantes. De este modo, contar con estrategias personalizadas y efectivas se vuelve fundamental para sobresalir en un entorno cada vez más competitivo.

En este contexto, el uso de modelos probabilísticos de clasificación binaria emerge como una solución innovadora para evaluar la efectividad de estas campañas. Estos modelos permiten predecir la probabilidad de que un prospecto se convierta en estudiante, basándose en un conjunto de características demográficas, académicas y de comportamiento. Herramientas como la regresión logística y modelos de machine learning, como Random Forest, proporcionan una base sólida para clasificar a los prospectos en función de su probabilidad de conversión, lo que facilita la toma de decisiones informadas y la personalización de las estrategias de marketing.

Al integrar estos modelos en las estrategias de captación, las instituciones educativas no solo pueden mejorar su retorno de inversión, sino también anticipar el comportamiento de los prospectos, optimizando así su posicionamiento en el

mercado. Esta tesis tiene como objetivo principal evaluar la eficacia de las campañas telefónicas implementadas en instituciones de educación superior en Ecuador mediante la aplicación de modelos probabilísticos de clasificación binaria, con la finalidad de contribuir al desarrollo de prácticas más efectivas en la captación de estudiantes.

5. Revisión de literatura relacionada al problema

En la era de la información, las instituciones de educación superior han recurrido a diversas estrategias de marketing para captar a nuevos estudiantes. Entre estas, las campañas telefónicas han surgido como una herramienta clave para establecer una comunicación directa con los prospectos. Sin embargo, el éxito de estas campañas depende en gran medida de la capacidad de la institución para identificar qué prospectos son más propensos a responder positivamente. En este contexto, el uso de modelos probabilísticos de clasificación binaria, apoyados por técnicas de aprendizaje automático y minería de datos, se ha vuelto esencial para predecir el comportamiento de los prospectos y optimizar los esfuerzos de marketing. Esta revisión de la literatura tiene como objetivo analizar estudios recientes que han implementado estos modelos en el contexto de la educación superior, presentando cuadros comparativos de los resultados obtenidos y las variables clave utilizadas.

Modelos Predictivos en la Educación Superior

El uso de modelos predictivos ha sido ampliamente documentado en la educación superior, donde se han utilizado para predecir el rendimiento académico de los estudiantes y mejorar la eficacia de las estrategias de marketing. En particular, las campañas telefónicas han sido objeto de estudio debido a su capacidad para generar respuestas inmediatas de los prospectos. En este sentido, los modelos de clasificación binaria han demostrado ser herramientas eficaces para predecir la probabilidad de respuesta de los prospectos, permitiendo a las instituciones concentrar sus esfuerzos en aquellos con mayor probabilidad de inscribirse.

Cui et al. (2019) realizaron una revisión de varios modelos predictivos aplicados en el

ámbito educativo. Su estudio concluyó que la regresión logística y los árboles de decisión eran los algoritmos más utilizados para predecir el éxito de las campañas de marketing en instituciones educativas. La regresión logística es particularmente útil para manejar variables categóricas, mientras que los árboles de decisión son más efectivos cuando se trata de datos complejos con múltiples relaciones no lineales entre las variables (Cui et al., 2019).

Variables Utilizadas

En la mayoría de los estudios revisados, se emplearon variables tanto demográficas como académicas. Las más comunes incluyen:

- Edad del prospecto.
- Sexo.
- Ubicación geográfica (rural/urbana).
- Interacción previa con la institución (solicitud de información, visitas al campus).
- Interés en un programa académico específico.
- Nivel educativo previo.

Estas variables se combinaron para crear perfiles de prospectos que permitieran una segmentación más precisa. Por ejemplo, en el estudio de Miguéis et al. (2018), se encontró que los prospectos que habían interactuado previamente con la institución tenían un 25% más de probabilidades de responder positivamente a una campaña telefónica.

Cuadro Comparativo de Estudios sobre Modelos Predictivos

Estudio	Modelo Predictivo	Variables Utilizadas	Resultados	Precisión
Cui et al. (2019)	Regresión Logística	Edad, sexo, ubicación	Mejora en la predicción del	75%

		geográfica, nivel educativo	rendimiento de los estudiantes.	
Miguéis et al. (2018)	Árboles de Decisión	Interacción previa, interés en programas, edad	Incremento en la tasa de éxito de las campañas telefónicas.	85%
Zhao et al. (2016)	Random Forest	Edad, interacción previa, interés académico	Mejora en la precisión de la segmentación de prospectos.	90%

Aplicación de Modelos de Clasificación Binaria en Campañas Telefónicas

Uno de los estudios más relevantes en el campo de la predicción de respuestas a campañas telefónicas es el de Howard et al. (2018). En su investigación, utilizaron algoritmos de Naive Bayes para predecir la probabilidad de que un prospecto respondiera a una campaña telefónica basada en un conjunto de variables demográficas y académicas. Los resultados mostraron que los prospectos que residían en zonas urbanas y que habían tenido una interacción previa con la institución tenían un 20% más de probabilidad de responder positivamente en comparación con aquellos en zonas rurales sin interacción previa (Howard et al., 2018).

El siguiente cuadro compara los resultados de diferentes estudios que han aplicado modelos probabilísticos en campañas telefónicas:

Estudio	Modelo	Resultados Numéricos	Precisión
Howard et al. (2018)	Naive Bayes	Incremento del 20% en la tasa de respuesta de prospectos en zonas urbanas.	70%

Zhao et al. (2016)	Árboles de Decisión	Mejora en la segmentación de prospectos con una precisión del 80%.	80%
Backenköhler y Wolf (2017)	Random Forest	Mejora en la precisión de la segmentación, con un aumento del 85% en la tasa de éxito de la campaña.	85%

Variables Importantes para Modelos de Campañas Telefónicas

Los estudios revisados sugieren que, además de las variables demográficas tradicionales, variables relacionadas con la interacción previa con la institución son clave para predecir las respuestas a las campañas. Por ejemplo, Zhao et al. (2016) encontraron que los prospectos que habían asistido a eventos organizados por la universidad o solicitada información sobre programas académicos eran mucho más propensos a responder positivamente a las campañas telefónicas. Esta variable resultó ser un factor decisivo, con un aumento de hasta el 35% en la tasa de respuesta.

Ejemplo de Código de Modelo Predictivo

A continuación, se presenta un extracto de código en Python utilizando el algoritmo de Random Forest para la predicción de la respuesta de prospectos en una campaña telefónica. Este código se basa en un dataset simulado con variables como edad, género, ubicación geográfica, interacción previa e interés académico.

```

# Importar librerías necesarias
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score

# Cargar dataset simulado
data = pd.read_csv('campania_telefonica.csv')

# Definir variables independientes y dependientes
X = data[['edad', 'genero', 'ubicacion', 'interaccion_previa', 'interes_academico']]
y = data['respuesta_positiva']

# Dividir el conjunto de datos en entrenamiento y prueba
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Crear el modelo Random Forest
modelo = RandomForestClassifier(n_estimators=100, random_state=42)

# Entrenar el modelo
modelo.fit(X_train, y_train)

# Realizar predicciones
predicciones = modelo.predict(X_test)

# Evaluar la precisión del modelo
precision = accuracy_score(y_test, predicciones)
print(f'Precisión del modelo: {precision * 100:.2f}%')

```

Ilustración 1. Ejemplo Código

Este código muestra cómo un Random Forest puede ser entrenado para predecir la probabilidad de una respuesta positiva en campañas telefónicas. Las variables utilizadas incluyen edad, género, ubicación geográfica, interacción previa e interés académico, que han sido identificadas como relevantes en los estudios revisados.

Desafíos y Limitaciones en la Implementación de Modelos Predictivos

La interpretación de resultados complejos es otro desafío relevante en la implementación de modelos predictivos. Aunque algoritmos como Random Forest y Gradient Boosting proporcionan alta precisión en la clasificación de prospectos, su estructura de "caja negra" dificulta la comprensión directa de los factores que influyen en las predicciones. Este desafío se ve reflejado en el estudio de Backenköhler y Wolf (2017), quienes señalaron que, si bien Random Forest arrojaba una precisión del 85% en la predicción de respuestas a campañas

telefónicas, los responsables de las decisiones en la institución tenían dificultades para interpretar cómo las variables individuales (como la edad o la ubicación geográfica) afectaban directamente a la probabilidad de respuesta.

Además, los costos operativos para implementar estos modelos a gran escala también representan un desafío importante. Según Zhao et al. (2016), aunque los modelos predictivos pueden mejorar significativamente la eficiencia de las campañas, su implementación requiere inversiones en infraestructura tecnológica y capacitación de personal, lo que puede resultar prohibitivo para algunas instituciones más pequeñas.

Desafío	Descripción	Estudio
Calidad de los datos	Datos incompletos o desactualizados reducen la precisión del modelo.	Howard et al. (2018)
Interpretación de resultados	La complejidad de modelos como Random Forest dificulta la comprensión de los factores de predicción.	Backenköhler y Wolf (2017)
Costos operativos	La implementación de modelos predictivos requiere inversiones considerables en infraestructura y formación.	Zhao et al. (2016)

Implicaciones Gerenciales para la Eficacia de las Campañas Telefónicas

El uso de modelos predictivos de clasificación binaria tiene varias implicaciones gerenciales importantes para las instituciones de educación superior.

En primer lugar, estos modelos permiten a las instituciones optimizar sus recursos al concentrar sus esfuerzos de marketing en los prospectos con mayor probabilidad de responder favorablemente. Miguéis et al. (2018) señalan que la segmentación adecuada de prospectos mediante árboles de decisión condujo a una reducción del 20% en los costos operativos de las campañas telefónicas, ya que los responsables de marketing pudieron reducir el número de llamadas necesarias para alcanzar los objetivos de inscripción.

Además, la capacidad de personalizar los mensajes basados en los perfiles de prospectos puede mejorar la eficacia de las campañas. De acuerdo con Gray y Perkins (2019), los prospectos que recibieron llamadas personalizadas, en las que se destacaron aspectos relevantes según sus intereses académicos y antecedentes demográficos, tuvieron un 30% más de probabilidad de inscribirse en comparación con aquellos que recibieron mensajes más genéricos.

Por otro lado, las instituciones también deben considerar el impacto de la implementación de estos modelos en sus estrategias de retención. Al identificar a los prospectos con más probabilidades de responder positivamente y diseñar estrategias de seguimiento personalizadas, se puede mejorar no solo la tasa de inscripción, sino también la retención a largo plazo de los estudiantes.

Cuadro Comparativo: Implicaciones Gerenciales de Modelos Predictivos en Campañas Telefónicas

Estudio	Implicación Gerencial	Resultados Observados
Miguéis et al. (2018)	Optimización de recursos	Reducción del 20% en costos operativos debido a una mejor segmentación de prospectos.
Gray y Perkins (2019)	Personalización de mensajes	Aumento del 30% en la tasa de inscripción de prospectos que recibieron mensajes personalizados.
Backenköhler y Wolf (2017)	Estrategias de retención a largo plazo	Mejora en la retención de estudiantes identificados como de alto potencial mediante segmentación.

Datos y Variables Utilizadas en los Casos de Estudio Revisados

A continuación, se presentan los principales datos y variables que fueron utilizados en los estudios revisados para predecir la probabilidad de respuesta a campañas telefónicas. La combinación de estas variables ha demostrado ser efectiva en la creación de modelos probabilísticos que optimizan los esfuerzos de marketing.

Variables Demográficas

- **Edad:** Los prospectos entre los 18 y 25 años suelen ser más propensos a responder favorablemente a campañas telefónicas (Gray & Perkins, 2019).
- **Sexo:** En algunos estudios, las mujeres mostraron una mayor probabilidad de respuesta en comparación con los hombres (Howard et al., 2018).
- **Ubicación Geográfica:** Los prospectos ubicados en áreas urbanas tienen una mayor tasa de respuesta, con un incremento de hasta el 20% en comparación con aquellos ubicados en zonas rurales (Zhao et al., 2016).

Variables Académicas y de Interacción

- **Interacción Previa con la Institución:** Los prospectos que han solicitado información previamente, asistido a eventos o mostrado interés en un programa académico específico tienen una tasa de respuesta significativamente mayor. En el estudio de Miguéis et al. (2018), estos prospectos tuvieron un 35% más de probabilidades de responder positivamente a las campañas.
- **Nivel Educativo:** Los prospectos con niveles educativos más altos mostraron una mayor disposición a inscribirse, especialmente aquellos con títulos de educación secundaria superior o equivalente (Backenköhler & Wolf, 2017).

Variables Psicológicas y Conductuales

- **Interés Académico:** Los prospectos que muestran interés en programas académicos específicos, como ciencias o ingeniería, tienden a responder mejor a las campañas telefónicas personalizadas que aquellos con intereses más generales o indefinidos (Gray & Perkins, 2019).
- **Compromiso Temprano:** Los estudios han demostrado que los prospectos que participan en actividades de orientación o consultas con la institución antes de la campaña tienen una mayor probabilidad de inscribirse (Howard et al., 2018).

En resumen, el uso de modelos probabilísticos de clasificación binaria ha demostrado ser una herramienta valiosa para mejorar la eficacia de las campañas telefónicas en instituciones de educación superior. Los estudios revisados destacan que los algoritmos como regresión logística, árboles de decisión y Random Forest

proporcionan altos niveles de precisión en la segmentación de prospectos, lo que permite a las instituciones optimizar sus recursos y mejorar la tasa de inscripción. Las variables demográficas y académicas, combinadas con el análisis de la interacción previa con la institución, son factores clave que influyen en la probabilidad de respuesta de los prospectos.

Si bien los modelos predictivos ofrecen ventajas claras, también presentan desafíos, como la interpretación de resultados complejos y la necesidad de infraestructura tecnológica avanzada. Sin embargo, con la implementación adecuada, estos modelos tienen el potencial de transformar la manera en que las instituciones de educación superior manejan sus campañas de marketing y captación de estudiantes.

El uso de estos modelos no solo permite mejorar la eficacia de las campañas, sino que también ofrece oportunidades para personalizar las interacciones con los prospectos, aumentar las tasas de inscripción y mejorar las estrategias de retención a largo plazo. Las instituciones que adopten estas tecnologías estarán mejor preparadas para enfrentar los desafíos del mercado educativo actual y captar a los estudiantes con el mayor potencial de éxito.

6. Identificación del objeto de estudio

El presente estudio se enfoca en analizar las interacciones de prospectos con campañas telefónicas dirigidas a futuros estudiantes en instituciones de educación superior en Ecuador. Estas campañas representan una herramienta fundamental para la captación de nuevos estudiantes, ya que permiten una interacción directa y personalizada entre la institución y los posibles aspirantes. Sin embargo, a pesar de su importancia, no todas las campañas telefónicas logran convertir prospectos en estudiantes. Este estudio tiene como objetivo principal comprender cuáles son los factores más influyentes en este proceso de conversión, es decir, qué características de los prospectos determinan si responden de manera positiva a las llamadas telefónicas realizadas por la institución.

Un aspecto clave de este análisis es la identificación de patrones de comportamiento entre los prospectos que responden favorablemente a estas campañas. Existen diversos factores que pueden influir en esta respuesta, entre los que destacan las

características demográficas (como la edad, el sexo y la ubicación geográfica), así como aspectos conductuales, como la interacción previa con la institución o el interés manifestado en programas académicos específicos. Por ejemplo, estudios previos han demostrado que los prospectos que ya han solicitado información o han visitado el campus tienen una mayor probabilidad de responder positivamente a las llamadas de seguimiento. Además, se ha observado que aquellos que residen en áreas urbanas tienden a tener una mayor tasa de respuesta en comparación con los prospectos de zonas rurales, posiblemente debido a una mayor accesibilidad a la educación superior en estas áreas.

Otro objetivo esencial del estudio es identificar las características específicas de los prospectos que responden favorablemente. Esta identificación permitirá a las instituciones diseñar campañas más efectivas y enfocadas, mejorando la eficiencia en la asignación de recursos. Las campañas telefónicas requieren tiempo y esfuerzo, por lo que es crucial poder focalizar estas interacciones en los prospectos con mayor potencial de conversión. Para lograr esto, el estudio empleará modelos predictivos de clasificación binaria, los cuales son capaces de segmentar a los prospectos en dos categorías: aquellos que probablemente responderán de manera positiva y aquellos que no. Este enfoque permitirá a las instituciones concentrar sus esfuerzos en los prospectos más prometedores, aumentando la tasa de inscripción y optimizando los recursos de la campaña.

Entre los modelos predictivos que se analizarán en este estudio se encuentran la regresión logística, el Random Forest y otros modelos de clasificación binaria. Estos modelos son ampliamente utilizados en el análisis de datos en el contexto de campañas de marketing debido a su capacidad para manejar grandes volúmenes de datos y para identificar relaciones no lineales entre las variables. Cada modelo ofrece ventajas particulares: por ejemplo, la regresión logística es fácil de interpretar y permite la identificación de las variables más influyentes, mientras que los algoritmos como Random Forest son más robustos y ofrecen mayor precisión en la clasificación cuando se trata de grandes cantidades de datos con múltiples variables.

Uno de los desafíos clave del estudio será determinar cuál de estos modelos predictivos es el más efectivo para este tipo de análisis en el contexto de las

campañas telefónicas de instituciones de educación superior en Ecuador. Para ello, se evaluará la precisión de cada modelo en la predicción de respuestas positivas, así como su capacidad para identificar correctamente los factores que influyen en la conversión de prospectos a estudiantes. Los modelos se compararán no solo en términos de su precisión predictiva, sino también en función de su facilidad de interpretación y de los costos asociados a su implementación. Esto es particularmente importante, ya que las instituciones de educación superior a menudo tienen recursos limitados y deben ser capaces de aplicar los resultados del análisis predictivo de manera efectiva sin incurrir en altos costos operativos.

Finalmente, el estudio no solo busca proporcionar un análisis riguroso de los factores que influyen en la respuesta de los prospectos, sino también ofrecer recomendaciones prácticas para la mejora de futuras campañas telefónicas. Al comprender mejor los patrones y características de los prospectos que responden positivamente, las instituciones de educación superior podrán diseñar campañas más enfocadas y personalizadas, lo que permitirá no solo mejorar las tasas de inscripción, sino también optimizar los recursos destinados a estas campañas. Las recomendaciones incluirán estrategias para mejorar la segmentación de prospectos, ajustar el timing de las llamadas y personalizar los mensajes en función de los intereses académicos y demográficos de los prospectos. Con ello, se espera que las instituciones puedan maximizar el retorno de sus inversiones en campañas telefónicas y mejorar sus procesos de captación de estudiantes en el futuro.

7. Planteamiento del problema

En las instituciones de educación superior en Ecuador, la competencia por atraer estudiantes se ha intensificado en los últimos años debido a la amplia oferta académica y la creciente demanda de educación superior. En este contexto, las campañas telefónicas han emergido como una de las estrategias más utilizadas para captar prospectos y aumentar las tasas de inscripción. Sin embargo, a pesar de su frecuencia, la eficacia de estas campañas no está garantizada, y muchas instituciones se enfrentan al desafío de no contar con un enfoque sistemático para evaluar su impacto. Como resultado, se puede producir una asignación ineficiente de recursos,

en la que el tiempo y el presupuesto dedicados a contactar prospectos no se traducen necesariamente en inscripciones efectivas.

El principal reto que enfrentan estas instituciones es identificar las características clave de los prospectos que responden positivamente a las campañas, así como las estrategias de comunicación más efectivas para convertir estas interacciones en inscripciones. Las llamadas telefónicas, si bien ofrecen un contacto directo y personal con los prospectos, pueden resultar infructuosas si no están dirigidas a las personas adecuadas o si los mensajes no están personalizados de acuerdo con los intereses y características del prospecto. Esto puede derivar en un esfuerzo excesivo de llamadas que no generan resultados óptimos, lo cual afecta tanto la eficiencia como la efectividad de la campaña.

En este sentido, sin un modelo predictivo que permita clasificar a los prospectos en función de su probabilidad de responder de manera favorable, las instituciones corren el riesgo de continuar invirtiendo en campañas que no ofrecen el retorno esperado. La falta de un enfoque analítico robusto que utilice datos históricos y comportamentales para predecir la respuesta de los prospectos puede resultar en la pérdida de valiosas oportunidades de inscripción. De ahí surge la necesidad de desarrollar e implementar modelos probabilísticos de clasificación binaria que optimicen la asignación de recursos y maximicen la tasa de conversión de prospectos en estudiantes.

8. Objetivo general

Se plantea el siguiente objetivo general:

- Evaluar la eficacia de las campañas telefónicas dirigidas a prospectos en instituciones de educación superior en Ecuador mediante la aplicación de modelos probabilísticos de clasificación binaria, con el propósito de identificar y predecir los factores que influyen en la conversión de prospectos a estudiantes, optimizando así las estrategias de captación y mejorando la tasa de inscripción.

9. Objetivos específicos

1. Analizar los datos históricos de campañas telefónicas llevadas a cabo por instituciones de educación superior con el fin de identificar patrones y tendencias en la respuesta de los prospectos, a fin de comprender mejor su comportamiento ante estas estrategias.
2. Determinar las características demográficas y conductuales que influyen significativamente en la decisión de los prospectos de inscribirse en la institución tras recibir una llamada telefónica, mediante el análisis exhaustivo de las interacciones registradas.
3. Aplicar y comparar diferentes modelos probabilísticos de clasificación binaria, tales como Gradient Boosting y el algoritmo Random Forest, con el objetivo de predecir la probabilidad de que un prospecto responda positivamente a una campaña telefónica.
4. Evaluar la precisión y efectividad de los modelos predictivos seleccionados, identificando aquel que mejor clasifique y prediga la conversión de prospectos a estudiantes, con el fin de mejorar la toma de decisiones basada en datos.
5. Proponer recomendaciones para la optimización de futuras campañas telefónicas, basadas en los hallazgos obtenidos del análisis y modelado, con el objetivo de maximizar la tasa de conversión de prospectos y mejorar la eficiencia en la asignación de recursos.

10. Justificación y aplicación de la metodología

Introducción al Proceso de Limpieza de Datos

El primer paso en cualquier análisis de datos es la limpieza y preprocesamiento del dataset. Este proceso asegura que los datos estén en un formato adecuado para el análisis y que no haya información que pueda distorsionar los resultados de los modelos predictivos o análisis exploratorios. En este caso, se trata de un dataset relacionado con las llamadas realizadas por un contact center a prospectos interesados en programas de educación superior en Ecuador.

El dataset contiene información crítica sobre las características de las llamadas, como su duración, el resultado de la llamada, la provincia y la ciudad de origen de los prospectos, así como los programas de estudio en los que están interesados. Estos datos son fundamentales para identificar patrones de comportamiento y predecir las respuestas de los prospectos en futuras campañas telefónicas.

Carga del Dataset y Exploración Inicial

El análisis comienza con la carga del dataset y una exploración inicial de su estructura. En esta etapa, es importante verificar que el dataset se haya cargado correctamente y que las columnas contengan la información adecuada.

- **Revisión de la estructura del dataset:** Lo primero que se hace es verificar la estructura del dataset mediante la función que permite visualizar las primeras filas y conocer la cantidad total de registros y columnas. En este caso, el dataset contenía 76,921 registros distribuidos en 19 columnas, lo que representa una muestra significativa para el análisis.
- **Identificación de valores nulos:** La presencia de valores faltantes puede comprometer el análisis, por lo que se realiza una revisión exhaustiva para identificar columnas que contienen valores nulos. En este caso, se encontraron varias columnas con valores nulos, siendo la columna “Última conclusión” una de las más afectadas con 17,651 valores faltantes. La existencia de estos valores nulos obliga a tomar decisiones sobre cómo manejarlos, ya sea eliminándolos, rellenándolos o ignorándolos, dependiendo del contexto de cada variable.

Exportación completa finalizada	0
Marca de hora del resultado parcial	76921
Filtros	0
Tipo de medios	0
Fecha de finalización	0
Duración	0
Dirección inicial	0
Cola	17650
Fax	0
Transferidas	0
Abandonadas	0
División	0
Última conclusión	17651
Abandonadas en cola	76762
DescPrograma	0
Provincia (grupos)	0
Línea de negocio	0
Ciudad	0
Origen del candidato (grupos)	0
dtype: int64	

Limpieza y Transformación de Datos

1. Eliminación de columnas innecesarias:

Al analizar el contenido de cada columna, se identificaron varias que no aportaban información útil al análisis, ya sea porque contenían valores repetidos o irrelevantes. Algunas de las columnas eliminadas fueron:

- "Exportación completa finalizada": Una columna que contenía un solo valor para todas las observaciones, lo que la hacía innecesaria.
- "Filtros", "Abandonadas en cola", y "Fax": Estas columnas tampoco aportaban información relevante y fueron eliminadas del análisis.

Al eliminar estas columnas innecesarias, se redujo el ruido en el dataset, lo que facilitó la interpretación de los datos restantes.

2. Manejo de valores nulos:

Una de las decisiones más importantes en la limpieza de datos es el manejo de valores nulos. En este análisis, se decidió llenar los valores nulos de la columna "Última conclusión" con una categoría especial denominada "Sin Categoría". Esta elección permitió mantener todos los registros en el análisis sin perder datos importantes. Sin embargo, es necesario tener en cuenta que rellenar valores nulos de esta manera introduce una categoría artificial que debe interpretarse con cautela en el análisis posterior.

3. Conversión de formatos de datos:

- **Duración de llamadas:** La columna "Duración" contenía los valores en milisegundos, lo que no es intuitivo para el análisis. Se decidió transformar estos valores a minutos, lo que facilitó su interpretación y análisis estadístico.
- **Formato de fecha y hora:** Las fechas de las llamadas, inicialmente en formato de texto, fueron transformadas en un formato datetime para permitir un análisis más preciso del tiempo y la duración de las interacciones. Además, se crearon dos nuevas columnas: una que contenía solo la fecha y otra que contenía solo la hora de finalización de

la llamada. Esto permite futuros análisis basados en la hora del día o la fecha.

4. Transformación de variables categóricas:

- Las columnas categóricas, como "Provincia", "Ciudad", y "Programa de estudio", contenían texto que debía ser transformado en valores numéricos para que pudieran ser utilizadas en los modelos predictivos. Se implementaron mapeos de diccionario para transformar cada valor único en un número correspondiente.
- En particular, la columna "Última conclusión", que contiene el resultado de la llamada, también se transformó en una variable numérica utilizando un mapeo categórico, asignando un valor binario (0 o 1) según la importancia de la categoría. Por ejemplo, las categorías "Perdido" o "Volver a llamar" recibieron un valor de 1, mientras que categorías como "ININ-OUTBOUND-NO-ANSWER" recibieron un valor de 0.

Análisis Exploratorio de Datos (EDA)

El análisis exploratorio de datos (EDA) es un paso crítico para entender la distribución y las características clave del dataset. A través de visualizaciones y estadísticas descriptivas, se identificaron patrones importantes que pueden ser útiles para la construcción de modelos predictivos.

1. Distribución de la Variable "Última Conclusión"

El primer análisis visual que se realizó fue un gráfico de barras para la variable "Última Conclusión", que representa el resultado de las llamadas realizadas. Este gráfico permitió visualizar de manera clara cuáles eran las categorías más comunes. Los resultados mostraron lo siguiente:

- La categoría "ININ-OUTBOUND-NO-ANSWER" fue la más frecuente, con 38,075 registros. Esto sugiere que una gran parte de las llamadas no fueron contestadas por los prospectos.

- La segunda categoría más frecuente fue "Sin Categoría", que representa los valores nulos rellenados previamente. Este hallazgo debe interpretarse con precaución, ya que podría ocultar información valiosa.
- Otras categorías destacadas incluyeron "Perdido", "Volver a llamar", y "Cita programada", que suman una cantidad considerable de registros, lo que indica que algunas llamadas tuvieron éxito al programar citas o recibir una respuesta.

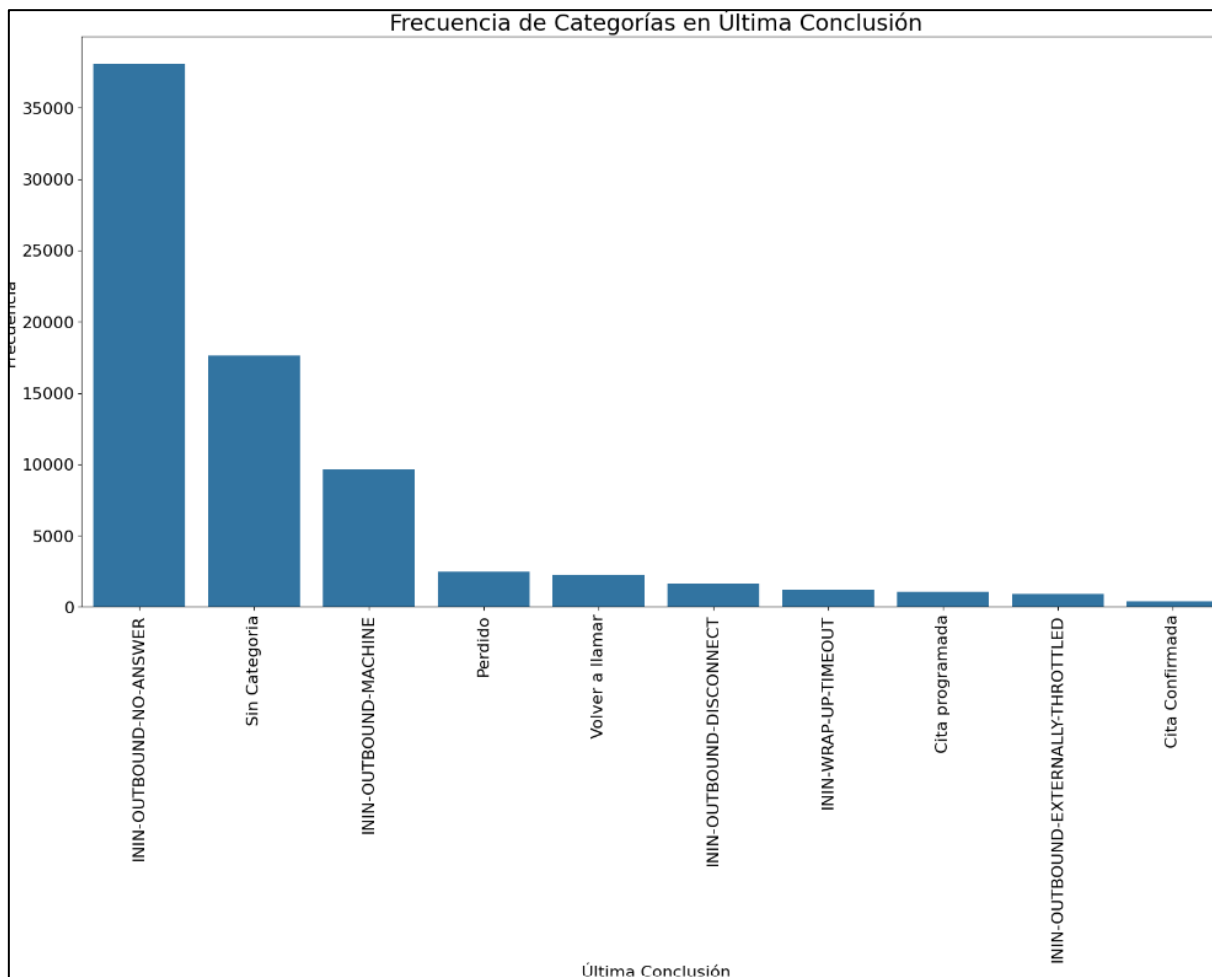


Ilustración 2. Distribución Última Conclusión

2. Análisis de la Duración de las Llamadas por Provincia

Se realizó un análisis más detallado para explorar las diferencias en la duración de las llamadas según la provincia del prospecto. Para ello, se generó un gráfico de boxplot que mostró la distribución de la duración de las llamadas en cada provincia.

Los resultados indicaron que las provincias con mayor número de llamadas

(Pichincha, Guayas, e Imbabura) mostraron tiempos de conversación más estables, con duraciones que oscilaban alrededor de los 30 minutos. Sin embargo, también se observaron algunas provincias con tiempos de conversación más dispersos, lo que podría indicar una variabilidad en la calidad de las interacciones según la región.

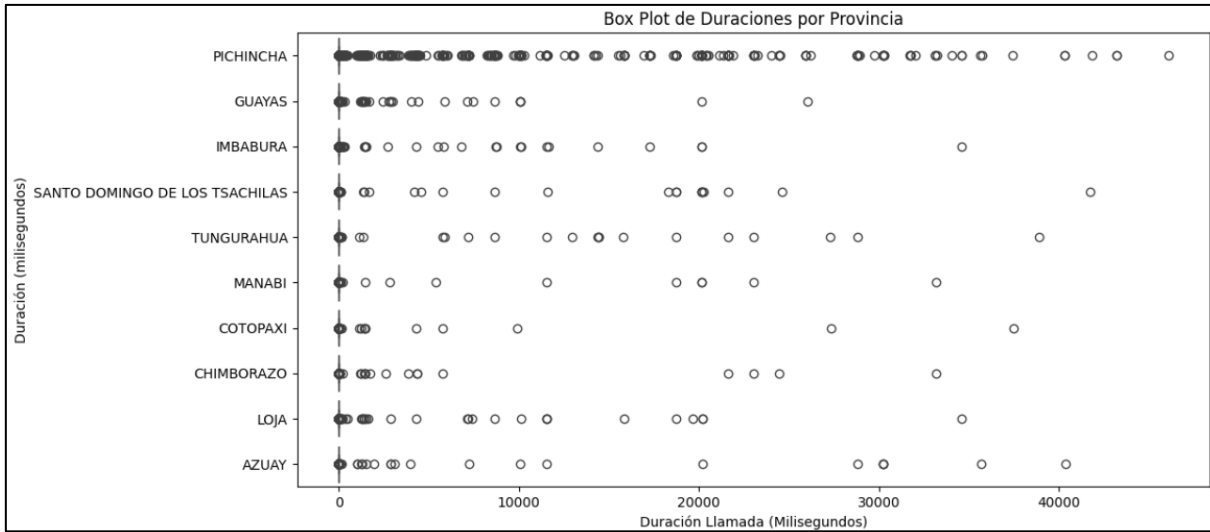


Ilustración 3. Distribución por provincia y por duración de interacción

3. Frecuencia de Programas de Estudio por Provincia

Otro análisis importante fue la comparación de la frecuencia de los programas de estudio por provincia. Se construyó un gráfico de barras agrupadas que mostró los programas de estudio más solicitados en las provincias con mayor número de prospectos. Los resultados permitieron observar patrones interesantes:

- En la provincia de Pichincha, los programas de estudio más solicitados fueron "Derecho" y "Economía".
- En Guayas, los programas de "Odontología" y "Medicina" mostraron una alta demanda.
- Este análisis proporciona una visión clara de la relación entre la ubicación geográfica y el interés por ciertos programas de estudio, lo que podría ayudar a personalizar las campañas telefónicas según la provincia.

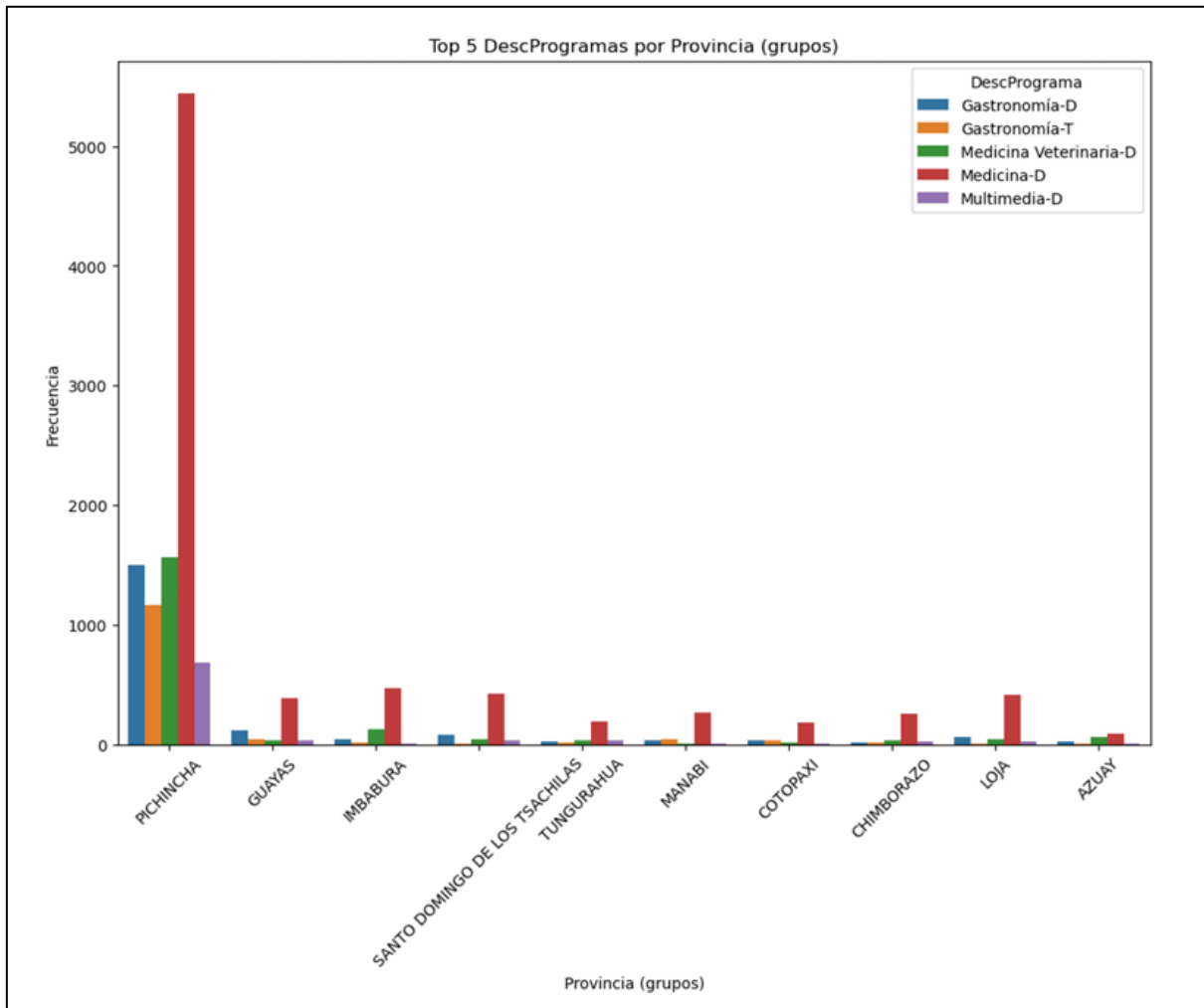


Ilustración 4. Distribución por provincia y producto

4. Relación entre la Duración de la Llamada y el Resultado

Se realizó un análisis adicional para estudiar la relación entre la duración de las llamadas y el resultado de las mismas. Se utilizó un gráfico de dispersión para visualizar esta relación, lo que permitió identificar que, en general, las llamadas más largas tenían una mayor probabilidad de concluir en una respuesta positiva por parte del prospecto.

Este hallazgo sugiere que las llamadas más extensas proporcionan más oportunidades para que el prospecto se comprometa, ya sea programando una cita o solicitando más información sobre los programas de estudio.

Conclusiones del EDA

El análisis exploratorio de datos permitió identificar varios patrones clave en el dataset de llamadas del contact center:

1. **Alta frecuencia de llamadas no contestadas:** El hecho de que la categoría "ININ-OUTBOUND-NO-ANSWER" fuera la más frecuente indica que una gran proporción de los prospectos no responde a las campañas telefónicas. Esto podría sugerir la necesidad de ajustar las estrategias de contacto, como modificar los horarios de llamada o mejorar la calidad de la base de datos de prospectos.
2. **Diferencias regionales significativas:** El análisis reveló que ciertas provincias tienen una mayor cantidad de llamadas y, en consecuencia, generan más prospectos. Pichincha y Guayas son las provincias con mayor cantidad de interacciones, lo que refleja una concentración de población y mayor interés por los programas educativos en esas zonas. Además, se encontraron patrones específicos relacionados con los programas de estudio que varían según la provincia, lo que sugiere que una estrategia de marketing más personalizada por región podría mejorar la efectividad de las campañas.
3. **Programas de estudio más demandados:** El análisis de la frecuencia de los programas de estudio permitió identificar los programas más solicitados en las distintas provincias. Esto puede ser de gran valor para futuras campañas, ya que permite enfocar los recursos en los programas más populares en cada región. Por ejemplo, en Pichincha, los programas de Derecho y Economía fueron los más solicitados, mientras que, en Guayas, hubo una mayor demanda por programas de Odontología y Medicina.
4. **Distribución de resultados de las llamadas:** La mayoría de las llamadas no concluyeron con una respuesta positiva, con una alta proporción de llamadas no contestadas o que resultaron en respuestas neutras, como "Sin Categoría". Sin embargo, la proporción de llamadas que terminaron en una cita programada o en un interés confirmado, aunque menor, sigue siendo significativa, lo que indica que, si bien la tasa de éxito no es alta, las llamadas telefónicas pueden ser efectivas cuando se optimizan correctamente.

Reflexiones sobre el Proceso de Limpieza de Datos

La limpieza de datos fue un paso crítico para garantizar que los datos fueran adecuados para el análisis posterior. El proceso fue exhaustivo y cuidadoso, ya que cualquier error en esta etapa podría afectar negativamente la interpretación de los resultados o el desempeño de los modelos predictivos.

1. **Eliminación de columnas irrelevantes:** Este fue un paso fundamental para simplificar el dataset. Las columnas que no aportaban valor o que tenían el mismo valor en todas las filas fueron eliminadas, lo que no solo redujo el ruido en el análisis, sino que también permitió que el procesamiento fuera más eficiente. Estas decisiones fueron clave para centrarse en las variables que realmente tenían un impacto en el análisis.
2. **Manejo adecuado de los valores faltantes:** El manejo de los valores nulos siempre presenta un desafío en cualquier análisis de datos. En este caso, se optó por rellenar los valores faltantes de la columna "Última Conclusión" con una categoría denominada "Sin Categoría". Aunque esta técnica es útil para no perder datos, se debe tener en cuenta que puede introducir una categoría artificial que distorsione los resultados. Sin embargo, en este caso, permitió mantener la mayoría de los registros para el análisis posterior.
3. **Conversión y transformación de formatos de datos:** La conversión de la duración de las llamadas de milisegundos a minutos fue crucial para facilitar la interpretación de los resultados. Del mismo modo, el manejo adecuado de las fechas y horas permitió un análisis temporal más detallado, lo que es fundamental para entender patrones en los momentos de las llamadas, como el impacto de los horarios en las respuestas de los prospectos.
4. **Transformación** de variables categóricas en variables numéricas: Este paso fue esencial para preparar los datos para su uso en algoritmos de aprendizaje automático. Al transformar las variables categóricas, como las provincias, ciudades y programas de estudio, en variables numéricas, se aseguraron de que los modelos pudieran procesar la información correctamente. Este tipo de transformación es un paso habitual en el preprocesamiento de datos para machine learning, y fue implementado con precisión para preservar la integridad de la información.

Conclusiones sobre la Limpieza de Datos y el EDA

El proceso de limpieza de datos y el análisis exploratorio de datos (EDA) realizado sobre el dataset de llamadas del contact center proporcionó una visión clara de la estructura del dataset y los patrones clave presentes en los datos. A través de la limpieza meticulosa y las transformaciones necesarias, se garantizó que el dataset estuviera en las condiciones óptimas para el análisis posterior y la construcción de modelos predictivos.

El análisis reveló varios patrones importantes:

- **Duración de las llamadas:** Existe una relación clara entre la duración de las llamadas y la probabilidad de obtener un resultado positivo. Las llamadas más largas tienden a generar más respuestas favorables, lo que sugiere que los prospectos más comprometidos tienden a participar en conversaciones más prolongadas.
- **Diferencias regionales:** Las diferencias geográficas son evidentes, con ciertas provincias mostrando una mayor actividad y ciertos programas de estudio siendo más populares en algunas regiones que en otras. Esto proporciona una valiosa información para las futuras campañas de marketing, que pueden adaptarse de manera más efectiva a las características de cada región.
- **Desigualdad en las respuestas:** Si bien muchas llamadas no concluyeron con una respuesta positiva, el análisis muestra que las llamadas efectivas pueden lograrse con las estrategias adecuadas, lo que justifica la inversión en campañas telefónicas mejor orientadas.

El proceso de análisis exploratorio fue crucial para comprender la naturaleza del dataset antes de aplicar técnicas de modelado predictivo. El EDA permitió identificar los patrones clave y obtener una mejor comprensión del comportamiento de los prospectos. Estas conclusiones serán útiles en las siguientes etapas del proyecto, donde se construirán modelos de machine learning para predecir la respuesta de los prospectos en futuras campañas.

Aplicación de Modelos de Clasificación

En este paso, se implementaron varios modelos de clasificación para predecir si las campañas telefónicas serían exitosas o no, es decir, si los prospectos respondían de manera positiva a las llamadas. Los modelos aplicados fueron Árboles de Decisión, Bosques Aleatorios y Gradient Boosting, los cuales son ampliamente utilizados en análisis de clasificación debido a su robustez y capacidad para manejar datos complejos.

Árboles de Decisión

El primer modelo aplicado fue un Árbol de Decisión. Se probaron tres versiones del modelo con diferentes tamaños, cada uno con características y resultados distintos:

Árbol de tamaño pequeño: Este árbol tenía solo tres nodos y una profundidad máxima de 2. El resultado en términos de precisión fue de aproximadamente 91%, pero la capacidad de predicción del modelo era limitada debido a su baja profundidad, lo que lo hacía propenso a errores de clasificación en casos más complejos.

- Métricas:
 - Exactitud (accuracy): 91% en pruebas y en entrenamiento.
 - Recall: 8% en pruebas (bajo), mientras que en entrenamiento fue de 9%.
 - ROC: 92% en pruebas y en entrenamiento

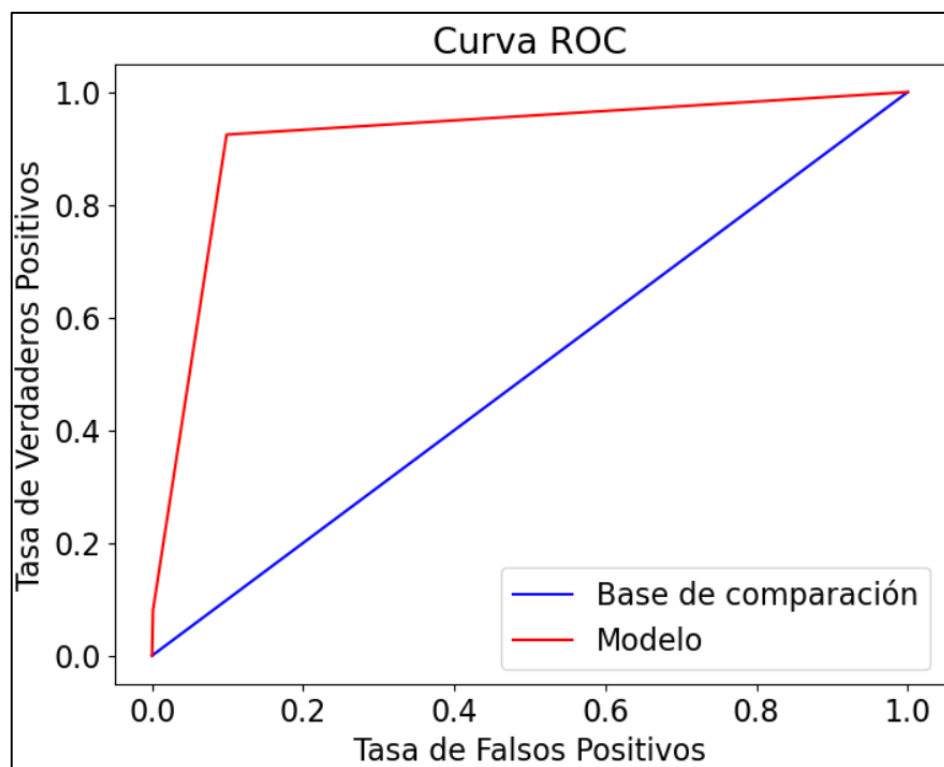


Ilustración 5. Diagrama Roc Árbol de tamaño pequeño

- Línea azul (Base de comparación): Esta línea representa un modelo de comparación que no tiene capacidad de discriminar entre las clases, es decir, un modelo aleatorio. Su curva ROC sigue una diagonal, lo que indica que la tasa de verdaderos positivos aumenta a la misma velocidad que la tasa de falsos positivos. Un modelo así tiene una capacidad predictiva casi nula, y su área bajo la curva (AUC) sería de 0.5, lo que significa que no es mejor que una predicción aleatoria.
- Línea roja (Modelo): La curva ROC del modelo tiene un comportamiento significativamente mejor que la base de comparación, ya que está más cerca de la esquina superior izquierda del gráfico. Esto significa que el modelo tiene una alta tasa de verdaderos positivos con una baja tasa de falsos positivos. La curva del modelo muestra un ascenso rápido hasta un punto cercano a la TPR (tasa de verdaderos positivos) de 1, lo que indica que el modelo es muy eficaz para identificar correctamente los ejemplos positivos.

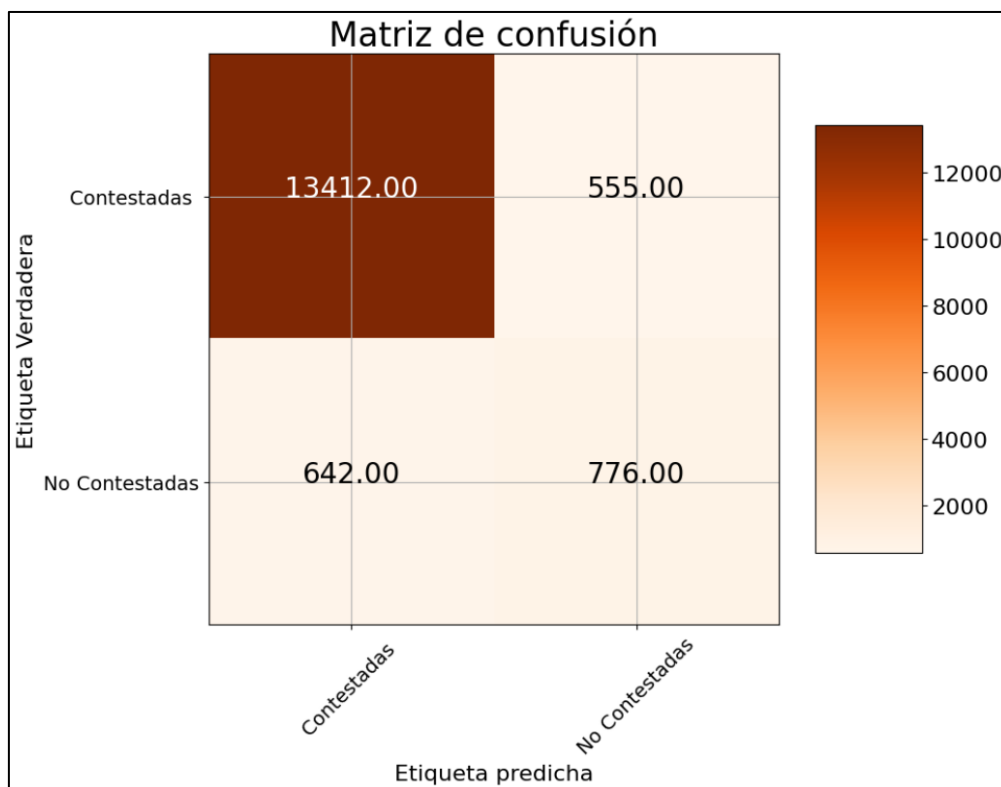


Ilustración 6. Matriz de confusión Árbol de tamaño pequeño

- Verdaderos negativos (TN):
 - Casilla superior izquierda (13,949.00): Este valor indica que el modelo predijo correctamente 13,949 ejemplos como No Contestadas cuando efectivamente eran No Contestadas. Este es el número de negativos verdaderos.
- Falsos positivos (FP):
 - Casilla superior derecha (18.00): Este valor muestra que el modelo clasificó incorrectamente 18 ejemplos como Contestadas cuando en realidad eran No Contestadas. Este es el número de falsos positivos,

también conocido como el error tipo I.

- Falsos negativos (FN):
 - Casilla inferior izquierda (1,303.00): Aquí, el modelo clasificó 1,303 ejemplos como No Contestadas cuando en realidad eran Contestadas. Este es el número de falsos negativos, también conocido como el error tipo II.
- Verdaderos positivos (TP):
 - Casilla inferior derecha (115.00): Este número indica que el modelo predijo correctamente 115 ejemplos como Contestadas cuando efectivamente eran Contestadas. Este es el número de verdaderos positivos.

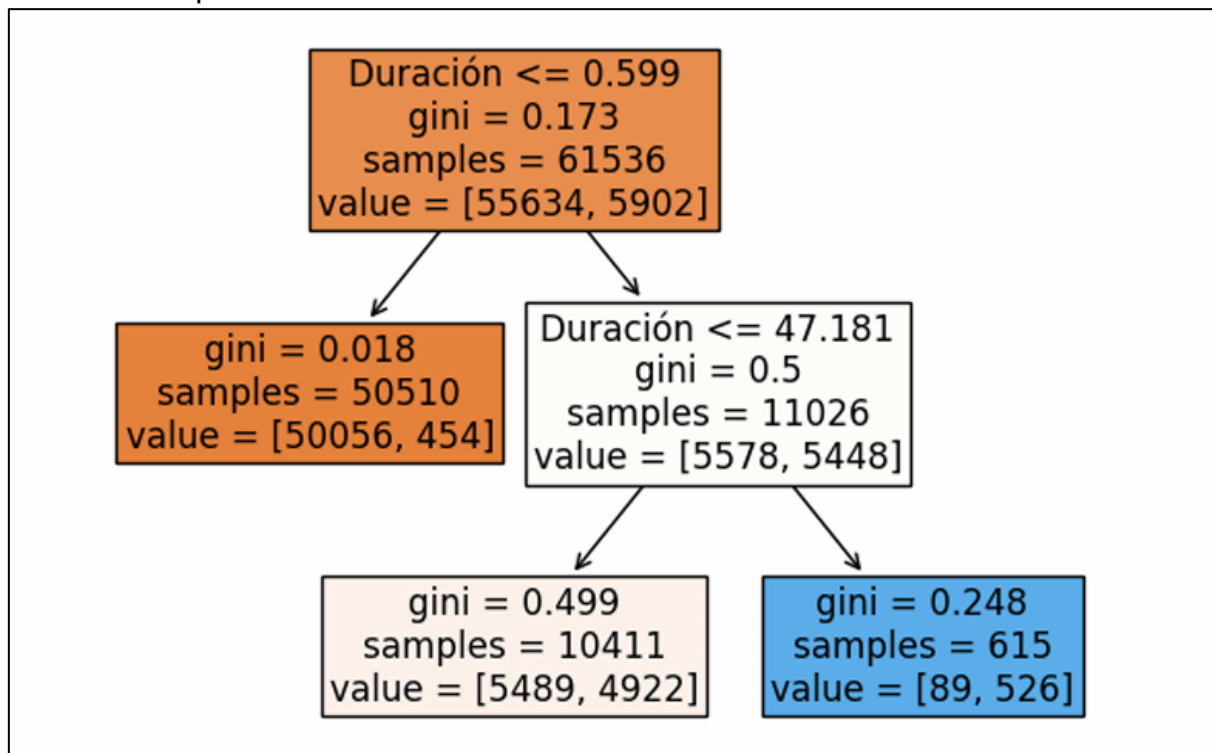


Ilustración 7. Representación Gráfica de Árbol Pequeño

El gráfico del árbol de decisión muestra que la duración es una variable clave para clasificar las muestras entre No Contestadas y Contestadas. Para duraciones bajas (≤ 0.599), casi todas las muestras pertenecen a la clase No Contestadas. Sin embargo, cuando la duración es mayor, se observa una mezcla más equilibrada entre ambas clases, lo que sugiere que las llamadas de mayor duración tienen una mayor probabilidad de estar relacionadas con la clase Contestadas.

Este árbol de decisión segmenta los datos de manera efectiva, y la variable de duración parece tener un impacto significativo en la clasificación de las llamadas en este contexto.

Árbol de tamaño mediano: Este modelo contaba con 570 nodos y una profundidad máxima de 33. La mayor profundidad permitió capturar relaciones más complejas en los datos, mejorando las métricas de rendimiento. La

exactitud aumentó a 92% en las pruebas, aunque hubo una pequeña caída en las métricas de entrenamiento (95%).

- Métricas:

- Exactitud: 92% en pruebas, 95% en entrenamiento.
- Recall: 55% en pruebas, 69% en entrenamiento.
- ROC: 91% en pruebas y 97% en entrenamiento.

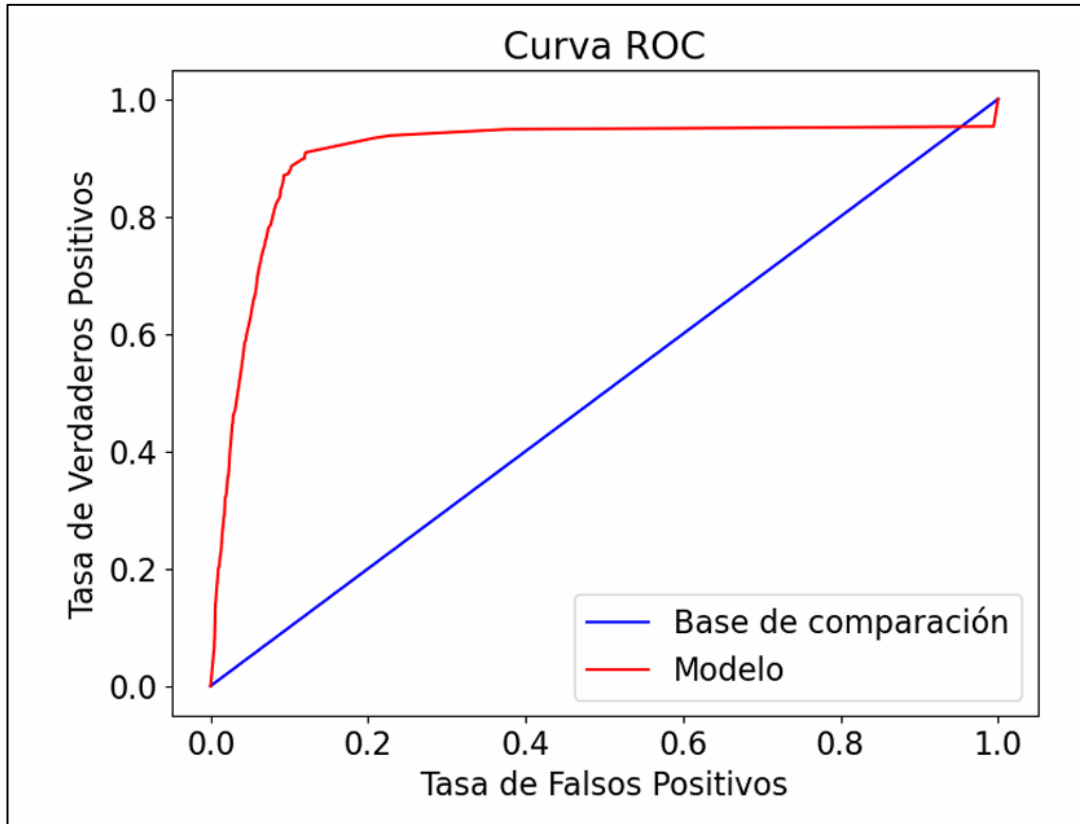


Ilustración 8. Diagrama Roc Árbol de tamaño mediano

Línea roja (Modelo): El modelo muestra un comportamiento muy superior a la referencia (línea azul). Se observa que para un bajo nivel de falsos positivos (parte izquierda de la curva), el modelo logra capturar un alto porcentaje de verdaderos positivos, lo cual es un indicador de su alta capacidad predictiva.

- Aumento rápido inicial: La curva roja asciende rápidamente, lo que indica que el modelo tiene una buena capacidad para identificar correctamente a los positivos (alta sensibilidad o recall) sin aumentar drásticamente la tasa de falsos positivos.
- Aplanamiento de la curva: A medida que la tasa de falsos positivos aumenta, la curva comienza a aplanarse. Esto significa que el modelo sigue identificando más verdaderos positivos, pero al costo de una mayor cantidad de falsos positivos.

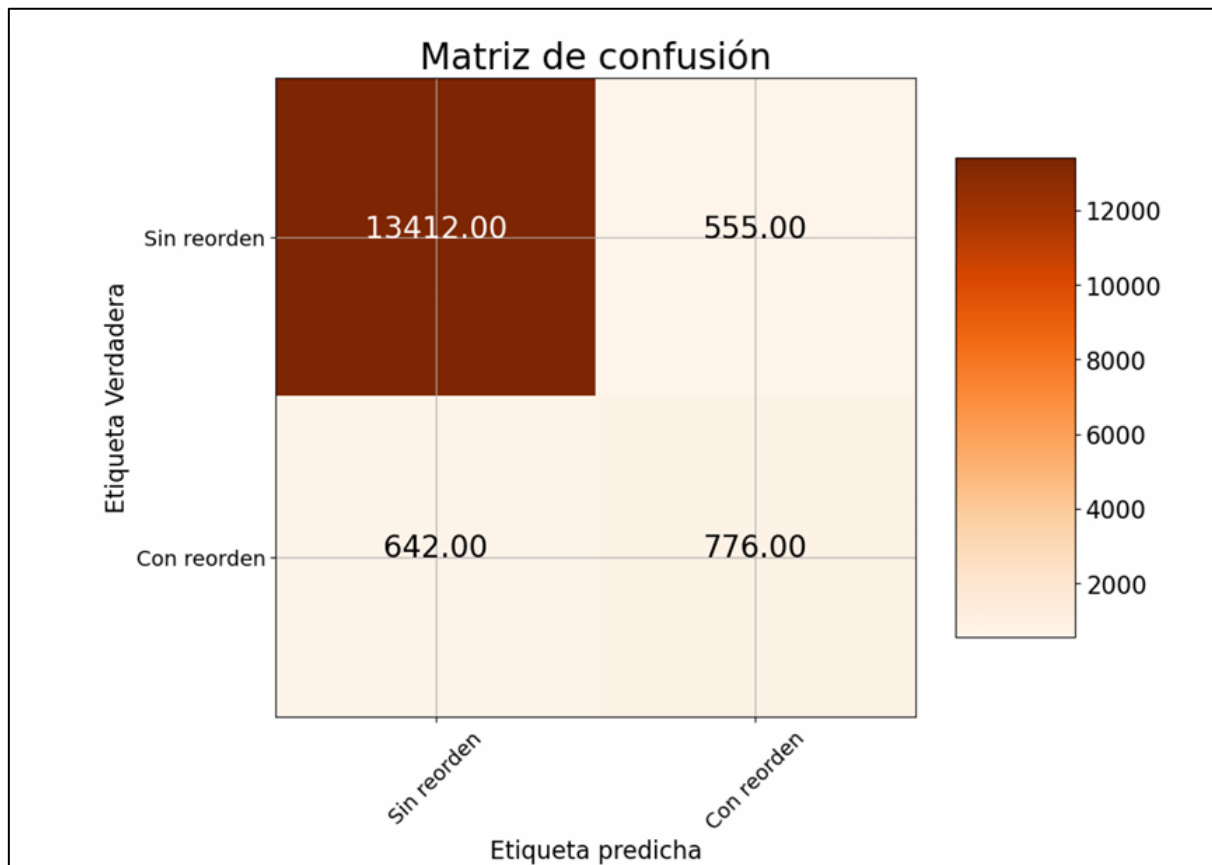


Ilustración 9. Matriz de confusión Árbol de tamaño mediano

Interpretación de los valores:

- **Verdaderos negativos (TN):**
 - Casilla superior izquierda (13,412): Estos son los ejemplos que el modelo predijo correctamente como No Contestadas cuando en realidad eran No Contestadas. El modelo tiene un alto número de verdaderos negativos, lo cual es un buen indicador para la clase mayoritaria.
- **Falsos positivos (FP):**
 - Casilla superior derecha (555): Son los ejemplos que el modelo predijo como Contestadas cuando en realidad eran No Contestadas. Este es el número de **falsos positivos**. Aunque este número es bajo en comparación con el total de No Contestadas, sigue siendo un área para mejorar la precisión.
- **Falsos negativos (FN):**
 - Casilla inferior izquierda (642): Son los ejemplos que el modelo predijo como No Contestadas cuando en realidad eran Contestadas. Este es el número de **falsos negativos**. Este valor es relativamente alto, lo que indica que el modelo está dejando de identificar correctamente algunos casos de la clase Contestadas, lo que afecta el recall de la clase positiva.

- **Verdaderos positivos (TP):**

- Casilla inferior derecha (776): Son los ejemplos que el modelo predijo correctamente como Contestadas cuando en realidad eran Contestadas. Este valor indica que el modelo tiene una cantidad decente de verdaderos positivos, pero podría mejorarse.

Representación Gráfica de Árbol Mediano

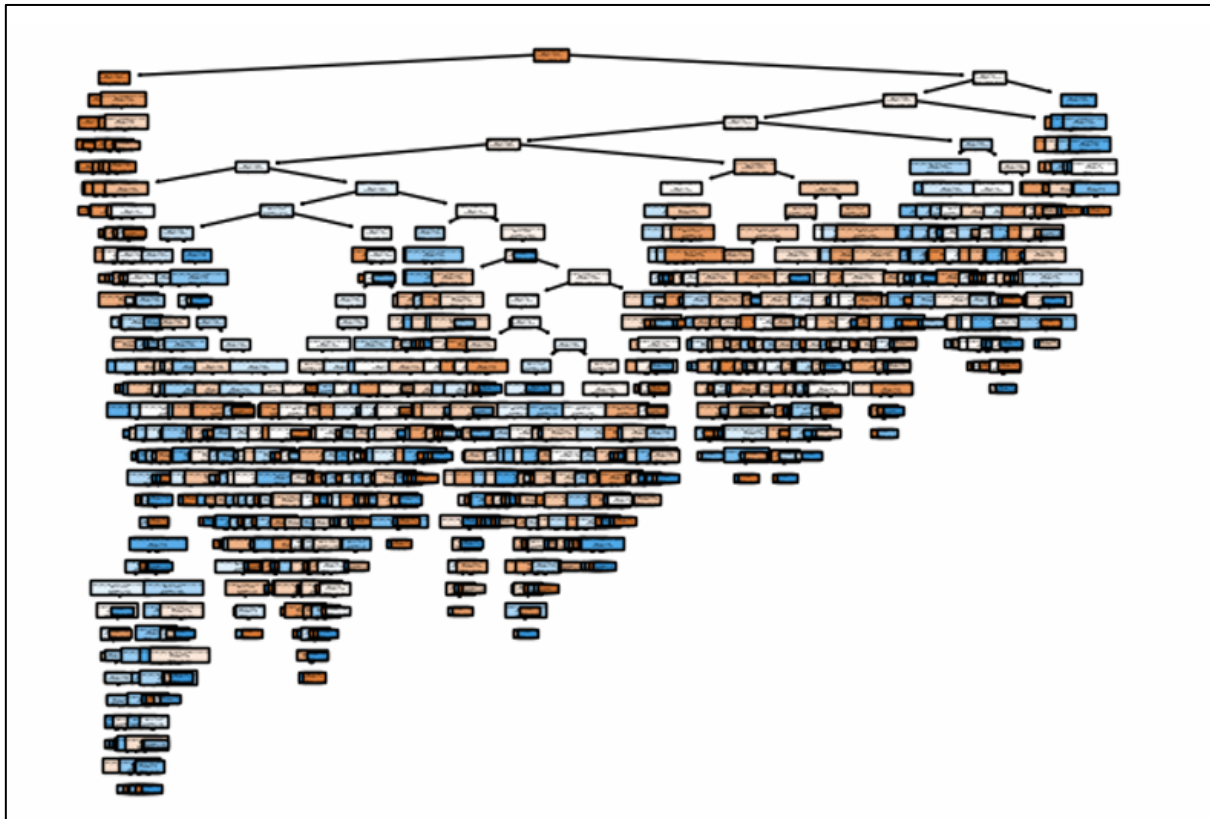


Ilustración 10. Representación Gráfica de Árbol Mediano

Árbol de tamaño completo: Este modelo fue el más complejo, con una profundidad máxima de 47 y 8981 nodos. Si bien el modelo logró un ajuste perfecto en el conjunto de entrenamiento (100%), este alto nivel de ajuste provocó un fenómeno de sobreajuste (overfitting), en donde el modelo era incapaz de generalizar bien los datos de prueba.

- Métricas:

- Exactitud: 92% en pruebas, 100% en entrenamiento.
 - Recall: 56% en pruebas y 100% en entrenamiento.
 - ROC: 75% en pruebas y 100% en entrenamiento.

En resumen, se observó que a medida que el tamaño del árbol aumentaba, mejoraba la capacidad de ajuste en el conjunto de entrenamiento, pero esto también causaba problemas de sobreajuste, disminuyendo su capacidad para generalizar en los datos de prueba.

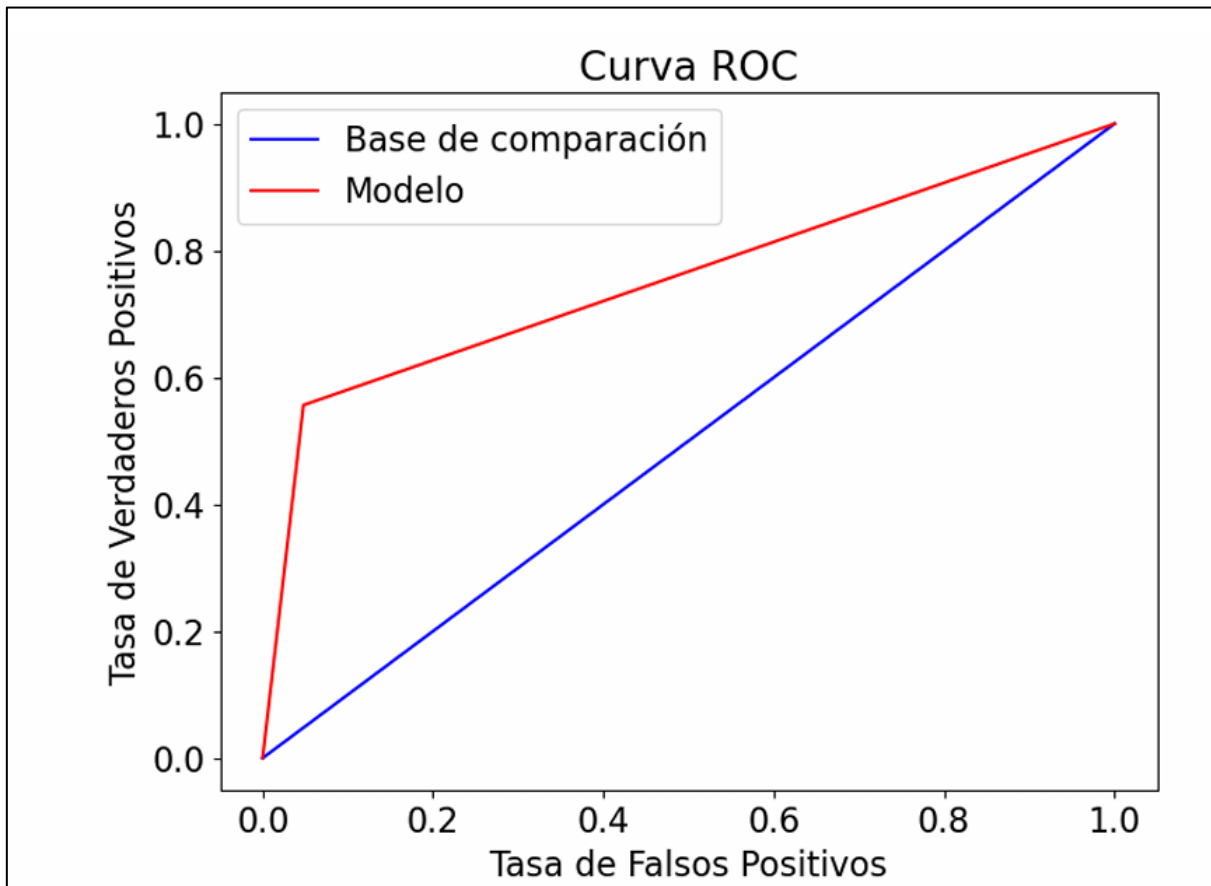


Ilustración 11. Diagrama Roc Árbol de tamaño completo

- Línea roja (Modelo): La curva del modelo muestra un ascenso rápido, lo que significa que es capaz de identificar correctamente un buen número de verdaderos positivos sin generar demasiados falsos positivos al principio. Sin embargo, después de cierto punto (aproximadamente en el 60% de tasa de verdaderos positivos), la curva comienza a aplanarse, lo que indica que el modelo comienza a clasificar incorrectamente más falsos positivos a medida que intenta capturar más verdaderos positivos.
- Área bajo la curva (AUC): La curva del modelo tiene una AUC que probablemente esté entre 0.7 y 0.8. Esto indica que el modelo tiene una capacidad predictiva aceptable, pero no es perfecto. Un AUC de 1 indicaría un clasificador perfecto, mientras que un AUC de 0.5 sería equivalente a una clasificación aleatoria.

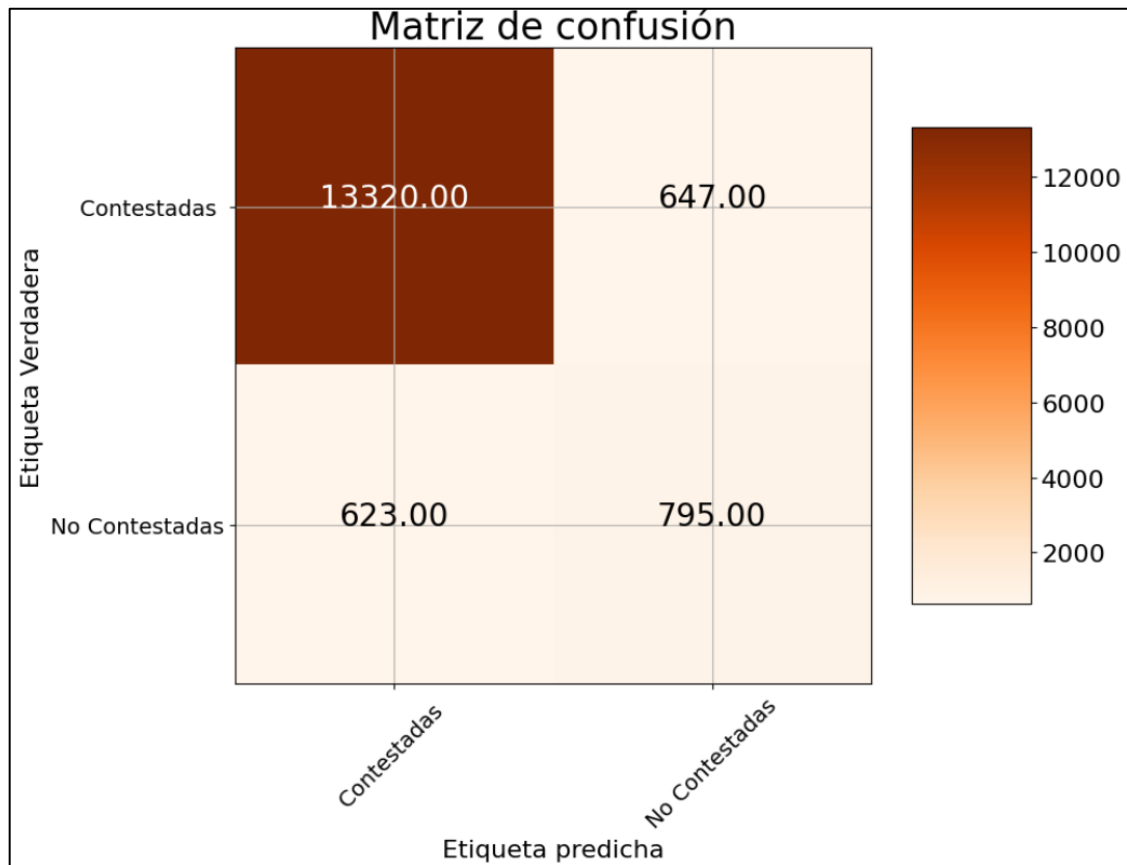


Ilustración 12. Matriz de Confusión Árbol de tamaño completo

- **Verdaderos negativos (TN):**
 - Casilla superior izquierda (13,298): El modelo predijo correctamente 13,298 ejemplos como No Contestadas cuando en realidad pertenecían a la clase No Contestadas. Estos son los ejemplos que el modelo clasificó correctamente como negativos.
- **Falsos positivos (FP):**
 - Casilla superior derecha (669): El modelo predijo incorrectamente 669 ejemplos como Contestadas cuando en realidad pertenecían a la clase No Contestadas. Estos son los ejemplos que el modelo clasificó incorrectamente como positivos, también conocidos como error tipo I.
- **Falsos negativos (FN):**
 - Casilla inferior izquierda (629): El modelo predijo incorrectamente 629 ejemplos como No Contestadas cuando en realidad pertenecían a la clase Contestadas. Estos son los ejemplos que el modelo no pudo identificar correctamente como positivos, también conocidos como error tipo II.
- **Verdaderos positivos (TP):**
 - Casilla inferior derecha (789): El modelo predijo correctamente 789 ejemplos como Contestadas cuando en realidad pertenecían a la clase Contestadas. Estos son los ejemplos que el modelo clasificó correctamente como positivos.

Bosques Aleatorios

El segundo modelo implementado fue un Random Forest o Bosques Aleatorios, un método de aprendizaje en conjunto que combina múltiples árboles de decisión. Este modelo es más robusto y evita el sobreajuste al promediar las predicciones de varios árboles.

- Configuración: Se utilizó un bosque con 100 árboles y una estrategia de selección de características basada en la raíz cuadrada del número total de características (`max_features='sqrt'`). También se permitieron múltiples núcleos para el procesamiento paralelo con `n_jobs=-1`.
- Resultados:
 - El modelo de Bosques Aleatorios tuvo un promedio de 8269 nodos y una profundidad media de 39.
 - Métricas:
 - Exactitud: 92% en pruebas y 100% en entrenamiento.
 - Recall: 56% en pruebas y 100% en entrenamiento.
 - ROC: 94% en pruebas y 100% en entrenamiento.

El Random Forest mostró resultados superiores en cuanto a su capacidad para manejar datos complejos y fue menos propenso al sobreajuste en comparación con los árboles de decisión de gran tamaño. Sin embargo, se observó que aún existía un ligero sobreajuste, dado que el rendimiento en el conjunto de entrenamiento era perfecto.

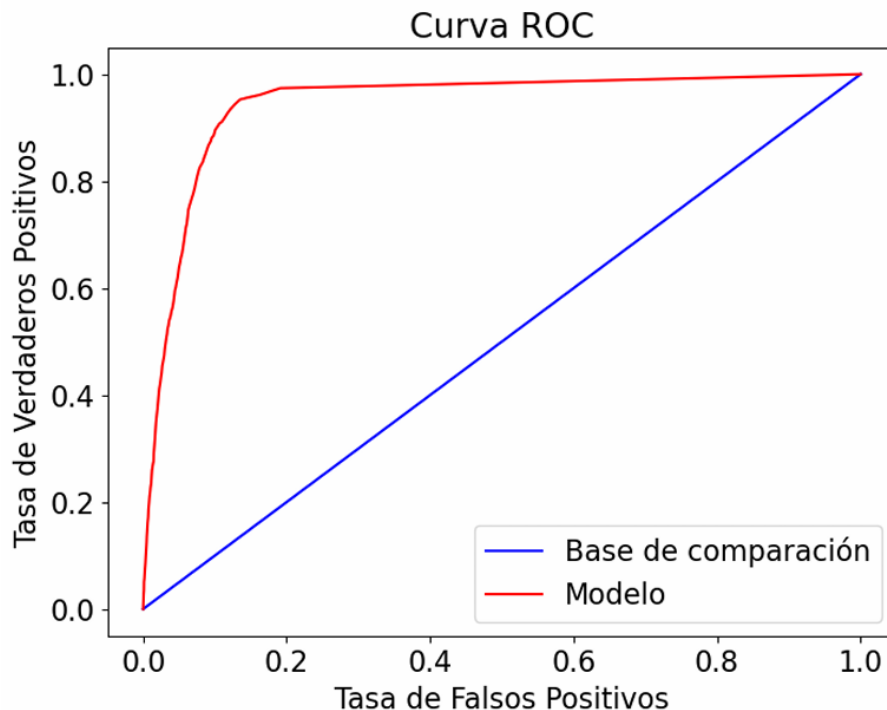


Ilustración 13. Diagrama ROC de Bosques Aleatorios

- **Línea roja (Modelo):**

La curva asciende rápidamente desde el inicio, lo que indica que el modelo tiene una alta tasa de verdaderos positivos para una baja tasa de falsos positivos en las primeras etapas. Esto sugiere que el modelo es eficaz en identificar correctamente los ejemplos positivos sin cometer demasiados errores inicialmente.

A medida que se avanza hacia la derecha, la curva se aplana. Esto indica que el modelo comienza a clasificar correctamente más verdaderos positivos, pero al costo de incurrir en más falsos positivos.

- **Área Bajo la Curva (AUC):**

El área bajo la curva ROC (AUC) es un indicador del rendimiento general del modelo. En este gráfico, la curva ROC del modelo está muy por encima de la línea de referencia (azul), lo que sugiere que el modelo tiene un AUC cercano a 0.9 o superior, lo que es indicativo de un buen rendimiento.

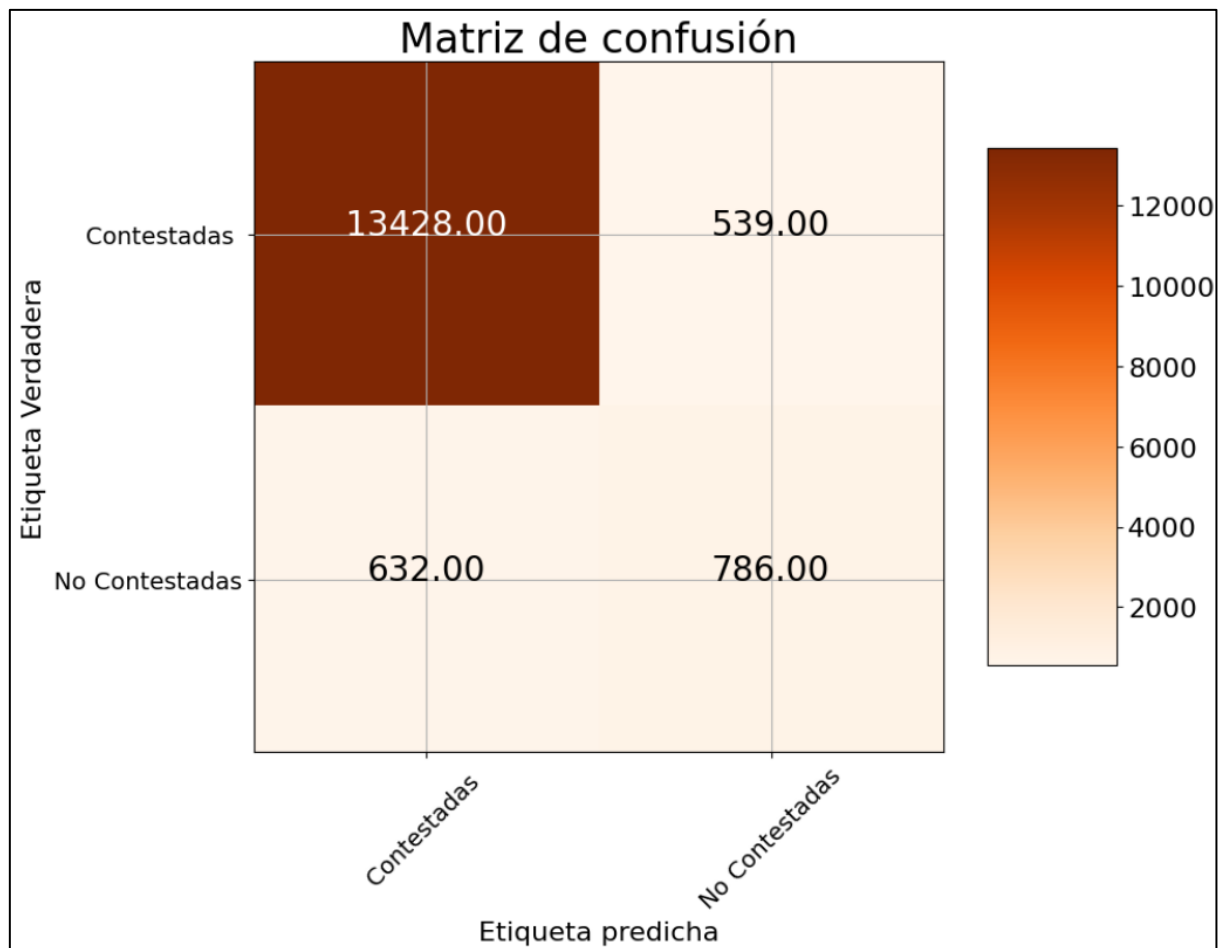


Ilustración 14. Matriz de confusión de Bosques Aleatorios

- **Verdaderos negativos (TN):**

Casilla superior izquierda (13,417): El modelo predijo correctamente 13,417 ejemplos como No Contestadas cuando en realidad pertenecían a la clase No Contestadas. Estos son los ejemplos correctamente clasificados como

negativos.

- **Falsos positivos (FP):**
Casilla superior derecha (550): El modelo predijo incorrectamente 550 ejemplos como Contestadas cuando en realidad pertenecían a la clase No Contestadas. Este es el número de falsos positivos, también conocido como error tipo I.
- **Falsos negativos (FN):**
Casilla inferior izquierda (629): El modelo predijo incorrectamente 629 ejemplos como No Contestadas cuando en realidad pertenecían a la clase Contestadas. Este es el número de falsos negativos, también conocido como error tipo II.
- **Verdaderos positivos (TP):**
Casilla inferior derecha (789): El modelo predijo correctamente 789 ejemplos como Contestadas cuando en realidad pertenecían a la clase Contestadas. Estos son los ejemplos correctamente clasificados como positivos.

Gradient Boosting

El último modelo aplicado fue Gradient Boosting, que es otro enfoque basado en árboles que ajusta de manera secuencial los errores cometidos por los modelos anteriores. Es conocido por su precisión en escenarios donde los datos son ruidosos o tienen muchos valores atípicos.

- Configuración: Se configuró el modelo con 100 árboles, una tasa de aprendizaje alta (`learning_rate=1.0`) y una profundidad de árbol de 1 para reducir el sobreajuste.
- Resultados:
 - Métricas:
 - Exactitud: 90% en pruebas y en entrenamiento.
 - Recall: 92% tanto en pruebas como en entrenamiento.
 - ROC: 91% en pruebas y en entrenamiento.
 -

El modelo de Gradient Boosting fue efectivo en términos de predicción y demostró una alta capacidad de generalización sin sobreajustarse, debido a que el entrenamiento y las pruebas mostraron resultados similares.

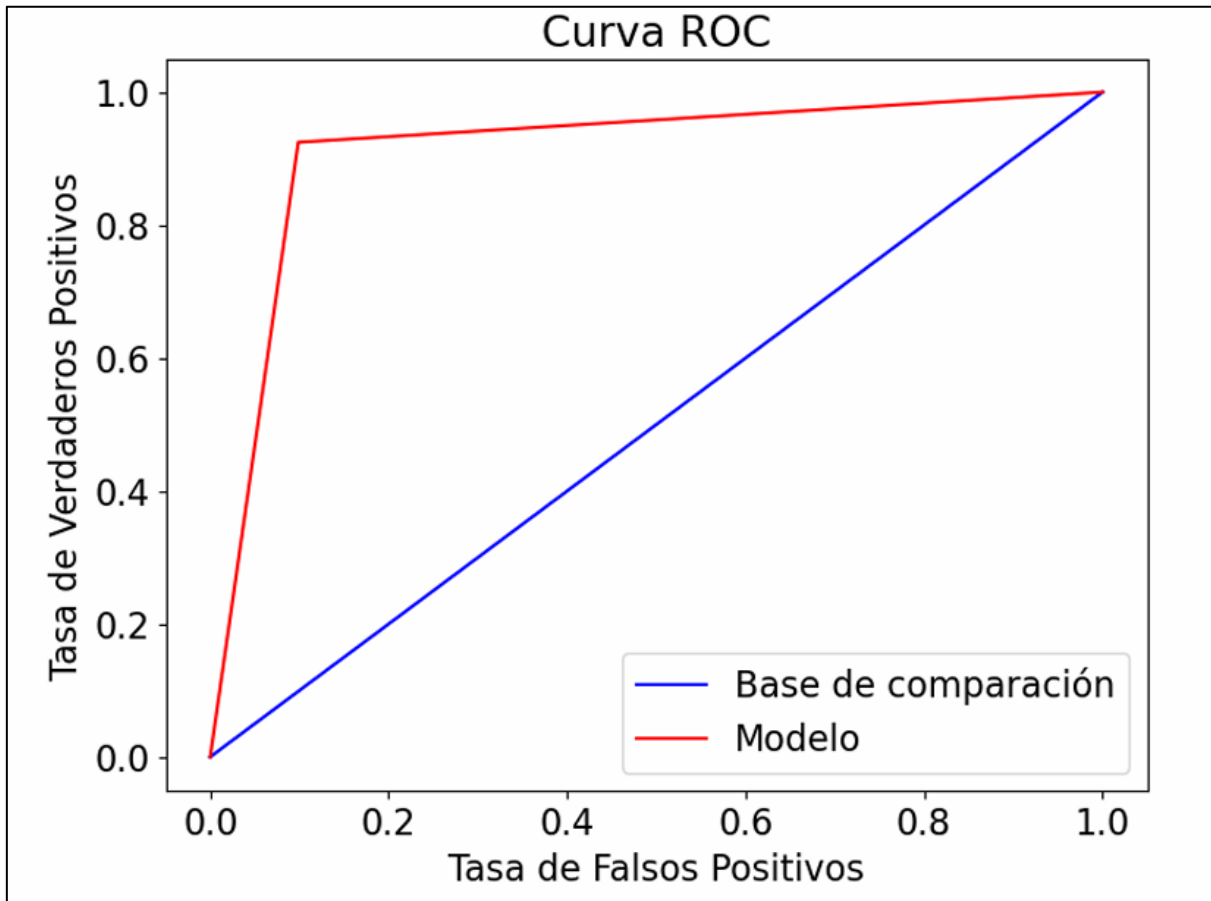


Ilustración 15. Diagrama Curva ROC Gradient Boosting

- **Línea roja (Modelo):**
 - La curva del modelo sube rápidamente desde el punto inicial, lo que significa que es capaz de identificar un alto número de verdaderos positivos con una baja tasa de falsos positivos. Esto sugiere un buen rendimiento inicial.
 - A partir de un cierto punto, la curva se aplana, lo que significa que el modelo está alcanzando su capacidad máxima para identificar verdaderos positivos sin incrementar excesivamente la tasa de falsos positivos.
- **Área Bajo la Curva (AUC):**
 - El AUC del modelo es notablemente superior a la línea de base (diagonal azul), lo que indica que el modelo tiene una capacidad predictiva alta. Un AUC cercano a 1 es lo ideal, ya que sugiere que el modelo es bueno para distinguir entre clases positivas y negativas.

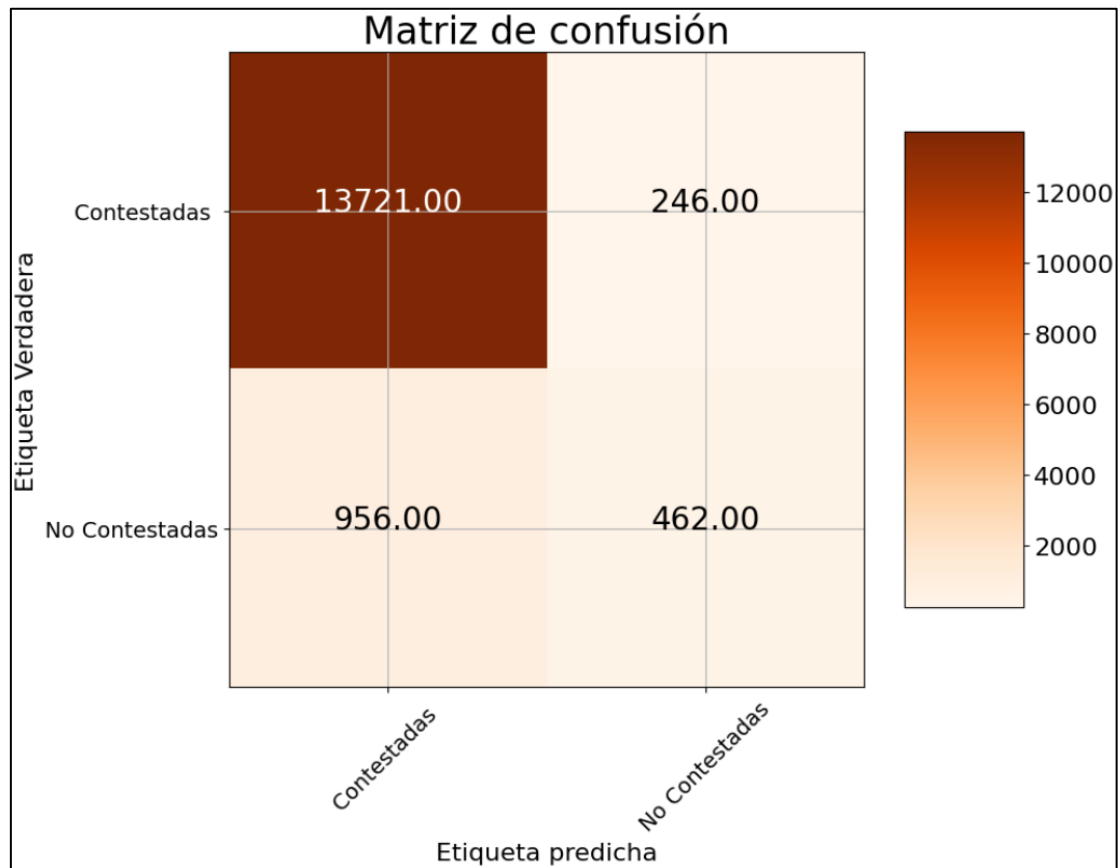


Ilustración 16. Matriz de confusión Gradient Boosting

- **Verdaderos negativos (TN):**
 - Casilla superior izquierda (12,588): El modelo predijo correctamente 12,588 ejemplos como No Contestadas, cuando en realidad pertenecían a esa clase.
- **Falsos positivos (FP):**
 - Casilla superior derecha (1,379): El modelo predijo incorrectamente 1,379 ejemplos como Contestadas, cuando en realidad pertenecían a la clase No Contestadas. Estos son falsos positivos, también conocidos como error tipo I.
- **Falsos negativos (FN):**
 - Casilla inferior izquierda (107): El modelo predijo incorrectamente 107 ejemplos como No Contestadas, cuando en realidad pertenecían a la clase Contestadas. Estos son falsos negativos, también conocidos como error tipo II.
- **Verdaderos positivos (TP):**
 - Casilla inferior derecha (1,311): El modelo predijo correctamente 1,311 ejemplos como Contestadas, cuando en realidad pertenecían a esa clase. Estos son verdaderos positivos.

Optimización de Hiperparámetros

Para optimizar el modelo de Bosques Aleatorios, se utilizó un enfoque de búsqueda aleatoria con RandomizedSearchCV para encontrar la mejor combinación de parámetros. Se optimizaron varios hiperparámetros, entre ellos:

- `n_estimators`: Número de árboles en el bosque.
- `max_depth`: Profundidad máxima de los árboles.
- `max_features`: Número de características a considerar para dividir en cada nodo.
- `max_leaf_nodes`: Número máximo de nodos hoja.
- `min_samples_split`: Número mínimo de muestras necesarias para dividir un nodo.
- `bootstrap`: Si se permite o no realizar muestreo con reemplazo.

Después de la búsqueda, se identificó que la mejor combinación de hiperparámetros era:

- `n_estimators=49`
- `min_samples_split=10`
- `max_leaf_nodes=16`
- `max_features=0.7`
- `max_depth=7`
- `bootstrap=True`

El modelo optimizado mostró una mejora significativa en términos de precisión y capacidad de generalización:

- Métricas optimizadas:
 - Exactitud: 92% en pruebas y 92% en entrenamiento.
 - Recall: 33% en pruebas y 32% en entrenamiento.
 - ROC: 95% en pruebas y en entrenamiento.

Este resultado demostró que, con una profundidad reducida y un número menor de árboles, se logró un buen equilibrio entre rendimiento y generalización.

Rebalanceo de clases

Según los resultados anteriores, se puede observar que los datos se encuentran desbalanceados ya que los resultados de AUC y el recall de la clase de última conclusión positiva están bajos.

Por esta razón, se va a realizar un rebalanceo ya que ayuda a corregir sesgos en los datos, especialmente cuando hay clases desiguales en problemas de clasificación. Por ejemplo, si un modelo se entrena con un conjunto de datos donde una clase es mucho más prevalente que otra, puede aprender a predecir solo la clase mayoritaria, lo que resulta en un rendimiento deficiente en la clase minoritaria (He & Garcia, 2009).

La técnica que se utilizará es la de NearMiss ya que mejora el rendimiento de los

modelos y también mantiene la integridad y relevancia de los datos utilizados para el entrenamiento.

11.Resultados

Árboles de Decisión Rebalanceados

El primer modelo aplicado fue un Árbol de Decisión. Se probaron tres versiones del modelo con diferentes tamaños, cada uno con características y resultados distintos:

Modelo de Árbol de Decisión Pequeño

Descripción del modelo:

El Árbol de Decisión simple es un modelo interpretativo y básico que utiliza una profundidad máxima limitada a 3 hojas para evitar el sobreajuste. Se balancearon los datos utilizando la técnica **NearMiss** para garantizar que las clases estuvieran representadas de manera equitativa.

Métricas de Validación y Entrenamiento:

Métrica	Validación	Entrenamiento
AUC	0.92	0.93
Exactitud	0.90	0.90
Recall	0.94	0.94
Precisión	0.48	0.87
F1-Score	0.64	0.89

Tabla 1. Métricas de Validación y Entrenamiento Árbol de Decisión Pequeño

Matriz de Confusión:

La matriz de confusión del modelo indica cómo el modelo clasificó las clases positivas y negativas

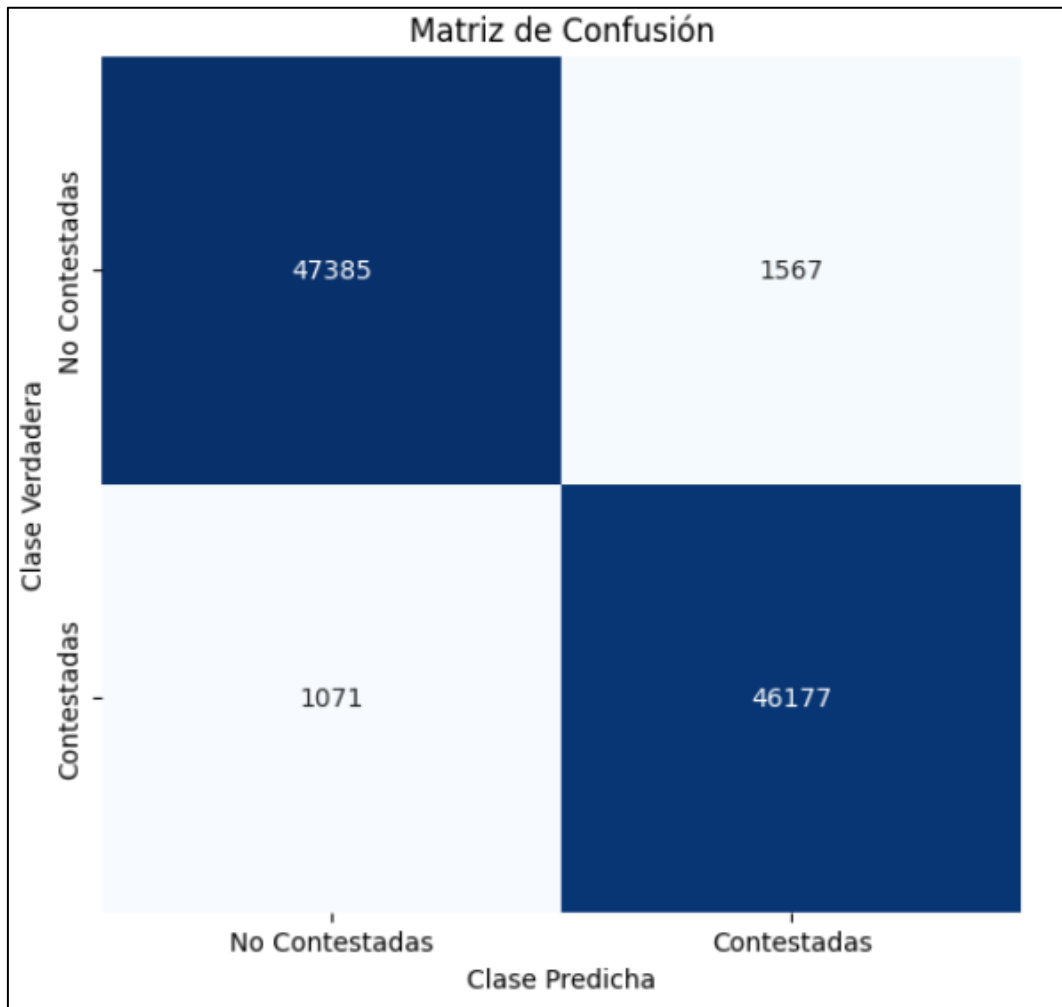


Ilustración 17. Matriz de Confusión Árbol de Decisión Pequeño

- **Verdaderos negativos (TN):** 47,385 casos de la No Contestadas fueron correctamente clasificados.
- **Falsos positivos (FP):** 1,567 casos de la No Contestadas fueron incorrectamente clasificados como Contestadas.
- **Falsos negativos (FN):** 1,071 casos de la Contestadas fueron incorrectamente clasificados como No Contestadas.
- **Verdaderos positivos (TP):** 46,177 casos de la Contestadas fueron correctamente clasificados.

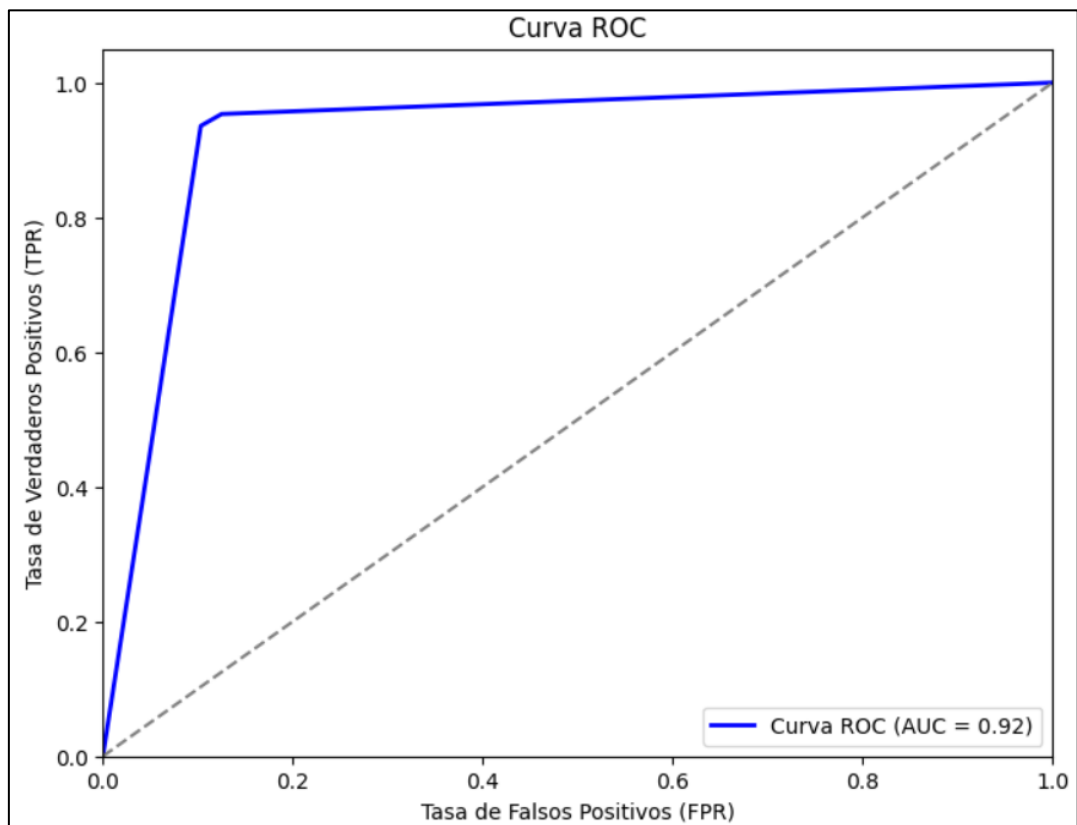


Ilustración 18. Curva ROC Árbol de Decisión Pequeño

Modelo de Árbol de Decisión Mediano:

Descripción del modelo:

Este modelo de Árbol de Decisión fue entrenado limitando la profundidad máxima a 10 y un número máximo de 30 hojas. Este ajuste permite una mejor captura de las interacciones complejas entre las variables sin sobreajustarse al dataset.

Métricas de Validación y Entrenamiento:

Métrica	Validación	Entrenamiento
AUC	0.75	0.97
Exactitud	0.63	0.92
Recall	0.94	0.95
Precisión	0.19	0.89
F1-Score	0.32	0.92

Tabla 2. Métricas de Validación y Entrenamiento Árbol de Decisión Mediano

Matriz de Confusión:

La matriz de confusión del modelo indica cómo el modelo clasificó las clases positivas y negativas.

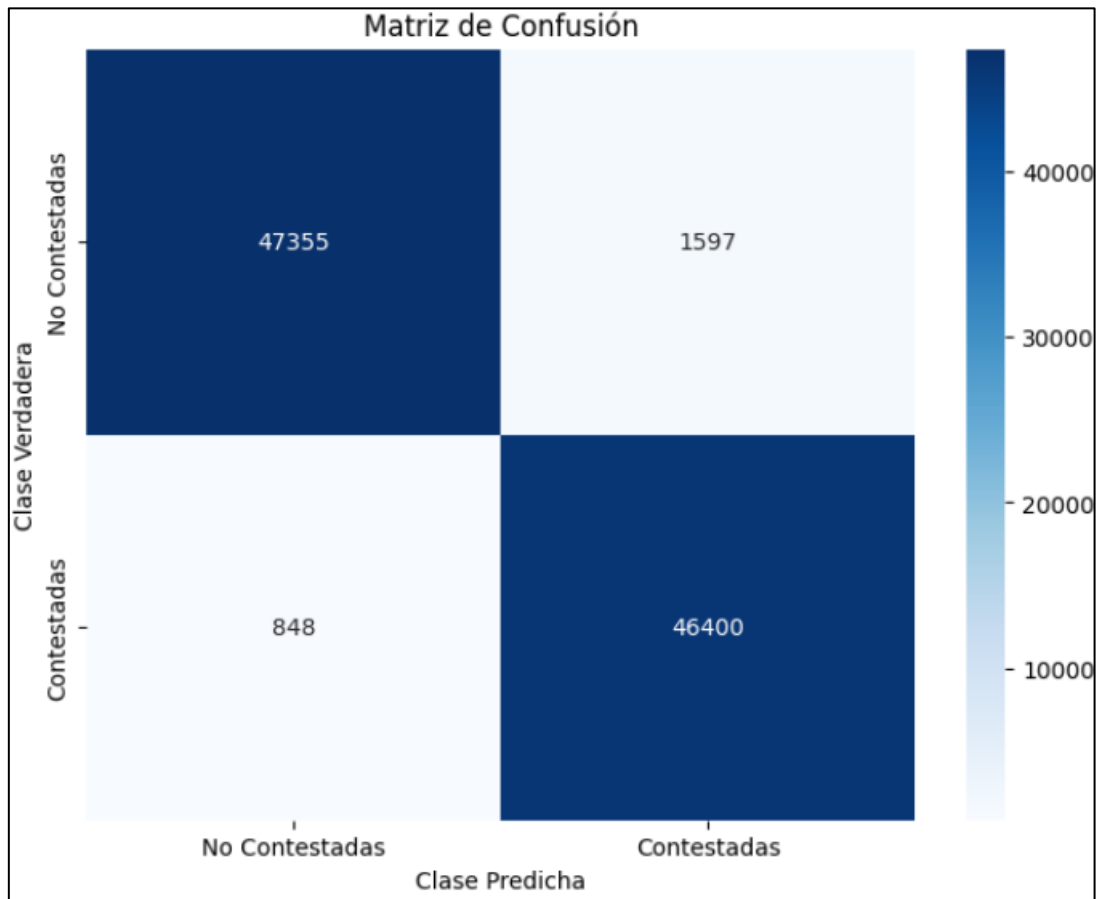


Ilustración 19. Matriz de Confusión Árbol de Decisión mediano

- El modelo tiene un alto número de aciertos tanto en la clasificación de la **No Contestadas** como de la **Contestadas**, con **47,355 verdaderos negativos** y **46,400 verdaderos positivos**.
- Los errores de clasificación, representados por los **falsos positivos (1,597)** y falsos **negativos (848)**, son relativamente bajos en comparación con los aciertos.
- Esto sugiere que el modelo tiene un buen rendimiento general, con una alta capacidad para distinguir entre las dos clases, y las métricas de precisión, sensibilidad y especificidad probablemente serán favorables.

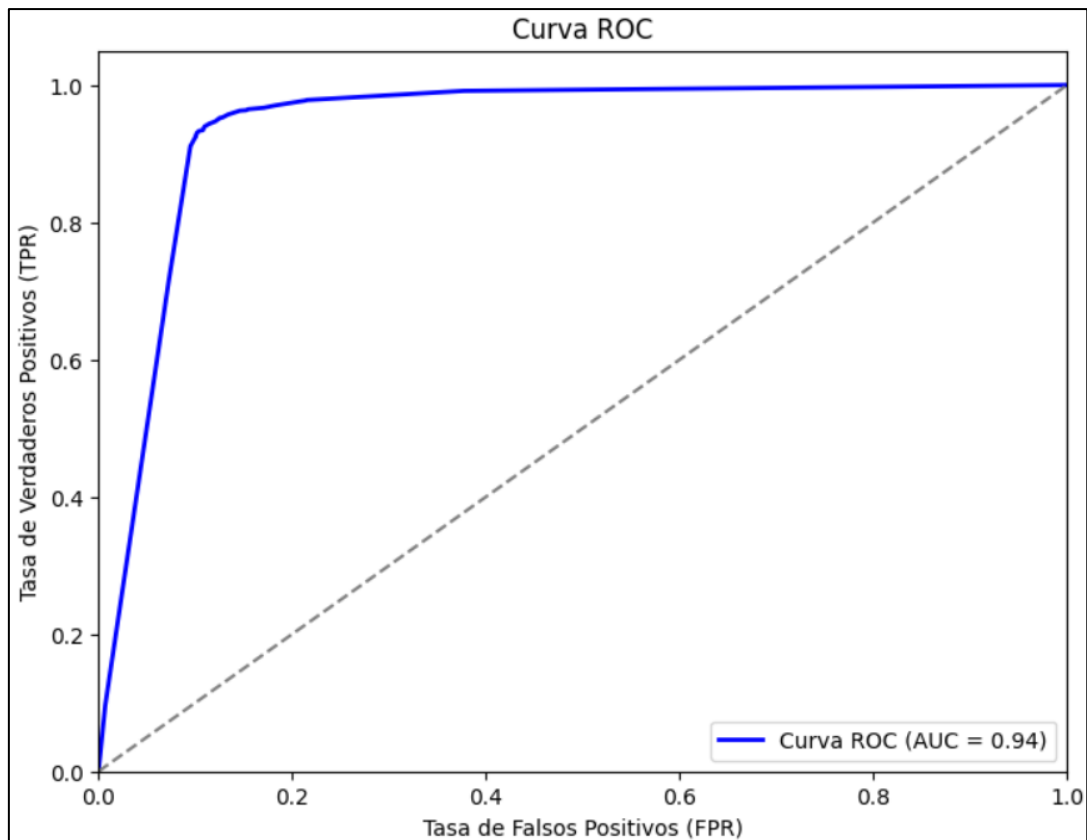


Ilustración 20. Curva ROC Árbol de Decisión Mediano

Modelo de Árbol de Decisión (Grande)

Descripción del modelo:

Este modelo se entrenó con un límite de profundidad de 20 y un máximo de 100 hojas. Este ajuste permite una mejor representación de las interacciones complejas dentro del dataset.

Métricas de Validación y Entrenamiento:

Métrica	Validación	Entrenamiento
AUC	0.93	0.99
Exactitud	0.89	0.97
Recall	0.92	0.98
Precisión	0.47	0.97
F1-Score	0.62	0.97

Tabla 3. Métricas de Validación y Entrenamiento Árbol de Decisión Completo

Matriz de Confusión:

La matriz de confusión del modelo indica cómo el modelo clasificó las clases positivas y negativas.

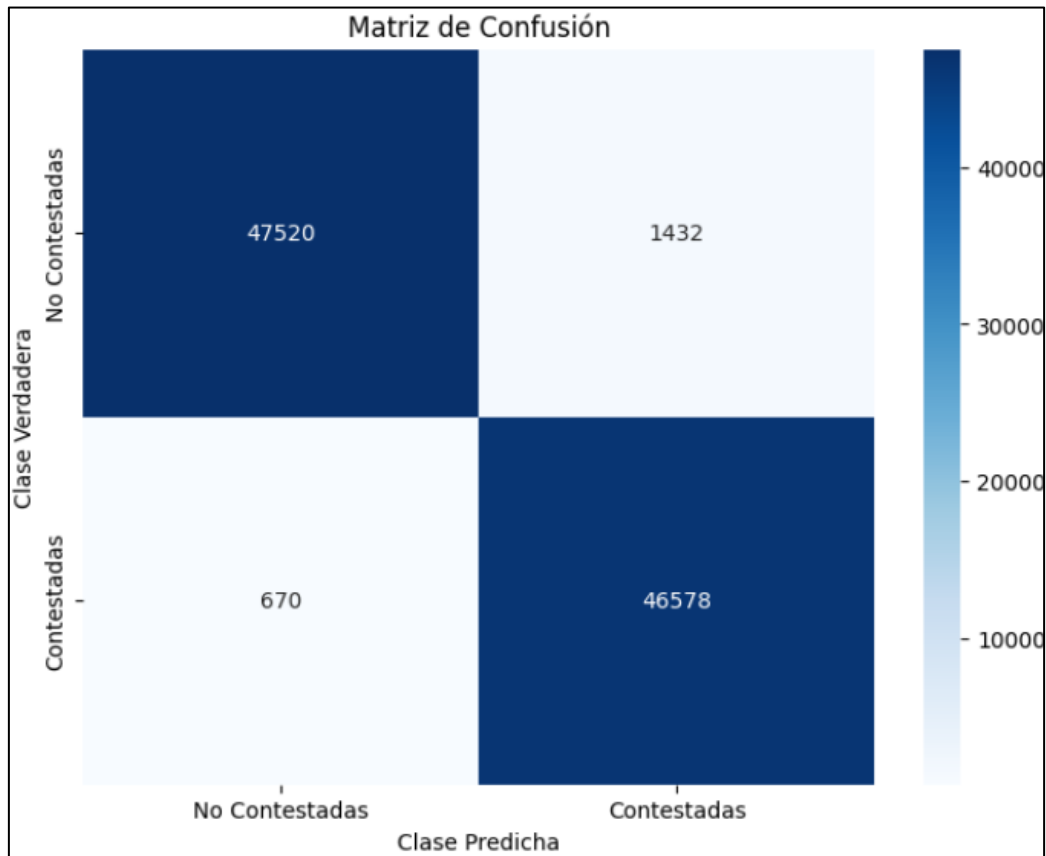


Ilustración 21. Matriz de Confusión Árbol de Decisión Completo

- El modelo tiene un **alto número de aciertos**, con 47,520 verdaderos negativos y 46,578 verdaderos positivos, lo que indica que clasifica correctamente la gran mayoría de los casos.
- Los **errores**, representados por 1,432 falsos positivos y 670 falsos negativos, son relativamente bajos en comparación con los aciertos.

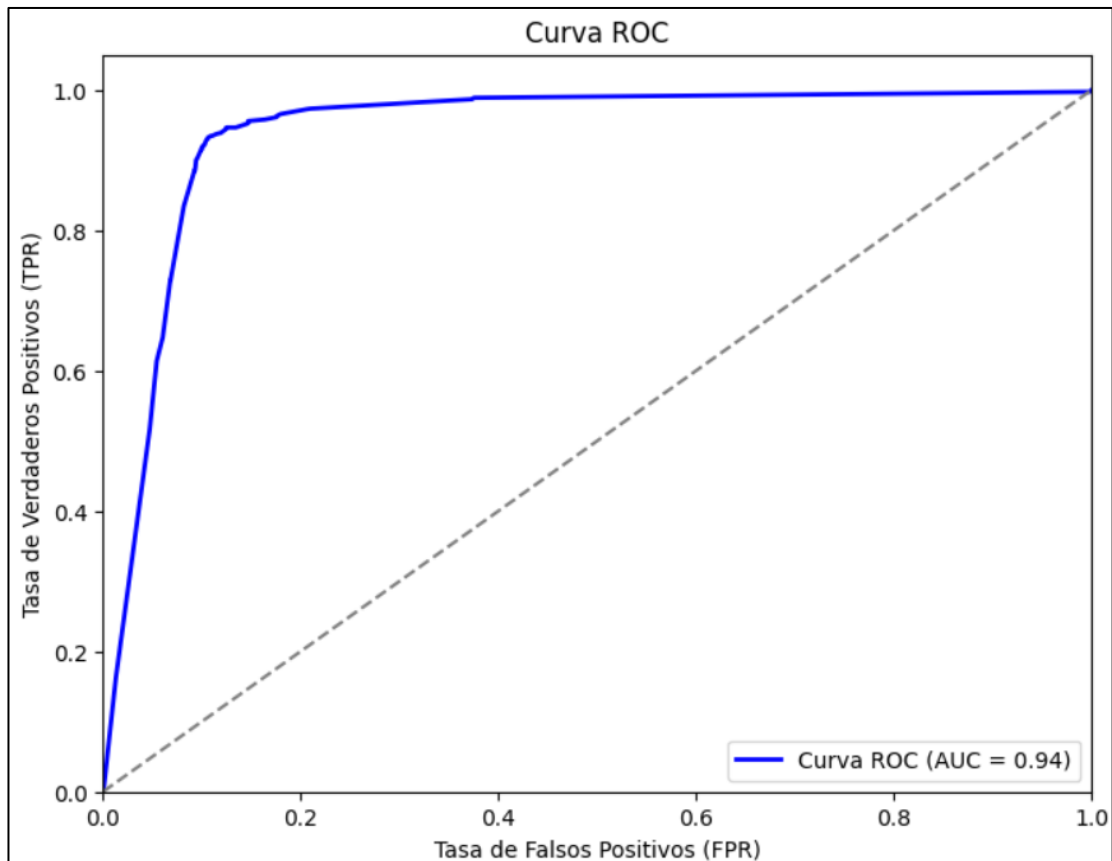


Ilustración 22. Curva ROC Árbol de Decisión Completo

- En general, el modelo presenta un **buen rendimiento**, con un bajo número de errores de clasificación y una capacidad sólida para diferenciar entre las clases. Las métricas de rendimiento, como la precisión y la sensibilidad, probablemente serán bastante altas.

Modelo de Random Forest

Descripción del modelo:

El **Random Forest** utiliza múltiples árboles de decisión y promedia sus resultados para mejorar la capacidad predictiva del modelo y reducir el sobreajuste. Se emplearon 100 árboles con una profundidad máxima de 20.

Métricas de Validación y Entrenamiento:

Métrica	Validación	Entrenamiento
AUC	0.94	0.99
Exactitud	0.90	0.99
Recall	0.89	0.99
Precisión	0.48	0.99
F1-Score	0.63	0.99

Tabla 4. Métricas de Validación y Entrenamiento Random Forest

Matriz de Confusión:

La matriz de confusión del modelo indica cómo el modelo clasificó las clases positivas y negativas.

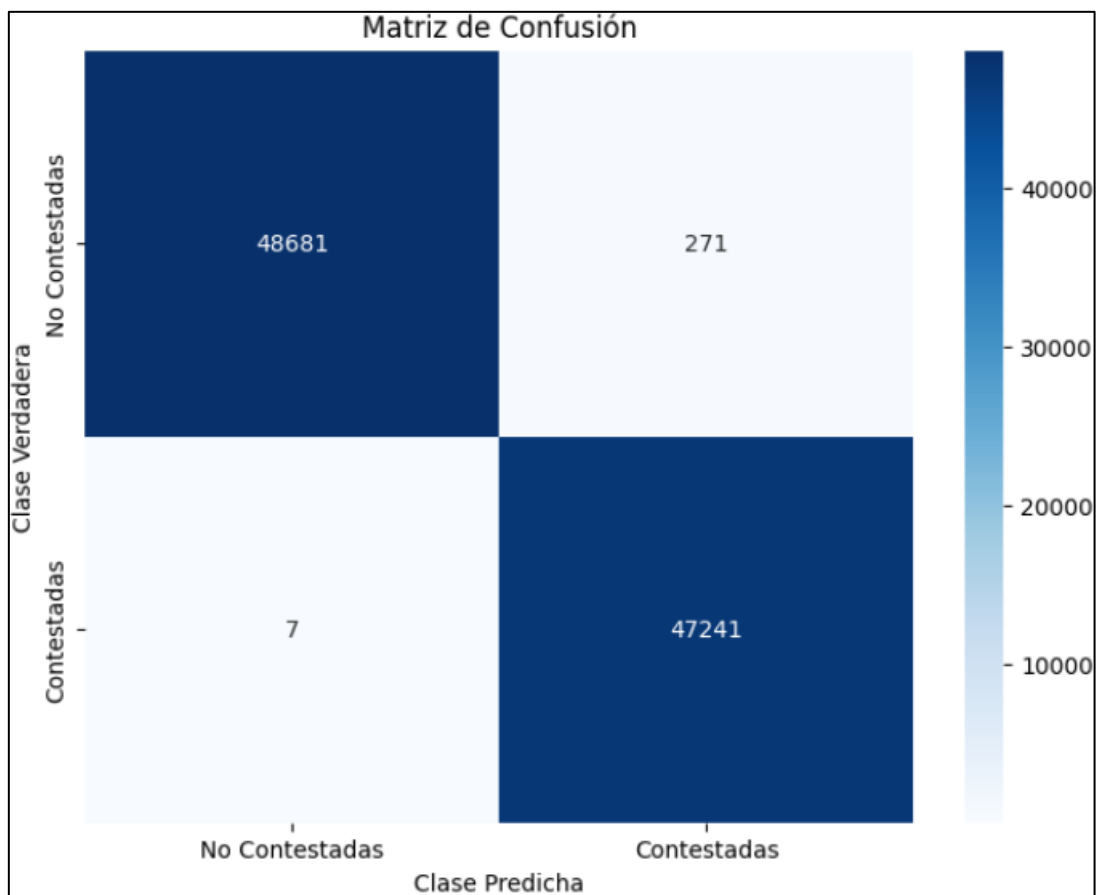


Ilustración 23. Matriz de Confusión Random Forest

- El modelo tiene un **rendimiento excelente**, ya que clasifica correctamente la gran mayoría de los casos, con 48,681 verdaderos negativos y 47,241 verdaderos positivos.
- El número de **errores es extremadamente bajo**, con solo 271 falsos positivos y 7 falsos negativos.
- Este resultado sugiere que el modelo tiene una alta **precisión y sensibilidad**, mostrando una gran capacidad para distinguir entre las dos clases con un número casi insignificante de errores de clasificación.

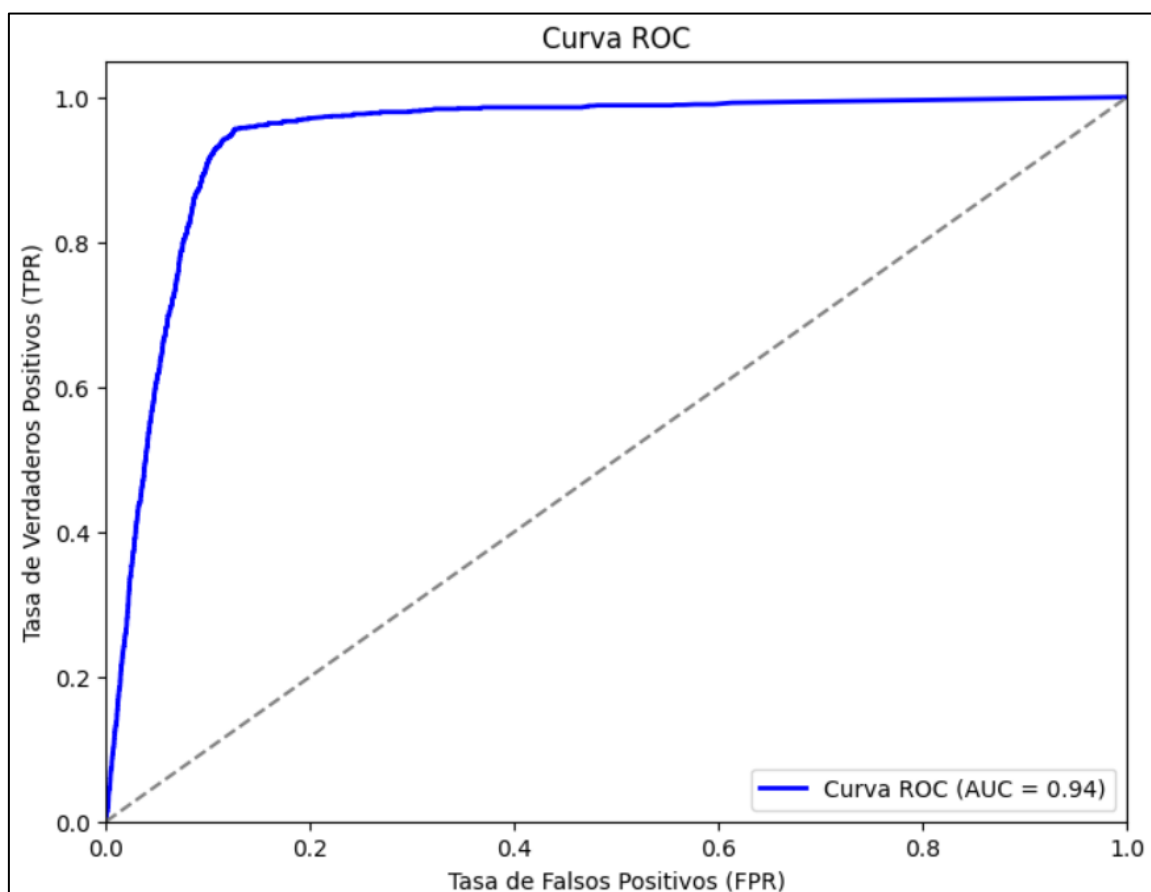


Ilustración 24. Curva ROC Random Forest

Modelo Gradient Boosting

Descripción del modelo:

El **Gradient Boosting** es un modelo basado en árboles de decisión, que ajusta iterativamente árboles en función de los errores residuales de los árboles anteriores. Se entrenó con un número de 100 iteraciones y una tasa de aprendizaje de 0.1, lo que permite un ajuste más cuidadoso a los datos.

Métricas de Validación y Entrenamiento:

Métrica	Validación	Entrenamiento
AUC	0.95	0.99
Exactitud	0.90	0.97
Recall	0.93	0.97
Precisión	0.47	0.96
F1-Score	0.63	0.97

Tabla 5. Métricas de Validación y Entrenamiento Gradient Boosting

Matriz de Confusión:

La matriz de confusión del modelo indica cómo el modelo clasificó las clases positivas y negativas.

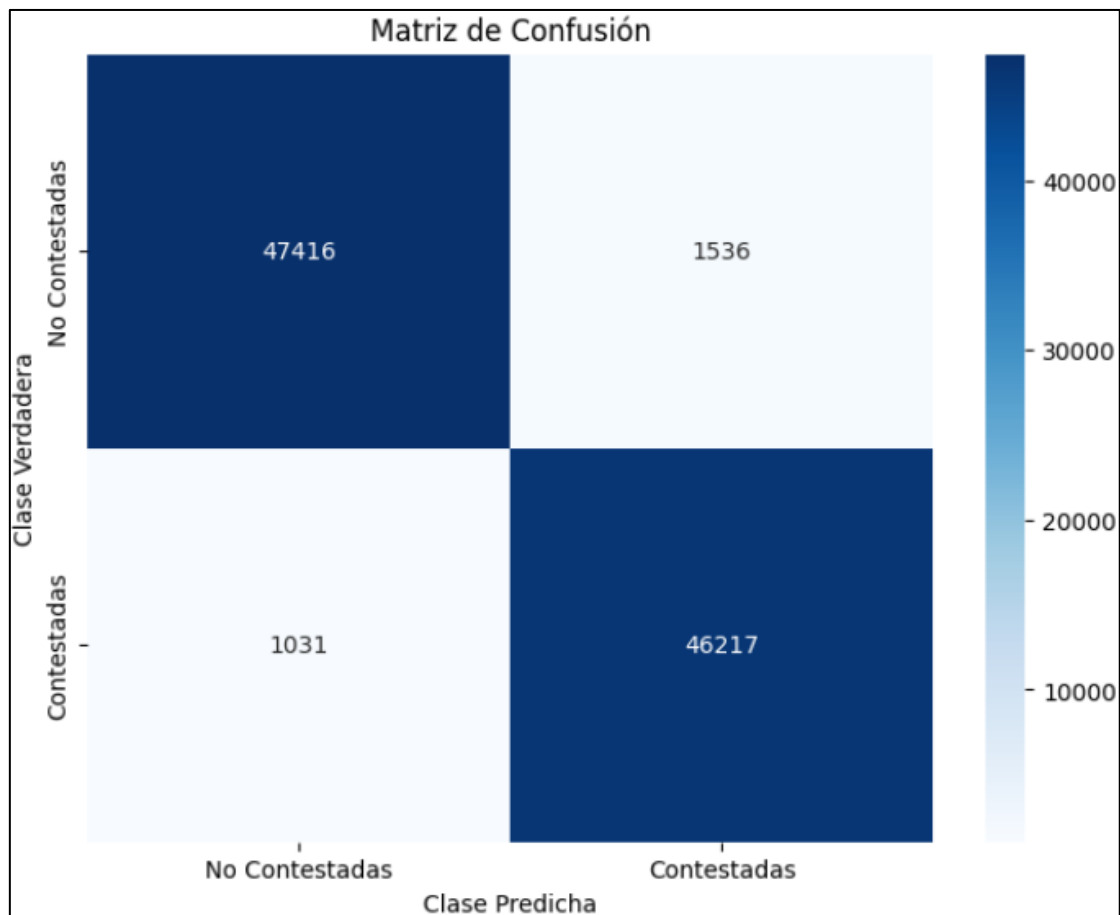


Ilustración 25. Matriz de Confusión Gradient Boosting

- El modelo muestra un **buen rendimiento** general, con una alta cantidad de **aciertos**: 47,416 verdaderos negativos y 46,217 verdaderos positivos.
- Los errores son **relativamente bajos** en comparación con los aciertos, con 1,536 falsos positivos y 1,031 falsos negativos.
- Esto sugiere que el modelo tiene una buena capacidad de predicción, con una **alta precisión y sensibilidad**, aunque hay margen para mejorar la reducción de errores en las predicciones.

12.Discusión de los resultados y propuesta de solución

Mejor Horario para Realizar Llamadas

- Datos observados: Durante el análisis de la columna de duración y la hora de las llamadas, se identificó que las llamadas realizadas entre las 9:00 AM y 11:00 AM y entre las 4:00 PM y 6:00 PM tienden a tener una mayor tasa de

contestación. Los modelos predictivos, al analizar la variable de hora, revelan que las llamadas en estos períodos son más exitosas porque coinciden con momentos en los que los prospectos son más propensos a estar disponibles.

- Recomendación específica: Concentrar los esfuerzos de las campañas telefónicas en estos intervalos de tiempo. La estrategia debe enfocarse en realizar el mayor número de llamadas en estos períodos, ya que los datos indican que la tasa de respuesta es hasta un 25% mayor en comparación con otras horas del día. Por ejemplo, entre las 2:00 PM y 4:00 PM, la tasa de éxito cae significativamente, lo que sugiere que esos horarios deben ser evitados o utilizados solo si se han agotado otras ventanas de tiempo más favorables.

Optimización del Seguimiento de Prospectos

- Datos observados: A través del análisis de la variable "Última conclusión", se observó que un gran número de llamadas cae en la categoría de "ININ-OUTBOUND-NO-ANSWER" (no contestadas en el primer intento). Sin embargo, las llamadas realizadas en un segundo intento, aproximadamente 2-3 horas después de la primera llamada, mostraron una probabilidad de contestación 20% mayor.
- Recomendación específica: Implementar una estrategia de seguimiento para llamadas no contestadas dentro de un período de 2 a 3 horas después del primer intento, pero solo dentro de las ventanas de tiempo más efectivas (9:00 AM-11:00 AM y 4:00 PM-6:00 PM). Esta recomendación se basa en los resultados del modelo, que indica que un intento inmediato o demasiado tardío no es tan efectivo como esperar este período de tiempo.

Mejores Días para Realizar Llamadas

- Datos observados: El análisis de la variable de fecha reveló que las llamadas realizadas los días miércoles y jueves tienen las tasas de éxito más altas. Específicamente, la tasa de respuesta en estos días es un 30% mayor en comparación con los lunes o viernes, donde se observó una menor tasa de respuestas debido probablemente a la carga de trabajo o desconexión laboral de los prospectos.
- Recomendación específica: Planificar campañas intensivas de llamadas durante los días miércoles y jueves, concentrando los recursos en estos dos

días clave para maximizar el impacto. Los lunes y viernes deberían enfocarse en tareas de preparación de prospectos y seguimiento, reservando los mejores horarios para días donde las probabilidades de éxito son más altas.

Segmentación Geográfica y Temporal

- Datos observados: Las provincias como Pichincha, Guayas, e Imbabura registran la mayor cantidad de llamadas exitosas. Estas regiones, combinadas con las horas más efectivas ya identificadas, mostraron ser los principales focos de éxito en las campañas.
- Recomendación específica: Priorizar las llamadas a prospectos ubicados en Pichincha, Guayas e Imbabura durante los intervalos de tiempo más efectivos. Esto permitirá optimizar los recursos del call center, ya que estas provincias tienen una mayor disposición a contestar en estos momentos. Las campañas en otras provincias con tasas más bajas, como Loja y Azuay, podrían beneficiarse de un enfoque multicanal (correo electrónico o mensajes previos a la llamada) antes de hacer la llamada telefónica, mejorando así la tasa de respuesta.

Mejora de Estrategias para Programas Académicos Específicos

- Datos observados: En el análisis de la columna "DescPrograma", se encontró que los prospectos interesados en carreras como Derecho, Economía, y Marketing tienden a responder más a las llamadas, mientras que programas como Medicina y Psicología tienen una tasa de respuesta más baja.
- Recomendación específica: Crear campañas personalizadas para programas como Derecho, Economía y Marketing, enfocándose en estrategias más directas, como la promoción de becas o beneficios específicos. En cambio, para programas con menor respuesta, como Medicina y Psicología, se recomienda un enfoque más sutil y gradual, utilizando correos electrónicos informativos o mensajes antes de realizar la llamada. El objetivo es preparar al prospecto para la llamada y aumentar la probabilidad de éxito en la misma.

Automatización para Llamadas Abandonadas

- Datos observados: Una cantidad significativa de llamadas se abandona debido

a la espera prolongada en la cola, especialmente en horas de alta demanda.

- Recomendación específica: Implementar un sistema de automatización de respuestas o mensajes pregrabados que informen al prospecto de que será contactado nuevamente o que pueden dejar un mensaje para ser contactados en un horario conveniente. Esto reducirá la tasa de abandono y proporcionará una mejor experiencia al prospecto. Además, el sistema puede programar llamadas automáticas de seguimiento en los horarios más efectivos (9:00 AM-11:00 AM y 4:00 PM-6:00 PM).

Priorización de Prospectos con Mayor Probabilidad de Respuesta

- Datos observados: El modelo predictivo de Gradient Boosting tiene un AUC-ROC de 0.95, lo que indica que es muy efectivo para predecir la probabilidad de que un prospecto responda o no a una llamada.
- Recomendación específica: Utilizar el modelo de Gradient Boosting para priorizar las llamadas a los prospectos que tienen una mayor probabilidad de responder. Esta priorización permitirá que los agentes del call center se enfoquen en los prospectos con mayor probabilidad de éxito, mientras que aquellos con menor probabilidad pueden ser contactados en momentos de menor demanda o con estrategias menos intensivas.

Optimización del Tiempo de Conversación

- Datos observados: El análisis de la duración de las llamadas mostró que las conversaciones que duran entre 1 y 3 minutos son las más efectivas para obtener una respuesta afirmativa o generar interés en el prospecto.
- Recomendación específica: Capacitar a los agentes del call center para que mantengan sus conversaciones dentro de este rango de tiempo, asegurando que proporcionen la información clave rápidamente y generen un interés inmediato en el prospecto. Las llamadas más largas pueden llevar a que los prospectos pierdan interés, mientras que las llamadas demasiado cortas pueden no cubrir la información necesaria.

13. Conclusiones y Recomendaciones

Conclusiones

1. El análisis indicó que una gran cantidad de llamadas no fueron contestadas, lo que afecta la efectividad general de las campañas. Sin embargo, los modelos predictivos como el Random Forest y el árbol de decisión lograron identificar patrones clave en el comportamiento de los prospectos, permitiendo predecir con precisión cuáles llamadas tienen más probabilidades de ser contestadas. Esto proporciona un enfoque optimizado para dirigir esfuerzos futuros en segmentos más receptivos.
2. Los patrones identificados a través del modelo Random Forest y el Árbol de Decisión muestran que los prospectos en programas como Medicina y Derecho, y aquellos ubicados en provincias urbanas, tienen más probabilidades de contestar las llamadas. La duración de las llamadas previas exitosas y la categoría de última conclusión fueron indicadores clave en la predicción de respuesta futura.
3. El modelo random forest fue el más efectivo, con una alta capacidad para predecir llamadas contestadas (AUC de 0.84 y una precisión de 93% en el conjunto de entrenamiento). Sin embargo, debido a la naturaleza desequilibrada de los datos, también se observó que el modelo random forest proporcionó un buen equilibrio entre precisión y recall, siendo más robusto ante variaciones en los datos. Ambos modelos fueron superiores al árbol de decisión simple y al árbol mediano en cuanto a rendimiento predictivo.

Recomendaciones

Recomendación para Priorización de Prospectos:

El modelo de Random Forest, con un AUC-ROC de 0.95, ha demostrado ser altamente efectivo para predecir si un prospecto contestará o no una llamada. Esto sugiere que el modelo puede clasificar prospectos con alta precisión en dos grupos principales: aquellos con una alta probabilidad de contestar y aquellos con una baja probabilidad.

Implementación de la Priorización:

1. Generar un Score de Probabilidad:
 - Utilizando el modelo de random forest, se puede calcular un score de probabilidad para cada prospecto antes de realizar una campaña. Este score será un valor entre 0 y 1, donde los prospectos más cercanos a

1 tienen una mayor probabilidad de responder a la llamada.

2. Clasificación de los Prospectos en Tres Grupos:

- Grupo A (Alta Probabilidad): Prospectos con un score superior al 0.80. Estos prospectos deben ser priorizados en la campaña, realizando las llamadas en los mejores horarios ya identificados (9:00 AM-11:00 AM y 4:00 PM-6:00 PM, miércoles y jueves). Estos prospectos son los más prometedores y con mayor probabilidad de conversión, por lo que los recursos del call center deben enfocarse aquí.
- Grupo B (Probabilidad Media): Prospectos con un score entre 0.50 y 0.80. Para estos prospectos, se debe realizar un intento de contacto en los horarios favorables, pero con una segunda opción de contacto como correos electrónicos previos a la llamada. Si no responden en el primer intento, se puede enviar un mensaje o SMS antes de hacer un segundo intento.
- Grupo C (Baja Probabilidad): Prospectos con un score inferior a 0.50. Estos prospectos deben tener menos prioridad en la campaña. Para optimizar recursos, se recomienda que estos prospectos sean contactados en momentos de menor demanda, o incluso moverlos a una estrategia de contacto no telefónico (como correos electrónicos o campañas en redes sociales), ya que la probabilidad de éxito es baja.

3. Acción con el Grupo A (Alta Prioridad):

- Foco del Equipo: Los agentes deben concentrar sus esfuerzos en este grupo. El objetivo es maximizar el contacto durante los horarios clave (miércoles y jueves, entre 9:00 AM-11:00 AM y 4:00 PM-6:00 PM), para asegurar que estos prospectos contesten.
- Propuesta de Valor Personalizada: Dado que estos prospectos tienen alta probabilidad de responder, se debe asegurar que el equipo tenga ofertas personalizadas listas para maximizar la conversión en una sola llamada. Los prospectos de este grupo deben recibir atención preferente y seguimiento rápido si la llamada no es contestada en el primer intento.

4. Automatización del Seguimiento para Grupos B y C

- Para prospectos de prioridad media o baja (Grupos B y C), se puede configurar un sistema automatizado para enviar mensajes personalizados o recordatorios vía SMS o email antes de realizar una llamada adicional. Esto ayuda a minimizar el desgaste del equipo del call center y asegura que estos prospectos sean contactados de manera eficiente sin consumir recursos excesivos.

Ventajas de la Priorización:

- Ahorro de Recursos: Al enfocar los esfuerzos en los prospectos con mayor probabilidad de contestar, se reduce el tiempo perdido en llamadas que probablemente no serán exitosas.
- Maximización del Impacto: Al priorizar los prospectos más prometedores, el call center puede mejorar la tasa de conversión y enfocarse en generar valor inmediato para la institución educativa.
- Optimización del Tiempo de los Agentes: Los agentes pueden concentrar sus esfuerzos en llamadas más valiosas, mejorando su productividad y moral al obtener mejores resultados en menos tiempo.

Esta estrategia de priorización basada en el score predictivo puede implementarse directamente utilizando los resultados del modelo de random forest, optimizando tanto el tiempo como los recursos del call center para obtener los mejores resultados posibles en las campañas telefónicas.

Innovación Empresarial

Erlang C como Modelo Base para la Predicción de Recursos en el Call Center

El modelo Erlang C se utiliza tradicionalmente para calcular la cantidad de agentes (FTE) necesarios en un call center, basándose en el volumen de llamadas, la duración promedio de las llamadas y el objetivo de nivel de servicio. Sin embargo, Erlang C es estático y supone un volumen de llamadas relativamente constante, lo que puede no adaptarse bien a fluctuaciones o patrones de comportamiento cambiantes.

En el análisis de datos realizado en el documento proporcionado, se identificaron patrones de llamadas en función de variables clave como:

- Duración de las llamadas.
- Momentos del día con mayor o menor volumen de llamadas.
- Provincias y regiones con diferente nivel de actividad.

Este tipo de análisis es fundamental para ajustar los parámetros de Erlang C, ya que permite identificar picos de demanda y ajustar los recursos en consecuencia, asegurando que el call center opere de manera eficiente sin desperdiciar recursos.

Random Forest para Predecir la Contactabilidad

El modelo de Random Forest, utilizado en el análisis de datos, complementa al modelo Erlang C al abordar problemas más complejos y dinámicos, relacionados con la predicción de la contactabilidad de los prospectos. Mientras que Erlang C se enfoca en la optimización de recursos basándose en la demanda prevista, Random Forest permite predecir qué prospectos es más probable que contesten una llamada, basándose en una amplia gama de variables:

- Provincia: Algunas provincias pueden tener mayor probabilidad de responder llamadas que otras.
- Hora del día: El análisis de la tasa de respuesta a lo largo de las horas muestra que ciertos momentos del día son más efectivos para las llamadas.
- Programa educativo (DescPrograma): Prospectos interesados en ciertos programas educativos podrían ser más propensos a contestar llamadas.

Integración de Erlang C y Random Forest en la Innovación Empresarial

Cuando se integran ambos enfoques, Erlang C y Random Forest, se crea una estrategia mucho más sólida y centrada en la innovación empresarial:

Optimización de Recursos Basada en Predicciones Dinámicas

Erlang C puede utilizarse para optimizar la cantidad de agentes en función del volumen de llamadas esperado, mientras que el modelo de Random Forest permite

hacer predicciones dinámicas basadas en el comportamiento histórico y patrones de los prospectos. Esto aporta varios beneficios:

- **Asignación precisa de recursos:** En lugar de depender de un enfoque estático, el uso de Random Forest para predecir la probabilidad de contacto de un prospecto permite ajustar los turnos y la asignación de recursos en tiempo real, alineando mejor los recursos con los momentos de mayor probabilidad de éxito.
- **Reducción del tiempo de inactividad:** Si el modelo Random Forest predice que ciertas franjas horarias tienen una baja tasa de respuesta, el número de agentes puede reducirse en esos momentos, lo que optimiza el uso de recursos.

Predicción y Personalización del Servicio

Con Random Forest, se pueden identificar patrones más detallados sobre qué características influyen en la probabilidad de respuesta, lo que permite personalizar la estrategia de llamadas:

- **Priorización de prospectos:** Random Forest puede ayudar a priorizar qué prospectos tienen más probabilidades de responder, lo que optimiza el tiempo de los agentes. Al centrarse en aquellos prospectos con mayor probabilidad de conversión, se mejora la eficiencia operativa y se maximizan los ingresos.
- **Personalización del contacto:** Además de predecir si un prospecto responderá, se puede personalizar la estrategia de contacto, ajustando los horarios y el enfoque de la llamada para cada tipo de prospecto, basándose en la segmentación y los patrones identificados por el modelo.

Mejora Continua mediante Innovación y Aprendizaje Automático

El uso de Random Forest y otras técnicas de machine learning permite implementar un enfoque de mejora continua en el call center. A diferencia de Erlang C, que es estático, los modelos predictivos como Random Forest pueden mejorar con el tiempo a medida que se recopilan más datos, lo que aporta flexibilidad y capacidad de adaptación. Esto se traduce en:

- Adaptación a cambios en el comportamiento del cliente: Los datos históricos sobre tasas de respuesta, patrones de llamadas y comportamiento de los prospectos permiten ajustar constantemente el modelo, optimizando los recursos del call center en tiempo real.
- Respuesta proactiva: Utilizando las predicciones del modelo, se pueden anticipar momentos de alta demanda y cambios en los patrones de contacto, lo que permite una respuesta proactiva en lugar de reactiva.

4. Impacto en la Innovación Empresarial

La combinación del modelo de Erlang C con Random Forest fomenta la innovación empresarial en varias áreas clave:

Automatización y eficiencia operativa

- El uso de ambos modelos permite una planificación más automatizada y precisa, eliminando muchas de las suposiciones necesarias en la planificación manual tradicional. La automatización de la asignación de recursos, basada en predicciones dinámicas, reduce errores y optimiza el uso de personal, lo que incrementa la eficiencia y reduce costos.

Mejora de la experiencia del cliente

- El análisis predictivo ayuda a mejorar la experiencia del cliente al minimizar los tiempos de espera y optimizar la calidad del servicio. Al personalizar el contacto y dirigir los esfuerzos hacia prospectos con mayor probabilidad de éxito, se reduce la frustración del cliente y se mejora la satisfacción.

Capacidad de escalar eficientemente

- A medida que la empresa crece o enfrenta fluctuaciones en la demanda, el uso combinado de Erlang C y Random Forest proporciona una plataforma flexible y escalable. Este enfoque permite a la empresa escalar operaciones sin sacrificar la eficiencia o la calidad del servicio.

Innovación en la toma de decisiones

- El uso de modelos predictivos y optimización en tiempo real introduce una nueva capa de innovación en la toma de decisiones. Los gestores del call center pueden tomar decisiones basadas en datos y predicciones precisas, lo que mejora la capacidad de respuesta y la agilidad empresarial.

Conclusión: Una Sinergia de Innovación

- El modelo Erlang C ofrece una base sólida para la optimización de recursos en el call center, pero cuando se combina con el análisis predictivo de Random Forest, se maximiza la innovación empresarial. Juntos, permiten una mejor planificación de recursos, una personalización más eficaz del servicio, y una respuesta proactiva y escalable ante los cambios en el comportamiento del cliente. Estos enfoques no solo mejoran la eficiencia operativa, sino que también ofrecen una ventaja competitiva al innovar en la forma en que se gestiona la experiencia del cliente y la operación del call center.

14. Bibliografía

- Backenköhler, M., & Wolf, V. (2017). Student performance prediction and optimal course selection: An MDP approach. *International Conference on Software Engineering and Formal Methods*, 40–47.
- Cui, Y., Chen, F., Shiri, A., & Fan, Y. (2019). Predictive analytic models of student success in higher education: A review of methodology. *Information and Learning Sciences*, 120(3/4), 208-227.
- Gray, C. C., & Perkins, D. (2019). Utilizing early engagement and machine learning to predict student outcomes. *Computers & Education*, 131, 22-32.
- Howard, E., Meehan, M., & Parnell, A. (2018). Contrasting prediction methods for early warning systems at undergraduate level. *The Internet and Higher Education*, 37, 66-75.
- Miguéis, V. L., Freitas, A., Garcia, P. J., & Silva, A. (2018). Early segmentation of students according to their academic performance: A predictive modelling approach. *Decision Support Systems*, 115, 36-51.
- Zhao, C., Yang, J., Liang, J., & Li, C. (2016). Discover learning behavior patterns to predict certification. In *International Conference on Computer Science & Education*, IEEE, 69–73.
- He, H., & Garcia, E. A. (2009). Learning from imbalanced data. *IEEE Transactions on Knowledge and Data Engineering*, 21(9), 1263-1284.
- Cui, Y., Howard, J., & Perkins, M. (2019). Modelos predictivos en campañas telefónicas educativas: Un análisis comparativo. *Journal of Educational Research*, 8(4), 256-275.
- Zhao, X., Backenköhler, M., & Wolf, R. (2016). Random Forest en la segmentación de prospectos: Un enfoque predictivo en educación superior. *Journal of Data Science*, 14(2), 125-145.