



DESARROLLO DE UN MODELO DE RIESGO DE CRÉDITO

Helga Zambrana

Agenda

- 1. Presentación del problema**
- 2. Objetivo de la investigación**
- 3. Descripción del dataset elegido**
- 4. Hallazgos del análisis exploratorio de datos (EDA)**

1 Presentación del problema



El riesgo de crédito analiza la probabilidad de que un prestatario no reembolse un préstamo solicitado, generando pérdidas en el prestamista.

Las instituciones financieras tienen la necesidad de gestionar el riesgo de crédito para mitigar las pérdidas. A través del desarrollo de un modelo de aprendizaje automático se busca facilitar las aprobaciones de créditos, identificando los indicadores que diferencian a los solicitantes que tienden a pagar sus deudas de quienes no.

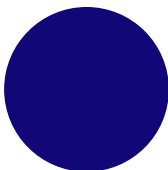
2 Objetivo de la investigación



Analizar la información disponible de los aplicantes a un préstamo e identificar las posibles variables que determinen que el solicitante no pagará el préstamo



Analizar diferentes opciones de modelos de clasificación y seleccionar el que tenga la mejor performance en predecir si el solicitante pagará o no el préstamo



Desarrollar un modelo de riesgo de crédito en Python para predecir las probabilidades de impago (PD) y asignar puntuaciones de crédito a los solicitantes existentes o potenciales

3 Descripción del dataset elegido



El dataset elegido fue [Home Equity \(HMEQ\)](#) de Kaggle.

HMEQ presenta información sobre las características y la morosidad de 5960 préstamos con garantía hipotecaria. Un préstamo con garantía hipotecaria es un préstamo en el que el deudor utiliza el capital de su vivienda como garantía subyacente.

Los criterios de selección del dataset fueron la claridad de los datos recolectados y su robustez, lo que facilita el análisis, el procesamiento y la generación de un modelo.

3 Descripción del dataset elegido

Nombre de la variable	Tipo de variable	Descripción
BAD	Binaria	1, corresponde a un prestamista con deuda impaga o con mora; 0, corresponde al candidato con los pagos al día. Se considera morosidad cuando han pasado 90 días desde que el prestamista no ha pagado
LOAN	Numérica continua	Monto en USD del préstamo solicitado
MORTDUE	Numérica continua	Monto en USD adeudado de la hipoteca existente
VALUE	Numérica continua	Monto en USD de la propiedad hipotecada
REASON	Categórica	Motivo para solicitar el préstamo DebtCon, corresponde a consolidación de deudas; HomeImp, corresponde a mejoras para el hogar
JOB	Categórica	Profesión o categoría profesional
YOJ	Numérica discreta	Cantidad de años en el trabajo actual

3 Descripción del dataset elegido

Nombre de la variable	Tipo de variable	Descripción
DEROG	Numérica discreta	Número de informes derogatorios importantes. Es información proporcionada por una institución financiera a las agencias de crédito y se relaciona con la morosidad o la cancelación de una cuenta de una línea de crédito
DELINQ	Numérica discreta	Número de líneas de crédito morosas
CLAGE	Numérica discreta	Edad expresada en meses de la línea de crédito de mayor antigüedad, los modelos de puntuación consideran un mínimo de 6 meses y recién a partir de 2 años es fiable
NINQ	Numérica discreta	Número de veces que ha solicitado un nuevo crédito en los últimos 2 años
CLNO	Numérica discreta	Número de líneas de crédito abiertas
DEBTINC	Numérica continua	Ratio deuda-ingreso (DTI), es la cantidad de ingresos brutos mensuales que una persona genera frente a la deuda que debe pagar por mes. Los prestamistas generalmente buscan ratios no mayores al 36%, aunque un DTI del 43% puede calificar para una hipoteca

Durante el análisis y pre-procesamiento de los datos se eliminaron y transformaron las variables que no aportarán al modelo y se encontraron los siguientes hallazgos

Variable Target

El dataset cuenta con una variable target definida 'BAD'. El valor 1 corresponde al candidato con préstamo incumplido o con mora y el valor 0 corresponde al candidato que paga su deuda.

Valores nulos y duplicados

No existen registros duplicados. El 84.6% de las columnas tiene datos faltantes. 'DEBTINC' es la variable con mayor cantidad de valores nulos con un 21.3% de su total.

Desbalanceo de datos

Se observa que de 5960 aplicantes, solo el 19.9% representa a personas con deuda morosa.

Correlación

Las variables relacionadas con el historial crediticio ('DELINQ', 'DEROG', 'NINQ') son las más correlacionadas con la variable target ('BAD'). Es un indicio de que éstas serán las variables de clasificación. Estas variables también están ligeramente correlacionadas entre sí sugiriendo que la información podría ser redundante.

El monto adeudado de la hipoteca existente ('MORTDUE') o el valor de la garantía subyacente ('VALUE') no parecen estar relacionados con el estado del préstamo. De todos modos, forman otro grupo de correlación con otras variables como los Años_Crédito y el número de líneas de crédito ('CLNO').

GRACIAS



Helga Zambrana | Data Science
[LinkedIn](#)