

ITESM- Campus Puebla

## **American Express - Default Prediction**

## Inteligencia artificial avanzada para la ciencia Datos

## **Integrantes Equipo 1:**

Myroslava Sánchez Andrade A01730712 José Antonio Bobadilla García A01734433 Karen Rugerio Armenta A01733228 Alejandro Castro Reus A01731065

Fecha: 26/08/2022

# I. Reto de limpieza de conjunto de datos

#### 1 Introducción

Las empresas bancarias tienen un peso muy grande en la economía global. Tan sólo en México se considera que al menos 41.1 millones de Mexicanos, que representa a un 49.1% del total de adultos de 18 a 70 años, tienen al menos una cuenta bancaria, según el Instituto Nacional de Estadística y Geografía (Inegi, 11 de mayo de 2022).

Los bancos desempeñan un papel importante, ya que se encargan de permitir el flujo de recursos financieros en todo el país, distribuir efectivo, pagar los cheques que se emiten, ofrecer servicio con tarjetas y transferencias bancarias, otorgar préstamos, entre otros.

Para este proyecto, nos enfocaremos en los préstamos bancarios que los bancos otorgan mediante tarjetas de crédito. El Portal de Educación Financiera, define el crédito como un préstamo de dinero que una parte otorga a otra, con el compromiso de que, en el futuro, quien lo recibe devolverá dicho préstamo en forma gradual (mediante el pago de cuotas) o en un solo pago y con un interés adicional que compensa a quien presta, por todo el tiempo que no tuvo ese dinero.

Para poder tener acceso a una tarjeta de crédito con American Express, se deben cumplir ciertos requisitos que garanticen que la persona que está aplicando al crédito pagará los montos acordados conforme a las tasas de interés. Es por ello, que este proyecto tratará de predecir si un usuario pagará el crédito de la tarjeta a la que aplicó, haciendo uso de un modelo de Machine Learning y entrenando un algoritmo con base en un historial de usuarios.

### 1.1 American express - Default Prediction

#### Desarrollo

Como primer acercamiento a la solución de esta primera fase del reto, se intentó cargar el dataset completo de datos en Python usando la librería de pandas, el cuál tiene un peso de 16gb y un contenido de aproximadamente 5 millones de filas. Sin embargo, debido a la magnitud del dataset las computadoras que se utilizaron no pudieron completar la operación debido a sus características. Por esta misma razón realizar el análisis y posteriormente la limpieza de datos, se complicó un poco más de lo que se esperaba, ya que esto involucra diversas operaciones que no se podían completar.

Solamente fue posible realizar la lectura y procesamiento del dataset completo en la computadora de un integrante del equipo, que cuenta con un procesador Arm M1 de Apple y

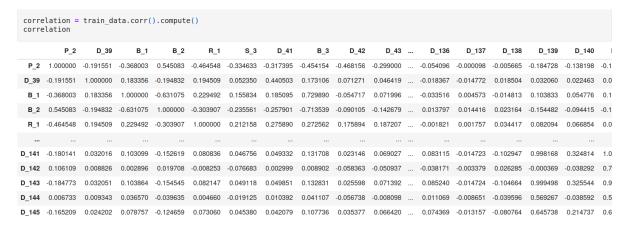
8gb de RAM. Esto fue un contratiempo ya que se tuvo que organizar de una mejor manera el equipo para avanzar con la limpieza de los datos como primer aproach y posteriormente investigar alternativas que ayudaran a analizar el dataset completo en cada una de las computadoras de los miembros del equipo, teniendo en cuenta el hardware de cada integrante.

Los datos de train y test fueron extraídos en formato CSV desde Kaggle. Dado el tamaño del dataset y la naturaleza de carga completa de este formato, decidimos guardar parte de la información de entrenamiento en formato *parquet*, particionando adicionalmente a través de compresión *snappy*. Posteriormente los datos fueron cargados en un script de python como dataframe de dask que nos permite dividir los datos para procesarlos de manera paralela.

```
# Get data
train_data = dd.read_parquet('train_data/snappy-parquet', engine='pyarrow')
train_labels = dd.read_csv('train_labels.csv')
# train_data = train_data.compute()

new_train_data.to_parquet("train_data/filtered-columns", engine="pyarrow", compression="snappy")
```

- Verificamos el contenido y estructura de las columnas notando que ya se encontraban estandarizadas y continuamos con el cálculo de la correlación.



Por cada conjunto de columnas cuya correlación es superior a 0.9, se conservó una sola columna de acuerdo a la justificación incorporada en el código.

```
for column in correlation.columns:
    local = correlation[column][correlation[column].abs() > 0.9]
    local = local[local.index != column]
         if local.count() > 0:
    print(local)
                print()
B_11 0.995574
B_37 0.992915
Name: B 1, dtype: float64
B 33
               0.912814
Name: B_2, dtype: float64
            0.904635
Name: S_3, dtype: float64
              0.995051
B 23
Name: B_7, dtype: float64
            0.995574
               0.987941
Name: B_11, dtype: float64
            0.904635
Name: S_7, dtype: float64
 # B_1 -> B_11, B_37, elegimos porque tenía mayor correlación y no tenía nulos # B_33 -> B_2, elegimos porque la distribución era ligeramente más compacta
  # B_33 -> B_2, elegimos porque la distribución era ligeramente mas compact

# S_7 -> 5_3, porque tiene menos valores extremos

# B_23 -> B_7, menos valores extremos

# D_58 -> D_74, D_75, pues tiene una distribución ligeramente más uniforme

# B_14 -> B_15, no tiene tantos valores faltantes

# D_62 -> D_77, no tiene tantos valores faltantes
  # S_24 -> S_22, no tiene tantos valores faltantes y tiene una distribución ligeramente más compacta
# D_103 -> D_104, la distribución es más compacta
# D_118 -> D_119, misma distribución
  # D_139 -> D_141, D_143, mejor correlación con ambas drop_columns = ['B_11', 'B_37', 'B_2', 'S_3', 'B_7', 'D_74', 'D_75', 'B_15', 'D_77', 'S_22', 'D_104', 'D_119', 'D_141', 'D_143']
```

- Algunas columnas tenían un porcentaje de nulos mayor al 30%. Dichas columnas fueron sustraídas.

```
max_val = train_data['customer_ID'].count().compute()
count_non_null = train_data.count().compute()
count_non_null
customer ID
               5531451
               5531451
S_2
P_2
               5485466
B 1
               5531451
               5429903
D 141
D 142
                944408
               5429903
D 144
               5490724
D 145
               5429903
Length: 190, dtype: int64
drop_columns += count_non_null[count_non_null < max_val * 0.7].index.to_list()</pre>
```

#### Conclusiones

Se nos indicó que, teniendo en cuenta las complicaciones con respecto a la potencia de nuestros equipos, podríamos usar una muestra más pequeña que respetara la distribución de clasificación. Seguiremos haciendo uso de Dask para una mejor computación de los datos. Posterior a que hayamos tomado una muestra con un tamaño adecuado para nuestros equipos, nos faltaría la última parte de esta etapa que es la estandarización de nuestros datos (estandarización one-hot encoding). Para un mejor análisis de los datos, nos gustaría, adicionalmente, aplicar PCA y un análisis de las variables a través del tiempo.