

Proyecto Final de Machine Learning

“Predicción de renovación de
pólizas de una empresa de
seguros”

Michelle Sophie Rosero Muriel &
Karol Vanessa Vitonco Burbano



Caso de estudio

1

Dataset: *InsuranceCompany.csv*

Para este proyecto trabajamos con el dataset **InsuranceCompany.csv**, proporcionado por una empresa del sector de seguros. Este archivo contiene la información histórica de sus clientes, incluyendo datos de sus pagos, comportamiento de renovación y características demográficas. El objetivo del caso de estudio es **predecir qué clientes tienen mayor probabilidad de renovar su póliza**. Esto permitirá a la empresa mejorar su estrategia comercial y optimizar los incentivos que reciben los agentes, de manera que enfoquen su esfuerzo en los asegurados con mayor potencial de renovación.

El dataset incluye información como:

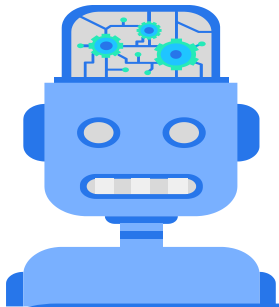
- Número de primas pagadas y primas atrasadas (3, 6 y 12 meses).
- Canal por el que llegó el cliente (sourcing channel).
- Edad, ingresos mensuales y tipo de área donde vive.
- Historial de renovación de pólizas.

Además, el cliente proporcionó datos adicionales sobre:

- Cuántas horas de esfuerzo debe dedicar un agente según los incentivos ofrecidos.
- Cómo ese esfuerzo aumenta la probabilidad de que un cliente renueve.

En resumen, se trata de un caso que combina datos reales del negocio, predicción de comportamiento del cliente y decisiones estratégicas para mejorar los ingresos de la compañía aseguradora.





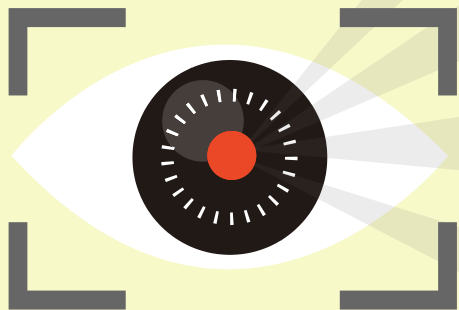
Objetivos del proyecto

Para este proyecto buscamos predecir si un cliente renovará su póliza en el siguiente periodo. La empresa necesita saberlo antes del vencimiento para decidir a quién ofrecer incentivos y así evitar pérdidas por falta de renovación. El modelo estimará esta probabilidad utilizando las relaciones que serán identificadas en el EDA, apoyando la toma de decisiones estratégicas en el plan de incentivos.

Para lograr esto, el proyecto combina varias etapas:

- **Entender el problema de negocio:** Analizar cómo la predicción de renovación ayuda a mejorar los ingresos netos y cómo los incentivos influyen en la decisión del cliente.
- **Explorar los datos (EDA):** Identificar patrones, relaciones y problemas en el dataset: atrasos de pago, canal de adquisición, nivel de ingresos, edad, etc.
- **Preparar los datos para modelar:** Limpiar, transformar, balancear y seleccionar las variables más importantes para que el modelo funcione bien.
- **Entrenar y comparar varios modelos:** Probar diferentes algoritmos (por ejemplo: regresión logística, KNN, modelos lineales, etc.) y elegir el que mejor prediga la renovación.
- **Evaluar el desempeño del modelo:** Usar métricas como Accuracy, F1 y ROC-AUC para comprobar qué tan bien predice y cómo ayuda al negocio.
- **Presentar resultados para la toma de decisiones:** Mostrar de manera clara qué encontró el análisis y cómo la empresa puede usar el modelo para mejorar su plan de incentivos.

**Este análisis nos ayuda
a responder preguntas
como:**



1

¿Qué factores demográficos, transaccionales o de póliza influyen en la renovación?

2

¿Cuáles son los patrones de atrasos en pagos y canales de abastecimiento más comunes asociados con no renovaciones?

3

¿Cómo impactan los ingresos, edad y puntuación de suscripción en la propensión a renovar?

4

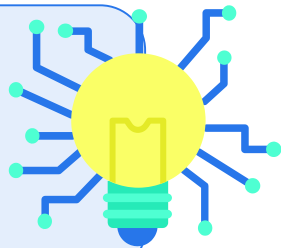
¿Qué variables deben eliminarse o transformarse para optimizar el modelo predictivo y el plan de incentivos?

Dataset

2

Descripción del dataset *insurance_company.csv*

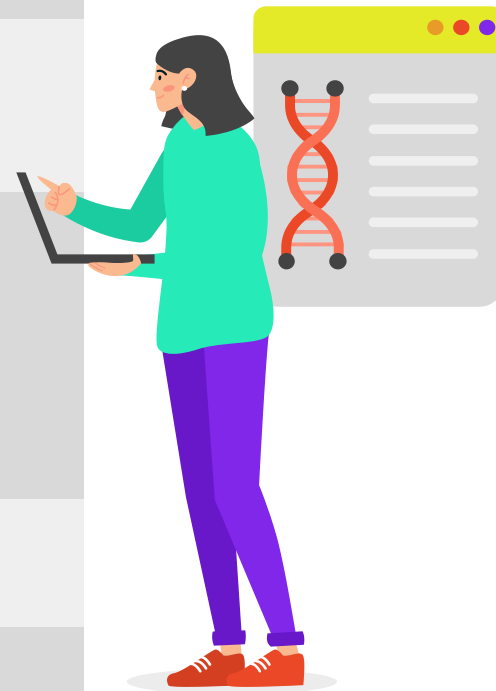
El dataset proviene de una compañía de seguros y contiene **79,853** registros y **13** variables.



Tipos de datos	
Categoricos (string)	<i>sourcing_channel, residence_area_type.</i>
Numéricos (float)	<i>perc_premium_paid_by_cash_credit, application_underwriting_score.</i>
Numéricos (integer)	<i>id, age_in_days, Income, Count_3-6_months_late, Count_6-12_months_late, Count_more_than_12_months_late, no_of_premiums_paid, premium, renewal.</i>

Variables del dataset por categoría

Identificadores	<ul style="list-style-type: none">• id
Demográficas	<ul style="list-style-type: none">• age_in_days• income• residence_area_type
Transaccionales y Pagos	<ul style="list-style-type: none">• perc_premium_paid_by_cash_credit• Count_3-6_months_late• Count_6-12_months_late• Count_more_than_12_months_late• no_of_premiums_paid• premium
Suscripción y Canal	<ul style="list-style-type: none">• application_underwriting_score• sourcing_channel
Variable Objetivo	<ul style="list-style-type: none">• renewal



Variables dependiente e independiente

Dependiente

La variable dependiente es ***renewal***, que indica si el cliente renovó su póliza (1) o no (0).

Es el resultado que buscamos predecir y refleja la decisión final del asegurado frente a su continuidad.

Independientes

Las variables que ayudan a predecir ***renewal*** se agrupan en tres categorías:

Pagos:

- ***perc_premium_paid_by_cash_credit***
- ***Count_3-6_months_late***
- ***Count_6-12_months_late***
- ***Count_more_than_12_months_late***
- ***no_of_premiums_paid***

Suscripción y canal:

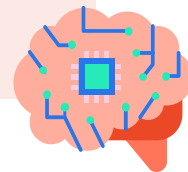
- ***application_underwriting_score***
- ***sourcing_channel***

Demográficas:

- ***age_in_days***
- ***Income***
- ***residence_area_type***

Aspectos a tener en cuenta

- Cero datos duplicados
- Variables como ***id*** fueron eliminadas por baja relevancia.
- Variables con valores faltantes: conteos de atraso (97 faltantes c/u) y ***application_underwriting_score*** (2,974).
- La variable objetivo presenta un desbalance importante (93.74% renovaciones), que será corregido en la fase de balanceo.



Hallazgos Iniciales del Preprocesamiento



- Los tipos de datos "**object**" **deben convertirse a "category"** para optimizar memoria y facilitar el modelado.
- El formato de los **nombres de las columnas debe estandarizarse a minúsculas** para mantener consistencia.
- Se identifican **valores faltantes** en las siguientes variables:
 - *count_3-6_months_late*: 97
 - *count_6-12_months_late*: 97
 - *count_more_than_12_months_late*: 97
 - *application_underwriting_score*: 2974
Estos se resuelven mediante (e.g., imputación con media o mediana).

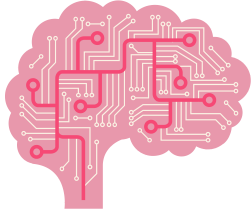
En general, los **tipos de datos numéricos son correctos** (*int64* para enteros, *float64* para decimales), y no requieren modificaciones mayores.

Análisis exploratorio de datos

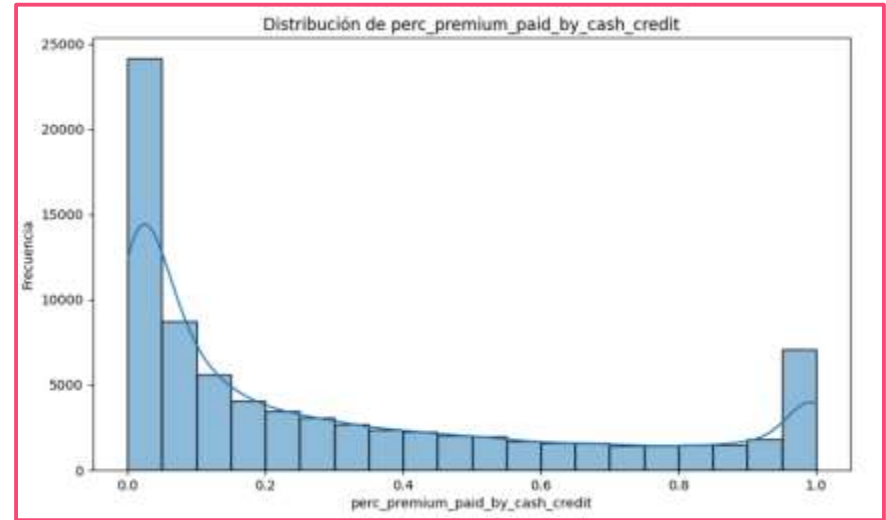
3

Análisis univariable: Variables numericas

3.1



perc_premium_paid_by_cash_credit

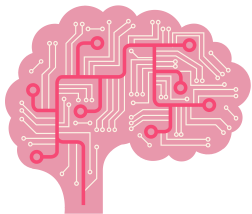


Hallazgos

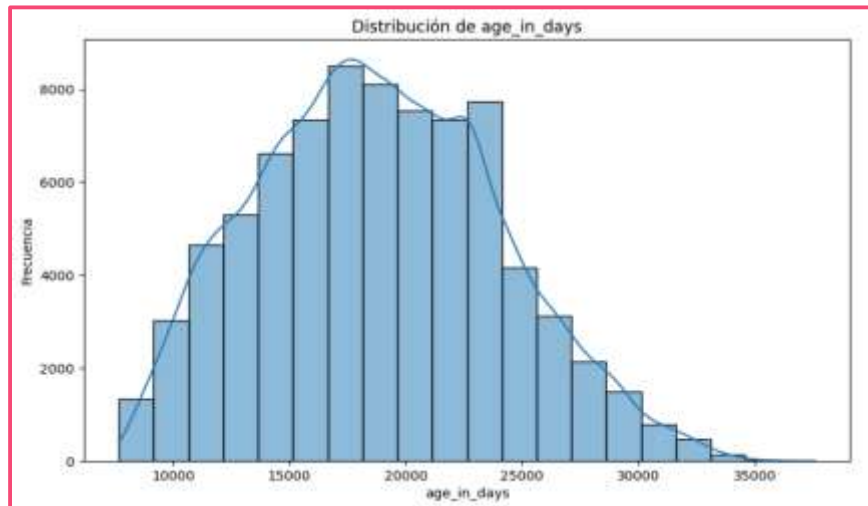
Rango 0–1, media ≈ 0.31 , mediana 0.17, con un 75% de los clientes por debajo de ≈ 0.54 . Distribución fuertemente sesgada a la derecha: muchos valores cercanos a 0 y una cola larga hacia valores altos.

Conclusión

La mayoría de clientes paga una fracción moderada o baja de la prima mediante efectivo/crédito, mientras que un grupo más pequeño paga casi el 100%. Este patrón sugiere comportamientos de pago heterogéneos que pueden estar asociados con diferentes niveles de compromiso o riesgo de no renovación.



age_in_days

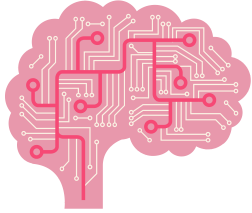


Hallazgos

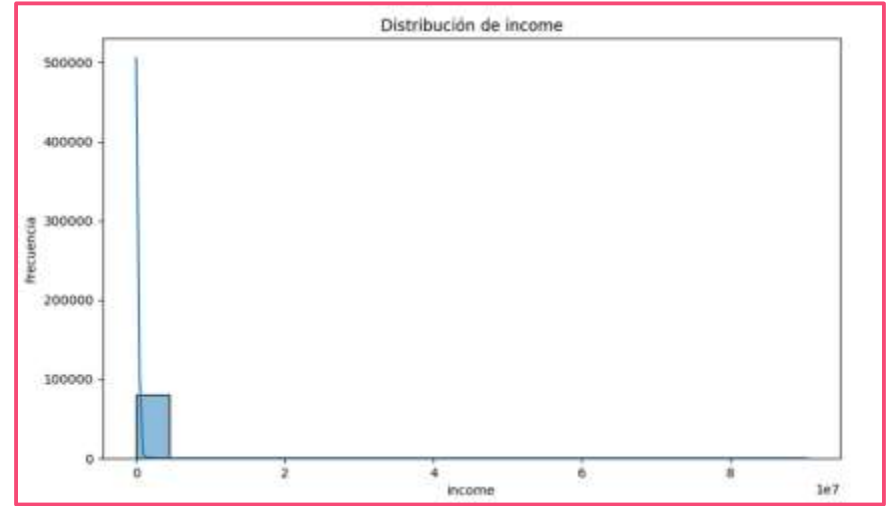
Rango entre $\approx 7,670$ y $\approx 37,602$ días (aprox. 21 a 103 años), media $\approx 18,847$ días ($\approx 51-52$ años) y mediana $\approx 18,625$ días. Distribución relativamente amplia pero sin valores extremos sospechosos.

Conclusión

La cartera cubre principalmente adultos de mediana edad y mayores, consistente con pólizas de vida/seguros a largo plazo. La edad es una variable relevante para entender patrones de renovación, pero no presenta problemas de calidad de datos.



income

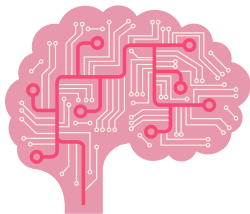


Hallazgos

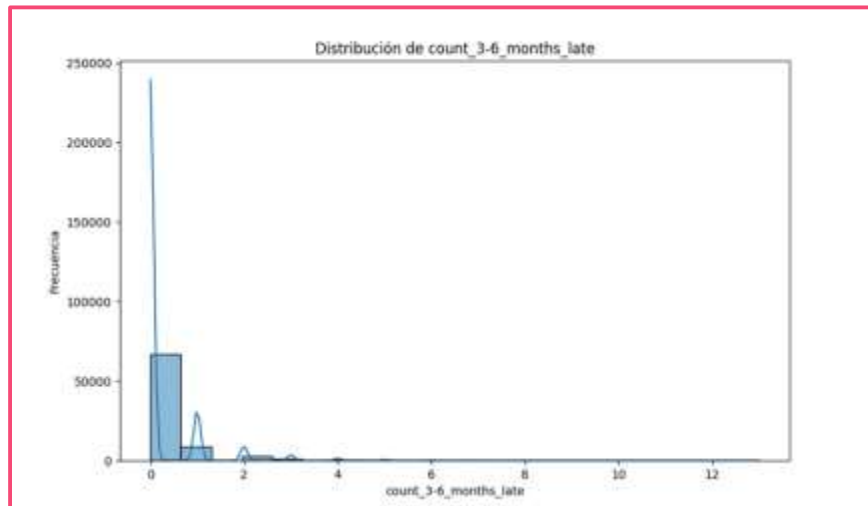
Rango muy amplio: mínimo $\approx 24,030$ y máximo superior a 90 millones. Desviación estándar muy alta, indicando la presencia de outliers fuertes en ingresos. Las medidas de posición (25%, 50%, 75%) se concentran bastante por debajo del máximo.

Conclusión

La distribución de ingresos está fuertemente sesgada por unos pocos clientes con ingresos extremadamente altos. Se requiere un tratamiento robusto para evitar que estos outliers distorsionen el modelado.



count_3-6_months_late

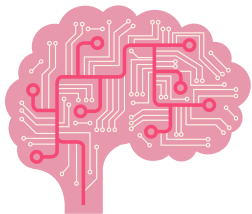


Hallazgos

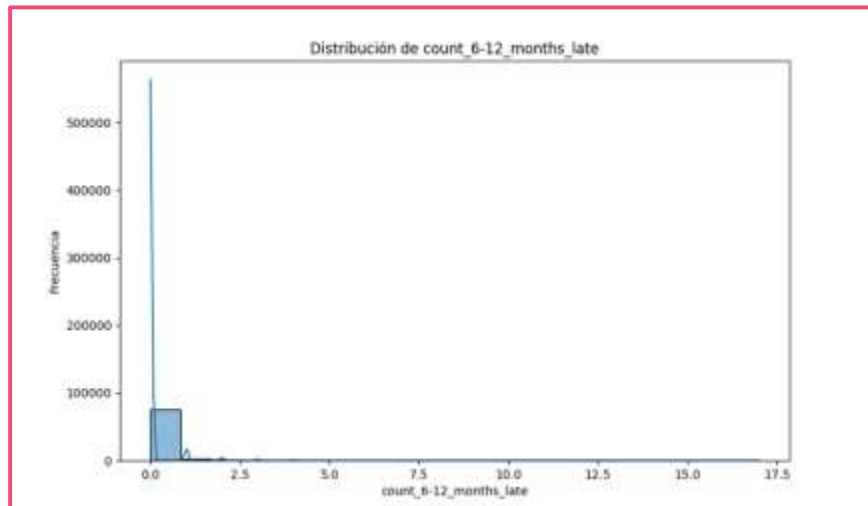
Media ≈ 0.25 , mediana 0, máximo 13. La mayoría de los clientes (percentiles 25%, 50% y 75%) tienen 0 atrasos en este rango; solo una minoría presenta varios atrasos.

Conclusión

Los atrasos de 3 a 6 meses son poco frecuentes, pero cuando aparecen lo hacen con valores altos, representando un segmento de clientes de mayor riesgo.



count_6-12_months_late

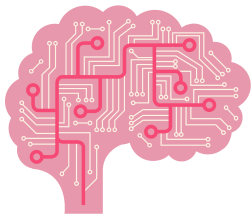


Hallazgos

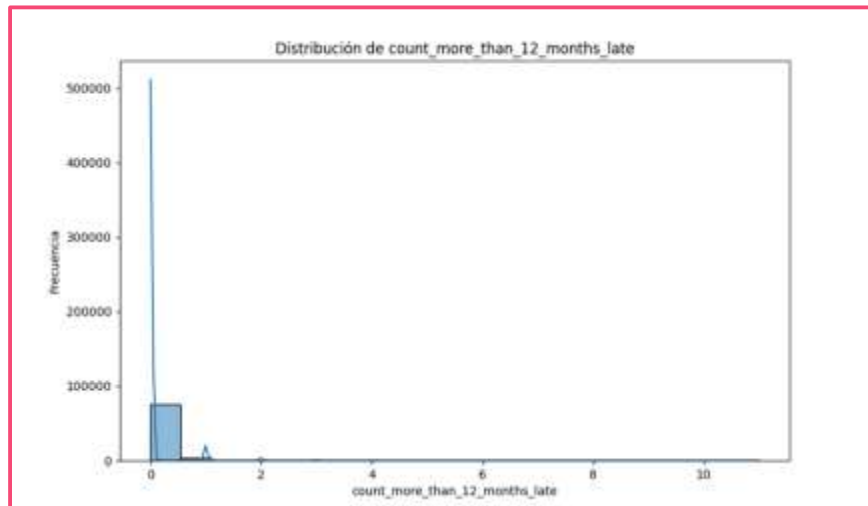
Media ≈ 0.08 , mediana 0, máximo 17. Casi todos los clientes no presentan atrasos en este rango; los que sí, conforman una cola larga de la distribución.

Conclusión

Los atrasos de 6 a 12 meses son aún más raros, pero su presencia indica comportamientos de impago prolongados que pueden estar fuertemente vinculados con la no renovación.



count_more_than_12_months_late

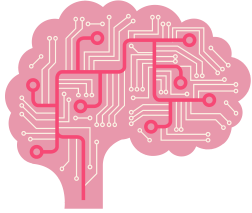


Hallazgos

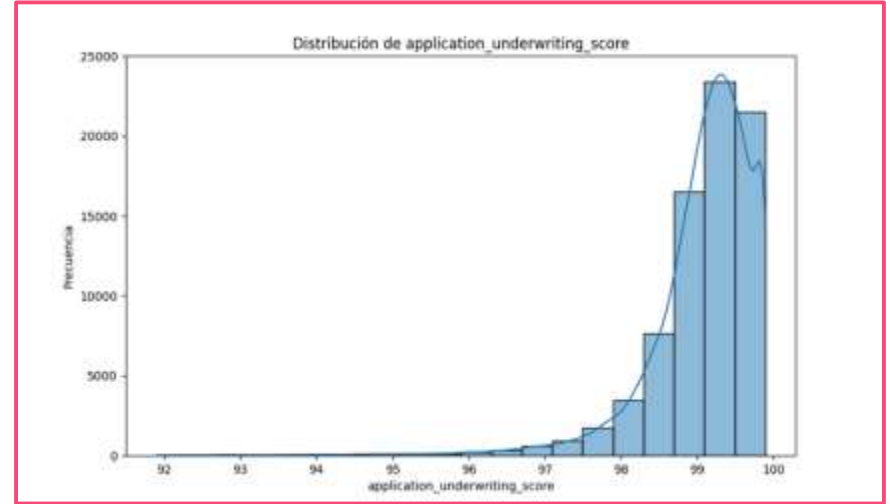
Media ≈ 0.06 , mediana 0, máximo 11. La distribución es muy concentrada en 0 atrasos con algunos casos extremos.

Conclusión

Atrasos mayores a 12 meses son eventos poco frecuentes pero muy críticos. Estos clientes representan un perfil de alto riesgo que puede ser clave para identificar probables no renovaciones.



application_underwriting_score

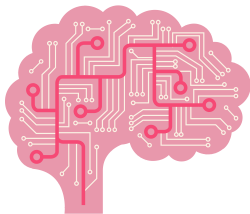


Hallazgos

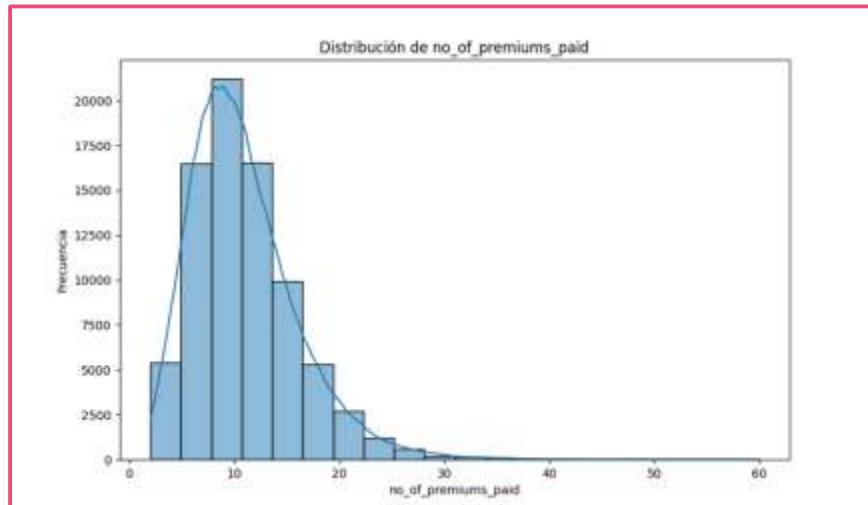
Media ≈ 99.07 , rango entre ≈ 91.9 y 99.89 , con percentiles muy concentrados alrededor de 99. La variabilidad es baja (desviación estándar ≈ 0.74).

Conclusión

La mayoría de clientes tiene un puntaje de suscripción alto y concentrado, lo que sugiere que la compañía trabaja principalmente con un perfil de riesgo relativamente bueno. A pesar de su baja dispersión, pequeñas diferencias en esta puntuación resultan estadísticamente significativas frente a la renovación, por lo que sigue siendo relevante para el modelo.



no_of_premiums_paid

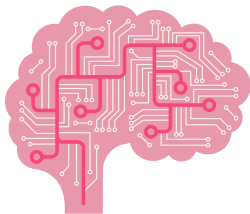


Hallazgos

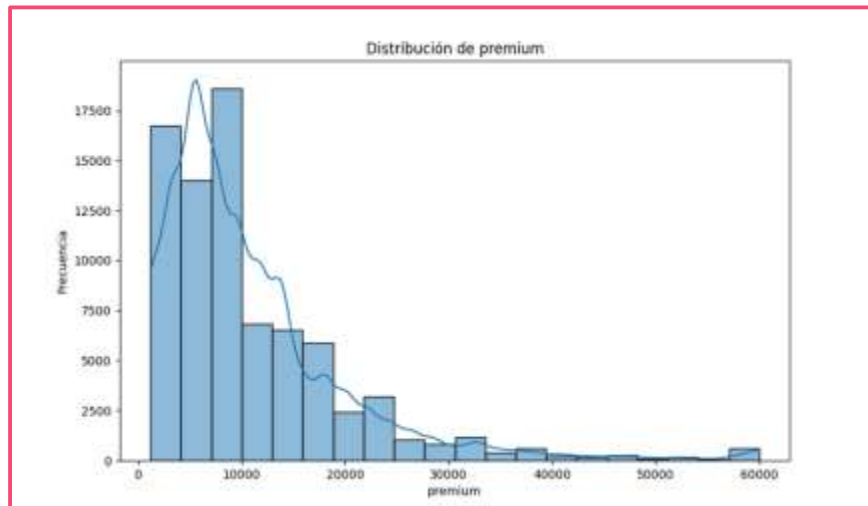
Media ≈ 10.86 , mediana 10, rango de 2 a 60 primas pagadas. El 50% de los clientes ha pagado entre 7 y 14 primas.

Conclusión

La mayoría de los asegurados tiene un historial moderado de pagos regulares, lo que refleja cierto nivel de antigüedad y compromiso con la póliza. Este historial puede influir en la decisión de continuar o no con la renovación.



premium



Hallazgos

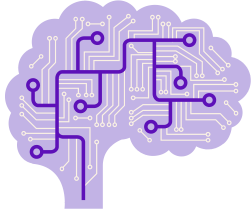
Rango entre 1,200 y 60,000, con media $\approx 10,925$ y mediana 7,500. Distribución sesgada a la derecha: muchas pólizas de primas bajas/medias y un grupo reducido con primas muy altas.

Conclusión

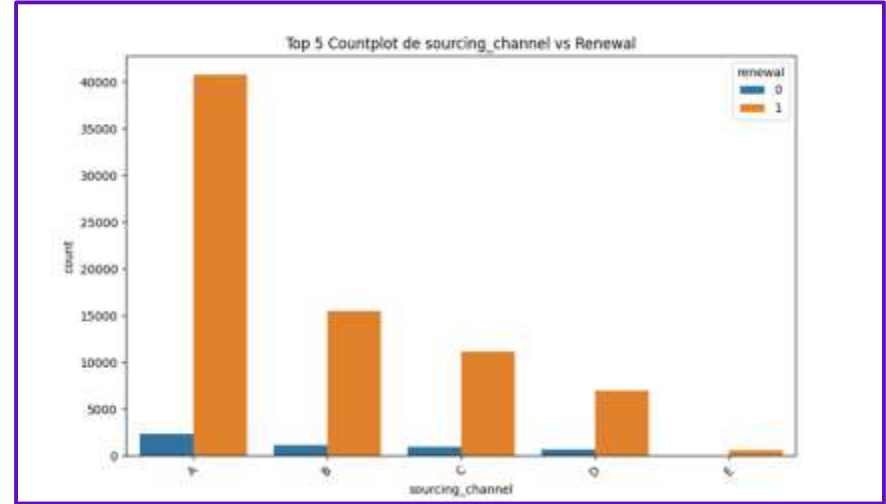
La cartera está compuesta principalmente por pólizas de montos moderados, mientras que las pólizas de alta prima constituyen un segmento pequeño pero relevante en ingresos. Esta variable es crítica para analizar el impacto económico de la renovación y la estrategia de incentivos.

Análisis univariable: Variables categoricas

3.2



sourcing_channel

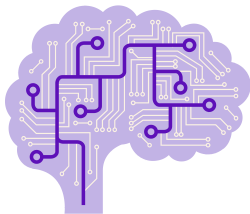


Hallazgos

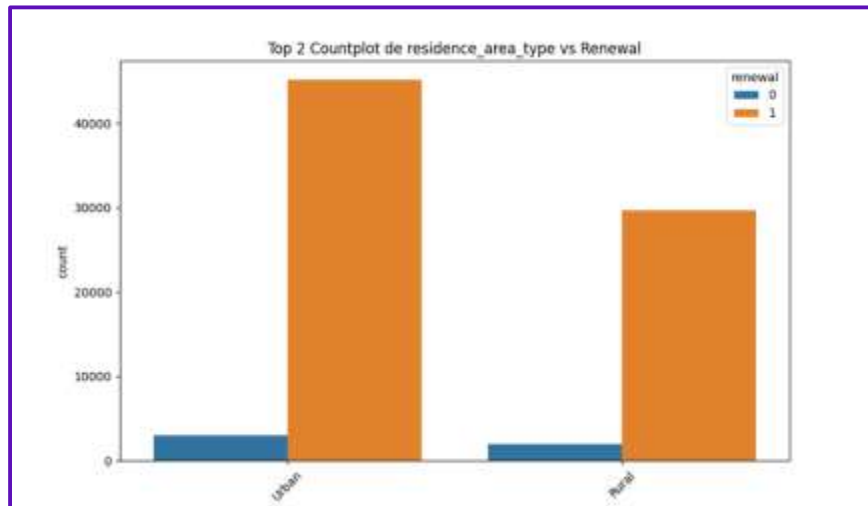
5 canales (A, B, C, D, E). El canal A concentra la mayoría de clientes ($\approx 43,000$), seguido por B y C; D y E son minoritarios, especialmente E.

Conclusión

El canal A es el principal motor de adquisición de clientes, por lo que cualquier estrategia de incentivos o gestión comercial debe prestar atención especial a este canal. Los canales pequeños, como E, podrían tener nichos específicos o menor eficiencia comercial.



residence_area_type

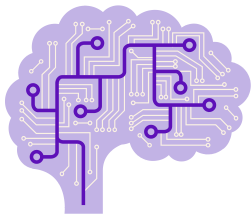


Hallazgos

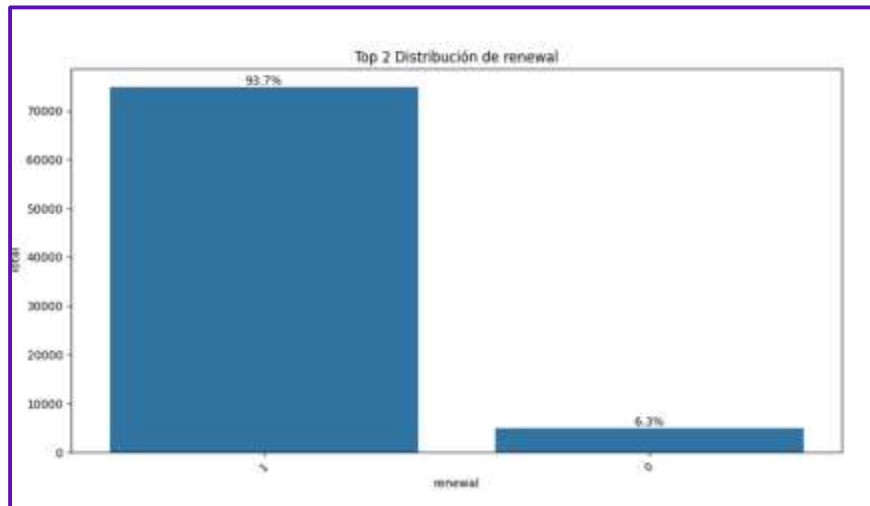
Dos categorías: *Urban* ($\approx 48,000$) y *Rural* (resto). La mayoría de asegurados reside en zonas urbanas.

Conclusión

La cartera está sesgada hacia clientes urbanos, lo cual puede reflejar la ubicación de la base de clientes de la compañía y el acceso a canales comerciales. Sin embargo, este desbalance no parece tener un impacto directo en la renovación.



renewal (variable objetivo)



Hallazgos

Variable binaria con fuerte desbalance: $\approx 93.7\%$ de clientes **renuevan** (1) y solo $\approx 6.3\%$ **no renuevan** (0).

Conclusión

El problema es una clasificación altamente desbalanceada. Deberán considerarse técnicas de balanceo para no sesgarse hacia la clase mayoritaria y realmente aprender a identificar casos de no renovación.

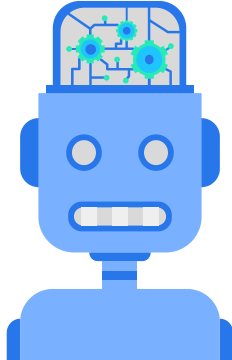
Análisis univariados: Conclusiones generales

Distribuciones sesgadas y outliers: Varias variables numéricas (income, premium, conteos de atrasos) están sesgadas a la derecha y presentan valores extremos. Esto exige escalado/normalización y posible tratamiento de outliers antes del modelado.

Histórico de pagos y comportamiento de atraso: La mayoría de clientes paga puntualmente, pero existe un subgrupo con varios atrasos en distintos rangos de tiempo. Este grupo es clave para entender el riesgo de no renovación.

Desbalance en la variable objetivo: La tasa de no renovación es baja (~6.3%). Este desbalance dificulta la detección de clientes en riesgo si no se corrige durante la fase de modelado.

Canales y residencia: El canal A y las zonas urbanas concentran la mayoría de la cartera, lo que marcará el foco de cualquier estrategia de segmentación e incentivos.



Análisis Bivariable

Numéricas vs Objetivo
(Renewal)

3.3

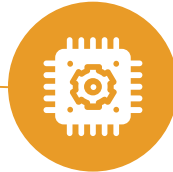
Numéricas vs Objetivo (Renewal)

Todas las variables numéricas analizadas muestran diferencias estadísticamente significativas entre clientes que **renuevan** y **no renuevan**, según la prueba de Kruskal-Wallis (p-value = 0.000 en todos los casos).

Comparando estadísticas por grupo (`renewal = 0` vs `renewal = 1`):



Perc_premium_paid_by_cash_credit: Los clientes que **no renuevan** tienen, en promedio, un mayor porcentaje pagado por efectivo/crédito (media ≈ 0.63 , mediana ≈ 0.73) comparado con quienes **sí renuevan** (media ≈ 0.29 , mediana ≈ 0.15).



age_in_days: Los clientes que **renuevan** tienden a ser ligeramente mayores ($\approx 18,975$) en promedio que los que **no renuevan** ($\approx 16,930$).



premium: Los clientes que **renuevan** tienen, en promedio, una prima algo más alta ($\approx 11,013$) que los que **no renuevan** ($\approx 9,600$).

Para las variables de atrasos y número de primas pagadas se observan diferencias significativas en la distribución (según Kruskal), aunque la mayoría de los valores sigue concentrada en rangos bajos.

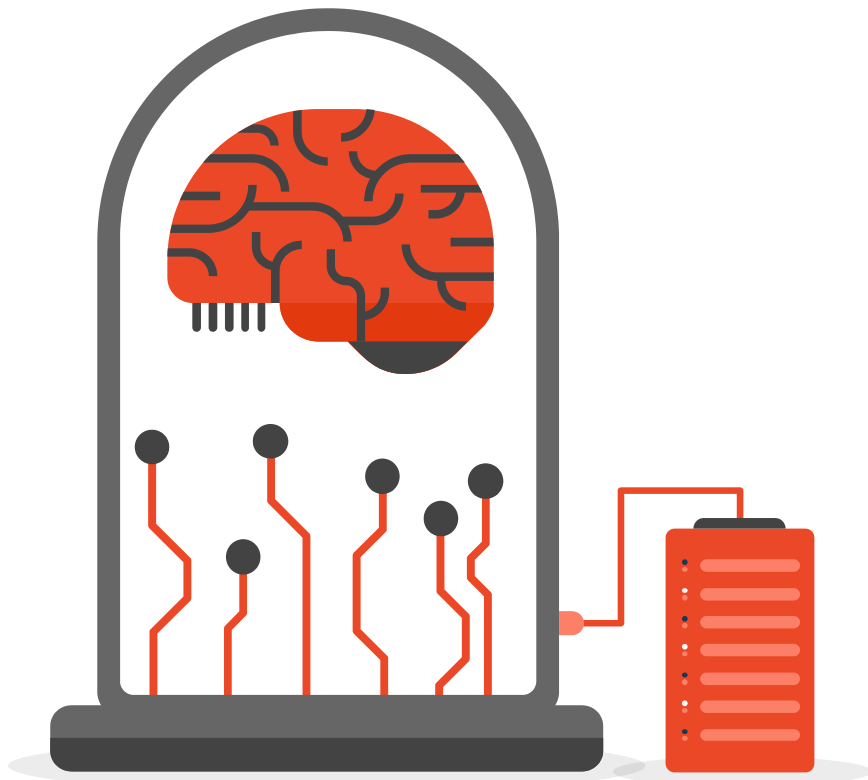
Conclusiones

Predictividad estadística: El hecho de que todas las variables numéricas resulten significativas ($p < 0.05$) confirma que aportan información relevante para distinguir entre clientes que renuevan y no renuevan.

Patrones de comportamiento de pago: Los clientes que no renuevan muestran, en promedio, un patrón de pagos distinto (mayor fracción de prima en efectivo/crédito y patrones de atraso diferentes). Esto sugiere que ciertas prácticas de pago pueden servir como señales tempranas de riesgo de fuga.

Perfil económico y de edad: Quienes renuevan tienden a tener primas más altas y ser ligeramente mayores, lo que podría indicar que clientes con más inversión en la póliza y más antigüedad/edad tienen mayor incentivo a continuar.

Implicaciones para el modelado: Todas estas variables deben conservarse para el modelado, aplicando transformaciones adecuadas (escalado/normalización y manejo de outliers). Su relevancia estadística justifica su inclusión como candidatos importantes para el modelo de propensión a renovar.



Análisis de Variables

Categóricas vs Objetivo
(Renewal)

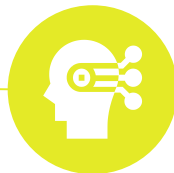
3.4

Categóricas vs Objetivo (Renewal)

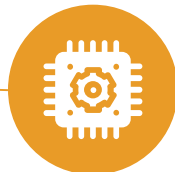
Todas las variables categóricas analizadas fueron comparadas frente a la variable objetivo renewal, usando pruebas de chi-cuadrado.

Los resultados muestran que:

- El canal de origen (sourcing_channel) sí tiene asociación significativa con la renovación (p-value = 0.000).
- El tipo de residencia (residence_area_type) NO muestra asociación con la renovación (p-value = 0.648).
- La distribución del objetivo (renewal) presenta un fuerte desbalance (93.7% renueva vs 6.3% no).



sourcing_channel: El canal **A** es el más frecuente tanto en clientes que renuevan como en los que no renuevan. Existe asociación estadísticamente significativa → algunos canales aportan clientes con **mayor propensión a renovar**.



residence_area_type: La mayoría de clientes reside en zona Urban, pero no se observan diferencias significativas en la renovación. *Es una variable descriptiva, pero no predictiva.*

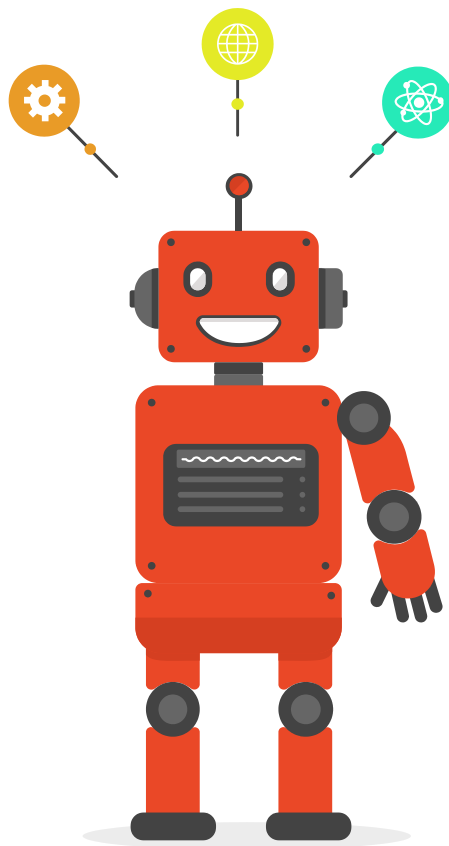


renewal: La variable presenta un fuerte desbalance: 93.7% renueva – 6.3% no renueva. Este desbalance debe ser tratado en la etapa de modelado para no sesgar los resultados.

La variable objetivo está fuertemente desbalanceada, por lo que será necesario aplicar técnicas de balanceo antes del modelado.

Conclusiones

- **Predictores clave de no renovación:** Las pruebas estadísticas (Kruskal-Wallis y chi-cuadrado) indican que todas las variables numéricas, así como el canal de adquisición, muestran diferencias significativas entre clientes que renuevan y no renuevan. En particular, destacan:
 - El patrón de pago (*perc_premium_paid_by_cash_credit*).
 - Los historiales de atraso (*count_3-6_months_late*, *count_6-12_months_late*, *count_more_than_12_months_late*).
 - El monto de la prima (*premium*).
 - El canal comercial (*sourcing_channel*).
- **Comportamiento de riesgo y compromiso:** Los clientes que no renuevan tienden a:
 - Tener patrones de pago diferenciados.
 - Mostrar comportamientos de atraso más marcados.
 - Presentar características económicas y demográficas ligeramente distintas. Esto respalda la hipótesis de que el comportamiento financiero e histórico de pago es un indicador fuerte del riesgo de no renovación.
- **Desbalance de clases y calidad de datos:**
 - La variable *renewal* está muy desbalanceada (≈93.7% renovaciones), lo que requiere aplicar técnicas como SMOTE u otras estrategias de balanceo.
 - Los niveles de nulos son bajos y concentrados en pocas variables (ya identificadas), por lo que la imputación es viable sin pérdida significativa de información.
- **Implicaciones para el modelo y el plan de incentivos:**
 - Las variables con mayor asociación con *renewal* deben priorizarse en el modelado y en la segmentación de clientes para el diseño del plan de incentivos.
 - El canal de adquisición se vuelve una palanca clave; puede utilizarse para ajustar incentivos según la calidad de los clientes que cada canal trae (propensión a renovar).
 - Dado el desbalance de la variable objetivo, la evaluación del modelo deberá ir más allá de la *accuracy*, priorizando métricas como *Recall* y *F1-score* para la clase de no renovación, ya que es el segmento de mayor interés para el negocio



Procesamiento de datos

4

Transformaciones Iniciales

01

Conversión de edad

Se transformó la variable ***age_in_days*** → ***age_in_years***
Para hacerla más interpretable
Reducir variabilidad
Evitar distorsiones al detectar outliers.



02

Tratamiento de datos nulos

Se identificó que solo 4 columnas tenían nulos:

- ***count_3-6_months_late***
- ***count_6-12_months_late***
- ***count_more_than_12_months_late***
- ***application_underwriting_score***

Se imputaron con la mediana porque:

- son distribuciones sesgadas
- la mediana maneja mejor los outliers
- no había suficiente correlación para usar KNNImputer

Resultado: **no quedaron nulos en el dataset.**

Análisis y tratamiento de outliers

Se detectaron outliers en casi todas las variables numéricas:

- Muy altos en **income** y **premium**
- Variaciones fuertes en los conteos de atrasos
- Algunos valores extremos en **application_underwriting_score**, **no_of_premiums_paid**, **age_in_years**

Qué se hizo:

Se aplicó transformación Log a:

- **income**
- **premium**

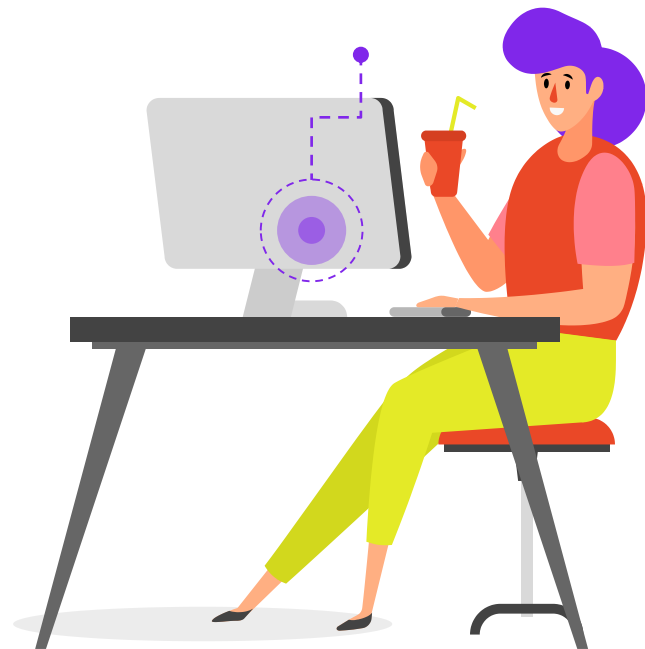
Para reducir sesgo extremo.

Se aplicó capping (winsorización) a:

- **application_underwriting_score**
- **no_of_premiums_paid**
- **Age_in_years**

Porque casi todos son 0 y el IQR no servía para recortarlos.

Se mantuvieron para conservar información de morosidad.



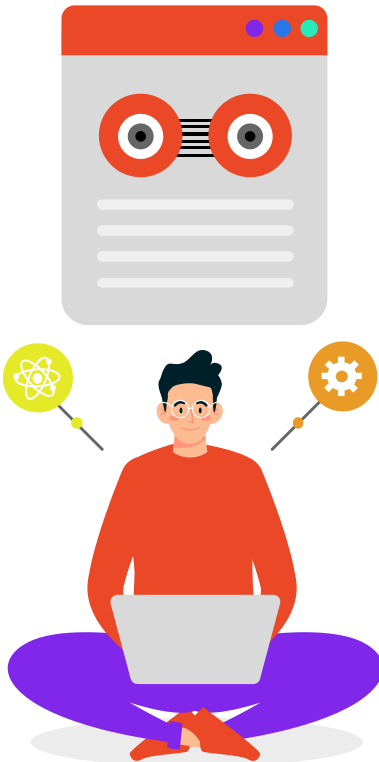
04

Tratamiento de variables categóricas

Se verificó que las categorías de:

- ***sourcing_channel***
- ***residence_area_type***

eran válidas (sin categorías raras)
No se eliminaron categorías ni se agruparon.



05

Codificación (One-Hot Encoding)

Se transformaron las categorías en variables binarias:

- ***sourcing_channel***
A / B / C / D / E
- ***residence_area_type***
Urban / Rural

No se eliminaron dummies
Se mantuvieron todas porque
algunos modelos (como KNN) no
se afectan por multicolinealidad.

Se crearon variables nuevas combinando los atrasos:

Nuevas columnas:

- **total_late_payments** → suma de todos los atrasos
- **has_late_payments** → si tiene al menos 1 atraso (0= no/1 si)

Para mejorar la capacidad predictiva del modelo.

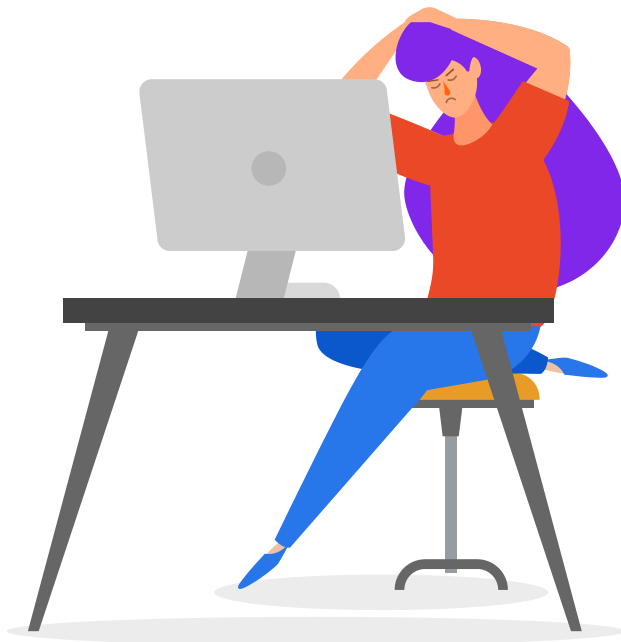
Luego se eliminaron:

- count_3-6_months_late
- count_6-12_months_late
- count_more_than_12_months_late

Para reducir dimensionalidad

Evitar multicolinealidad

Conservar sólo la información realmente útil



07

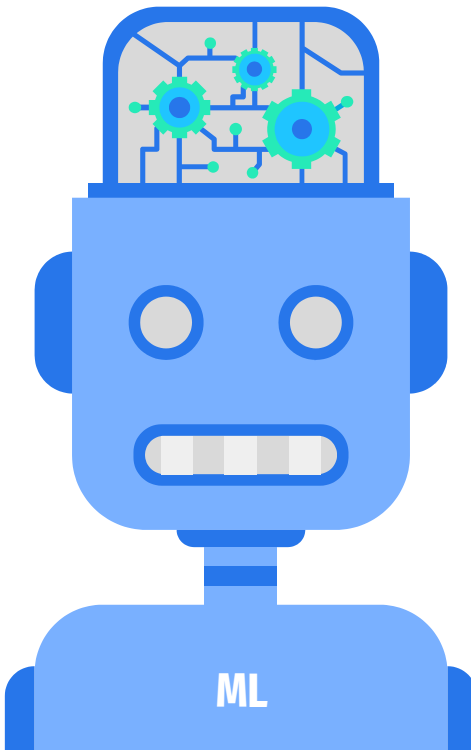
Escalado / Normalización

Se aplicó MinMaxScaler a todas las variables numéricas:

Porque KNN y modelos lineales necesitan variables en misma escala

Porque ya se manejaron outliers antes (log y capping)

Para mantener todo en rango [0,1]



08

Eliminación de columnas redundantes

Se eliminaron dummies y variables que no serían usadas en modelado:

- Todas las columnas dummy de canales y residencia
- premium
- no_of_premiums_paid

Porque eran redundantes o tenían correlación demasiado baja.

Selección de variables finales

Se seleccionaron sólo variables con **|correlación| > 0.05** contra renewal.

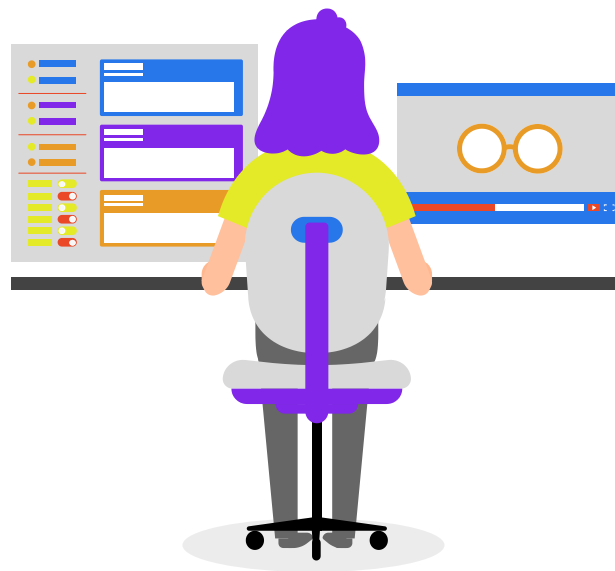
VARIABLES FINALES:

- *perc_premium_paid_by_cash_credit* (-0.24)
- *income* (0.06)
- *application_underwriting_score* (0.07)
- *age_in_years* (0.095)
- *total_late_payments* (-0.35)
- *has_late_payments* (-0.27)

Estas explican mejor la renovación

Evitan redundancia

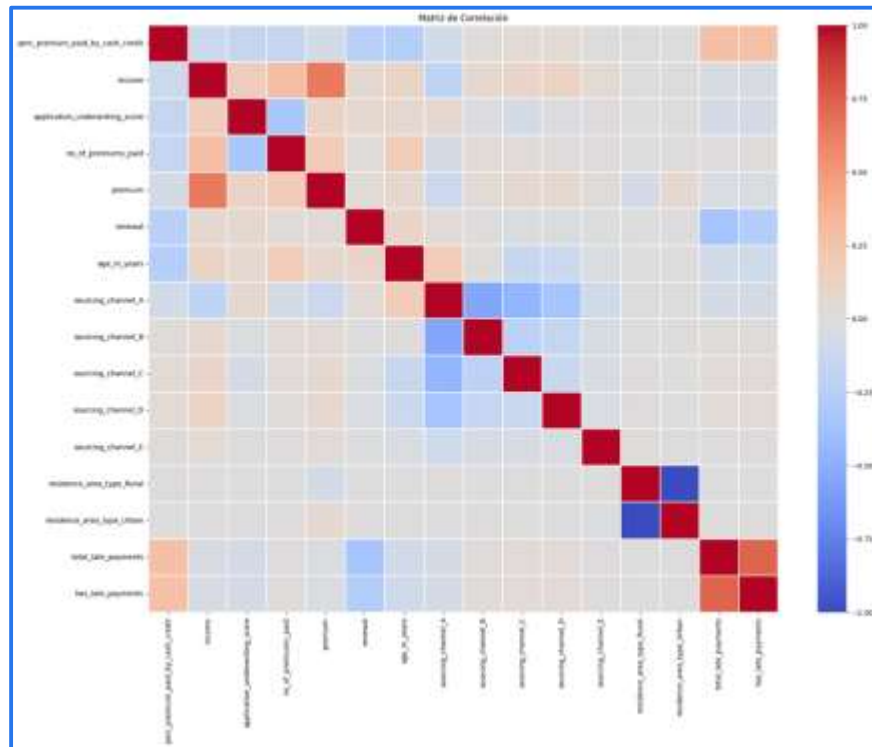
Combinan factores financieros, demográficos y de riesgo



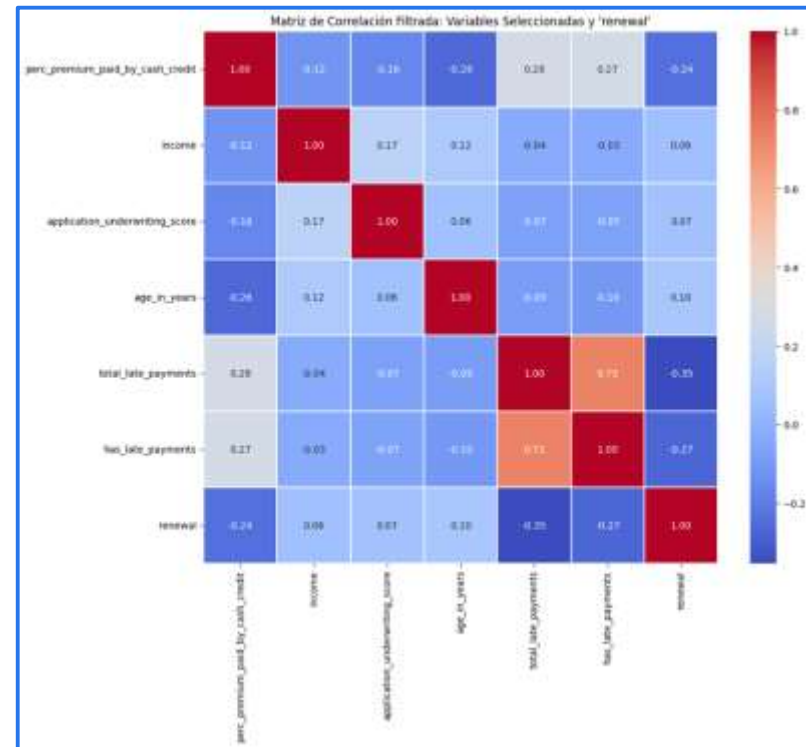
Procesamiento de datos:

Análisis de correlación

4.1



Matriz de correlación completa → muestra redundancia y baja correlación de canales/residencia

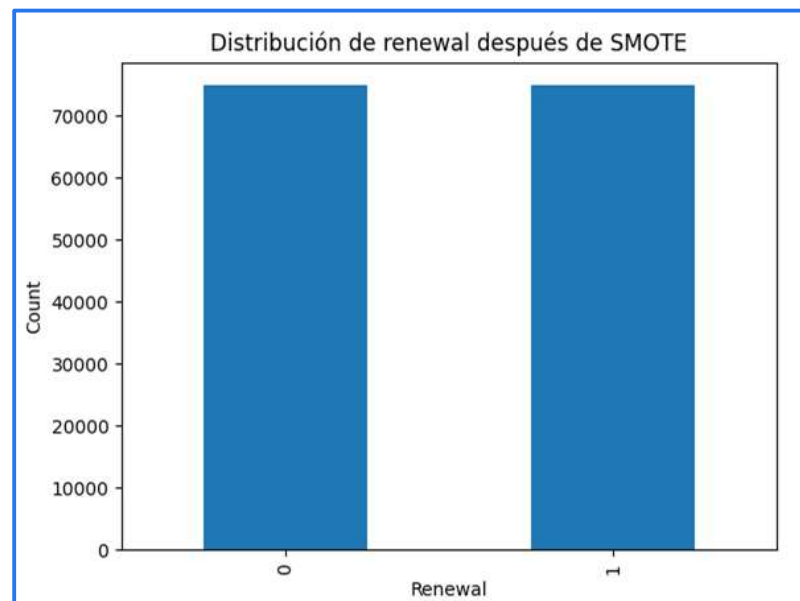
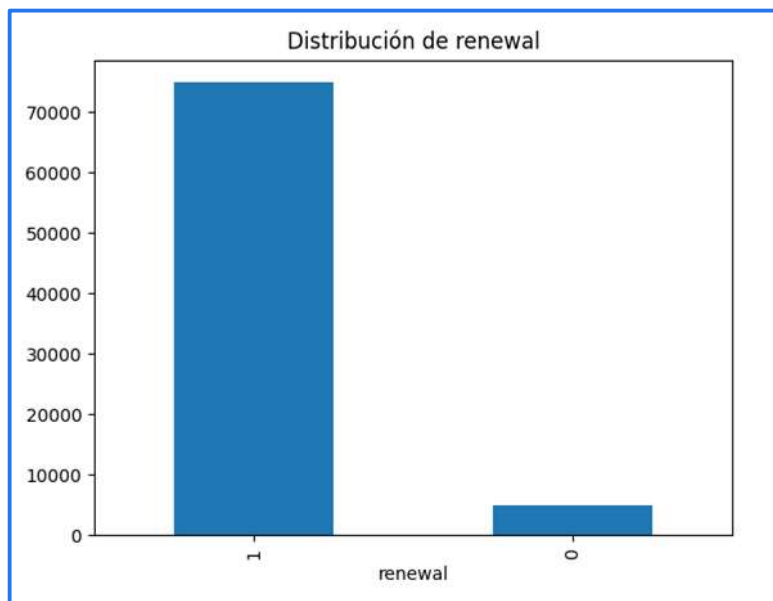


Matriz de correlación filtrada → muestra las variables realmente predictivas

Procesamiento de datos:

Balanceo de datos

4.2



Histograma renewal → evidencia desbalance 93.7% vs 6.3%

SMOTE → evidencia balance perfecto para modelado

Modelado / Comparación de algoritmos

5

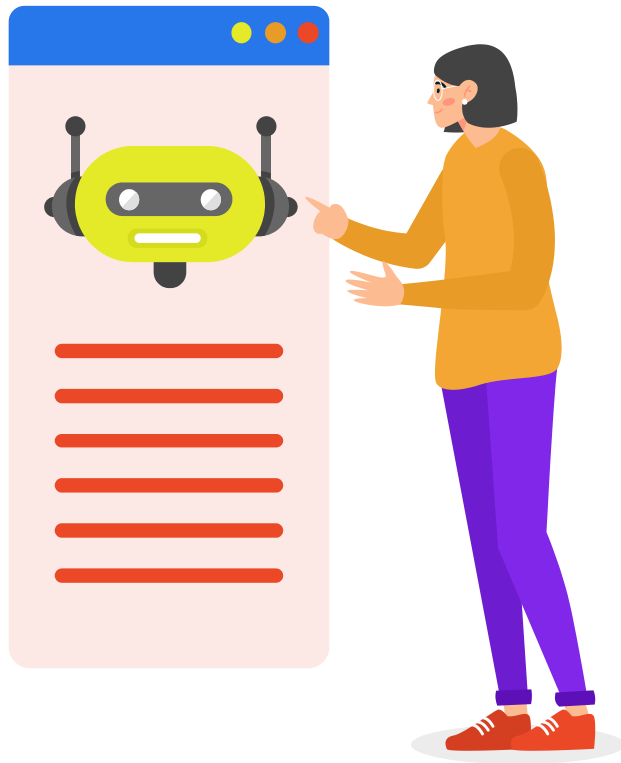
Modelado Predictivo

¿Qué se hizo en el modelado?

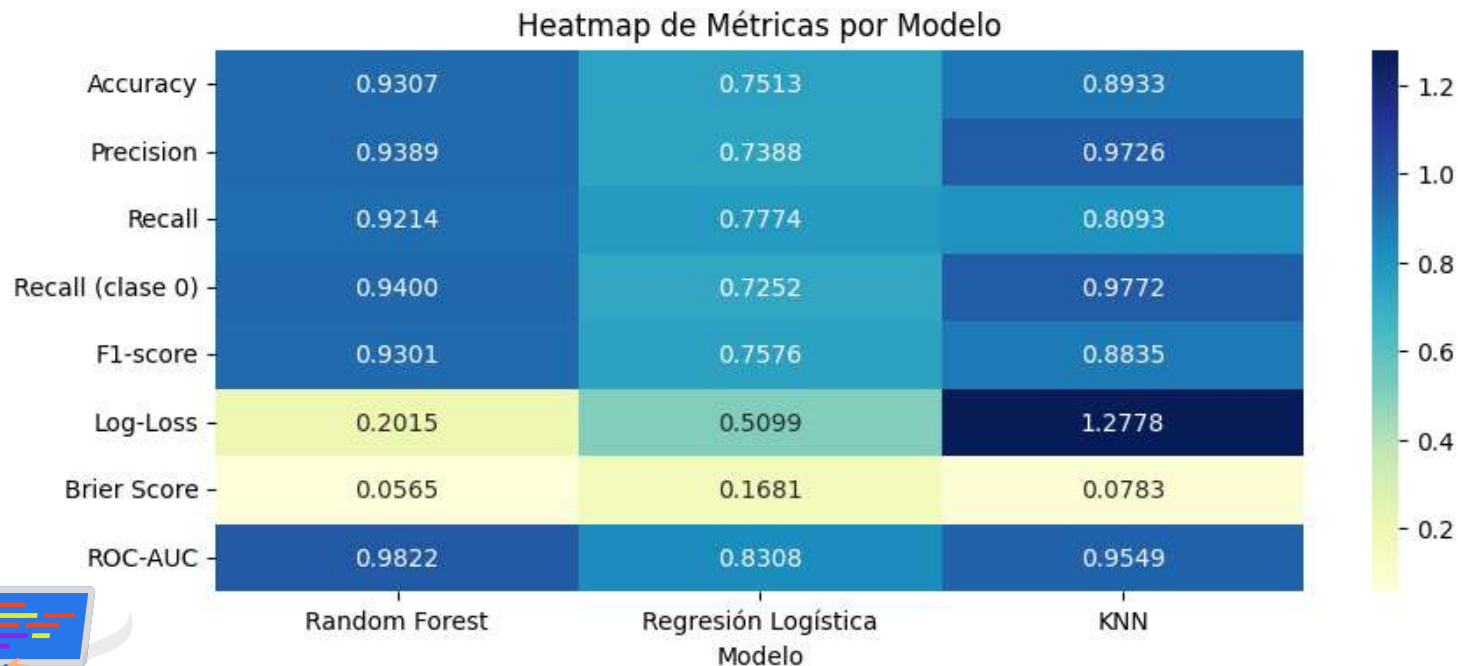
- Se entrenaron y compararon **3 algoritmos de clasificación**:
 - **Random Forest**
 - **KNN**
 - **Regresión Logística**
- Se usó el dataset **preprocesado, normalizado y balanceado con SMOTE**.
- Se realizó un **Train/Test Split 80/20**.
- Se evaluaron con métricas clave:
 - **Recall (clase 0 – clientes que NO renuevan)** → métrica más importante.
 - **ROC-AUC** → capacidad de discriminación.
 - **F1-score** → equilibrio general.
 - **Log-Loss / Brier Score** → calidad de probabilidades.

¿Por qué estas métricas?

Porque el negocio necesita **identificar clientes en riesgo de NO renovar** para ofrecerles incentivos y evitar pérdidas.



Comparación de métricas de los 3 modelos



Esta tabla comparativa de las métricas es la representación visual para explicar por qué un modelo ganó.

Resultados de los Modelos



Random Forest — Mejor desempeño general

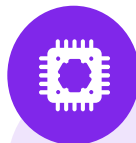
Accuracy: 0.9307

Precision: 0.9389

Recall: 0.9214 (excelente para clase 0)

ROC-AUC: **0.9822**

Modelo más robusto, estable y con mayor capacidad de identificar no renovaciones.



KNN — Buen desempeño, pero menos estable

Accuracy: 0.8933

Precision muy alta, pero

Recall (0.8093)

Buen modelo, pero menos confiable para retención.

&



Regresión Logística — Desempeño menor

Accuracy: 0.7513

Recall: 0.7774

ROC-AUC: **0.8308**

Menos adecuada para este caso.

&

Selección del modelo seleccionado

6

Modelo Seleccionado: Random Forest Classifier



¿Por qué ganó
Random
Forest?

Mejor Recall de clase 0 (0.94)

→ Identifica casi todos los clientes que NO renovarán.

Mayor ROC-AUC (0.9822)

→ Excelente capacidad para separar renovadores vs. no renovadores.

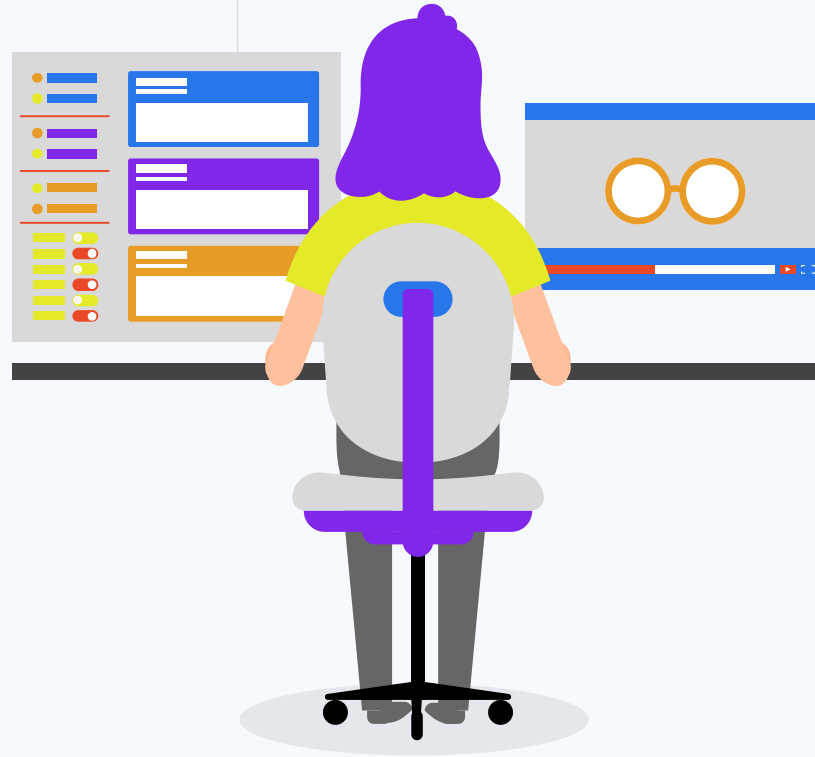
Probabilidades mejor calibradas (bajo Log-Loss y Brier Score)

→ Perfecto para priorizar incentivos según riesgo.

Es estable, robusto y maneja relaciones no lineales.

Impacto en el negocio

- Permite **detectar clientes en riesgo** con alta precisión.
- Optimiza los **planes de incentivos** reduciendo costos innecesarios.
- Incrementa la **tasa de retención** y protege los ingresos



FIN DE LA PRESENTACIÓN